

# LASSO Clássico e Bayesiano

Kévin Allan Sales Rodrigues

*kevin.asr@outlook.com*

Universidade de São Paulo  
IME - Instituto de Matemática e Estatística  
Programa de Pós-Graduação em Estatística

31 de outubro de 2018



IME-USP

- 1 Introdução
- 2 Modelo Linear
- 3 LASSO Clássico
- 4 Escolhendo o Valor de  $\lambda$
- 5 Inferência no LASSO
- 6 LASSO no R
- 7 Abordagem Bayesiana: BLASSO, the Bayesian LASSO
  - Amostrador de Gibbs
  - Escolhendo o  $\lambda$
- 8 Extensões do LASSO



O que significa **LASSO**?

**Least Absolute Shrinkage and Selection Operator.**

Em tradução livre, significa operador de seleção e redução de variáveis/coeficientes via norma  $L_1$ .

O Lasso foi proposto por Tibshirani (1996, JRSS) e sua versão bayesiana foi descrita e aprofundada em Park e Casella (2008, JASA).



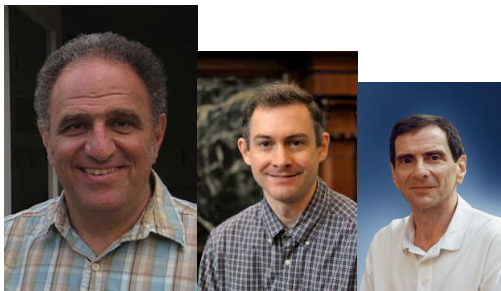


Figura: Tibshirani (LASSO), Park (BLASSO) e Casella (BLASSO).

## Porque utilizar o LASSO?

- Pode ser usado quando  $p > n$  (alta dimensionalidade);
- Pode estimar alguns  $\beta_j$  como 0;
- Evita a associação espúria;
- É um método parcimonioso;
- Predição: diminui a variância e aumenta viés.

Obs: as estimativas do LASSO não são invariantes por escala, por isso as covariáveis,  $\mathbf{X}$ , devem ser padronizadas antes de usar o LASSO.



- 1 Introdução
- 2 Modelo Linear**
- 3 LASSO Clássico
- 4 Escolhendo o Valor de  $\lambda$
- 5 Inferência no LASSO
- 6 LASSO no R
- 7 Abordagem Bayesiana: BLASSO, the Bayesian LASSO
  - Amostrador de Gibbs
  - Escolhendo o  $\lambda$
- 8 Extensões do LASSO



O modelo de regressão linear é dado pela Equação (1),

$$Y_i = \beta_0 + \beta_1 X_{1i} + \cdots + \beta_p X_{pi} + \epsilon_i, \quad i = 1, 2, \dots, n \quad (1)$$

em que  $Y_i$  é a  $i$ -ésima variável resposta,  $\beta_0$  é o intercepto do modelo,  $\beta_j$ ,  $j = 1, \dots, p$  é o parâmetro associado à  $j$ -ésima covariável,  $X_{ji}$ ,  $j = 1, \dots, p$ ,  $i = 1, \dots, n$  é a  $j$ -ésima covariável da  $i$ -ésima observação. e  $\epsilon_i$  é o erro associado à  $i$ -ésima observação.



Podemos expressar o modelo de forma mais compacta, como:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

em que,  $\mathbf{Y}_{n \times 1} = (Y_1, Y_2, \dots, Y_n)^\top$ ,  $\mathbf{X}_{n \times (p+1)} = (\mathbf{X}_1^\top, \mathbf{X}_2^\top, \dots, \mathbf{X}_n^\top)^\top$ ,  
 $\mathbf{X}_k = (1, X_{1k}, \dots, X_{pk})$ ,  $\boldsymbol{\beta}_{(p+1) \times 1} = (\beta_0, \beta_1, \dots, \beta_p)^\top$  e  
 $\boldsymbol{\epsilon}_{n \times 1} = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)^\top$ . Uma suposição comumente utilizada é que vale  
 $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{I}_n \sigma^2)$ .





O ajuste do modelo linear via mínimos quadrados é dado pela solução do problema

$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n (Y_i - \mathbf{X}_i \beta)^2,$$

se  $p + 1 \leq n$  então a matriz  $\mathbf{X}$  terá posto completo, consequentemente  $\mathbf{X}^T \mathbf{X}$  tem posto completo e sua inversa existe. Então obtemos a conhecida solução de mínimos quadrados

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y},$$

mas perceba que a solução só é única se  $p \leq n$ , caso contrário a solução de mínimos quadrados não será única. Isto é, quando temos dados em que  $p + 1 > n$  não temos como ajustar todos os modelos (no sentido de obter ajuste único para cada modelo) pois nem sempre existirá a inversa de  $\mathbf{X}^T \mathbf{X}$ .



- 1 Introdução
- 2 Modelo Linear
- 3 LASSO Clássico**
- 4 Escolhendo o Valor de  $\lambda$
- 5 Inferência no LASSO
- 6 LASSO no R
- 7 Abordagem Bayesiana: BLASSO, the Bayesian LASSO
  - Amostrador de Gibbs
  - Escolhendo o  $\lambda$
- 8 Extensões do LASSO



O estimador via LASSO no contexto da regressão linear é dado pela resolução do problema

$$\hat{\beta}_L = \arg \min_{\beta} \left( \sum_{i=1}^n (Y_i - \mathbf{x}_i \beta)^2 + \lambda \sum_{j=1}^p |\beta_j| \right), \lambda \geq 0 \text{ (fixo)}, \quad (2)$$

ou equivalentemente, por

$$\hat{\beta}_L = \arg \min_{\beta} \left( \sum_{i=1}^n (Y_i - \mathbf{x}_i \beta)^2 \right) \text{ restrito à } \sum_{j=1}^p |\beta_j| \leq s, s \geq 0. \quad (3)$$



O que ocorre no estimador se na Equação (2)  $\lambda = 0$  e se  $\lambda \rightarrow \infty$  ?

O que ocorre no estimador se na Equação (3)  $s = 0$  e se  $s \geq \sum_{j=1}^p \hat{\beta}_j$   
( $\hat{\beta}_j, j = 1, \dots, p$  são os estimadores de mínimos quadrados) ?

Então como escolher o valor de  $\lambda$  ?



O estimador via Ridge regression,  $\hat{\beta}_R$ , é dado por

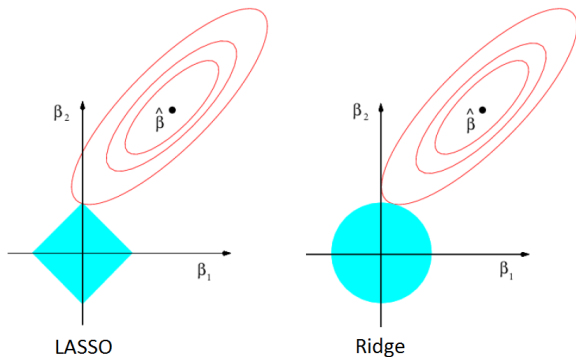
$$\hat{\beta}_R = \arg \min_{\beta} \left( \sum_{i=1}^n (Y_i - \mathbf{x}_i \beta)^2 + \eta \sum_{j=1}^p (\beta_j)^2 \right), \eta \geq 0 \text{ (fixo)}, \quad (4)$$

ou equivalentemente, por

$$\hat{\beta}_R = \arg \min_{\beta} \left( \sum_{i=1}^n (Y_i - \mathbf{x}_i \beta)^2 \right) \text{ restrito à } \sum_{j=1}^p (\beta_j)^2 \leq r, r \geq 0. \quad (5)$$



# LASSO e Ridge Regression



**Figura:** Esquema gráfico do LASSO e Ridge regression no caso em que existem apenas duas variáveis explicativas.

Fonte: Hastie, Tibshirani e Friedman (2008) (com adaptações).



IME-USP

- 1 Introdução
- 2 Modelo Linear
- 3 LASSO Clássico
- 4 Escolhendo o Valor de  $\lambda$**
- 5 Inferência no LASSO
- 6 LASSO no R
- 7 Abordagem Bayesiana: BLASSO, the Bayesian LASSO
  - Amostrador de Gibbs
  - Escolhendo o  $\lambda$
- 8 Extensões do LASSO



## Escolhendo o Valor de $\lambda$

Escolher um bom valor de  $\lambda$  é fundamental para que o método funcione adequadamente, um modo de escolher o valor de  $\lambda$  é a validação cruzada que é um método bastante fácil de entender e de implementar.

Basicamente dividiremos, aleatoriamente, a amostra em  $k$  partes iguais ou pelo menos aproximadamente iguais. Escolheremos a primeira parte para ser os dados de “validação” e as demais para ser dados de “treinamento”, ajustaremos o LASSO com  $\lambda = \lambda_0$  aos dados de “treinamento” e usaremos esse modelo para tentar prever os dados de “validação”, então calcularemos o erro de predição. Repetiremos esse mesmo procedimento mais  $k - 1$  vezes para as outras partes restantes. Após terminar as  $k$  iterações teremos  $k$  erros de predição, calcularemos a média dos  $k$  erros de predição. Esse procedimento será feito para vários valores  $\lambda_0$  distintos. E escolheremos o valor  $\lambda_0$  que minimize o erros de predição médio.





## Escolhendo o Valor de $\lambda$

Tibshirani (1996) recomenda  $k = 5$  ou  $k = 10$ , mas naturalmente quaisquer valores inteiros entre cinco e dez podem ser usados, não há uma regra para escolha de  $k$ . Claro que o número de observações deve ser “suficientemente grande” para ser dividido entre as  $k$  partes.

É importante salientar que este método de escolha do valor de  $\lambda$  não depende do conhecimento do número de parâmetros (“graus de liberdade”) do modelo e também não depende de uma estimativa do parâmetro de escala do modelo. Embora na situação tradicional em que se usa o estimador de mínimos quadrados o número de parâmetros do modelo ajustado esteja bem definido e seja fácil achar uma estimativa para o parâmetro de escala, no caso do LASSO isso não é trivial, pois  $\lambda$  também é uma quantidade que tem impacto no ajuste do modelo (embora  $\lambda$  não apareça de forma explícita no modelo ajustado). Logo esse método é proposto por Tibshirani para contornar esses problemas.



O erro de predição relacionado à  $\lambda_0$  é dado por

$$EP_{\lambda_0} = \frac{1}{n} \sum_{i=1}^k n_i EQM_i,$$

em que  $n_i$  é o número de observações da  $i$ -ésima parte dos dados e

$$EQM_i = \frac{1}{n_i} \sum_{j=1}^{n_i} (Y_j - \hat{Y}_j)^2,$$

ou seja,  $EQM_i$  é o erro quadrático médio quando usamos a  $i$ -ésima parte dos dados como “validação” e o erro de predição é simplesmente uma média ponderada dos EQM's. Essa é a quantidade que nos guiará na escolha do valor de  $\lambda$ , quanto menor o erro de predição, mais adequado é o valor de  $\lambda$ .



Em geral não se utiliza o valor de  $\lambda$  que minimiza o erro de predição, em vez disso se utiliza o maior valor de  $\lambda$  cujo erro de predição esteja a um desvio padrão do erro de predição associado ao valor de  $\lambda$  que minimiza o erro de predição. Tibshirani (1996) argumenta que como nós não conhecemos a real forma da curva do erro de predição em termos do valor de  $\lambda$  (afinal temos apenas uma curva aleatória, pois a validação cruzada seleciona aleatoriamente quais observações pertencerão a dada parte dos dados) devemos escolher o valor de  $\lambda$  como definido anteriormente, favorecendo um valor de  $\lambda$  maior e conseqüentemente um modelo mais simples (menos parâmetros).



- 1 Introdução
- 2 Modelo Linear
- 3 LASSO Clássico
- 4 Escolhendo o Valor de  $\lambda$
- 5 Inferência no LASSO**
- 6 LASSO no R
- 7 Abordagem Bayesiana: BLASSO, the Bayesian LASSO
  - Amostrador de Gibbs
  - Escolhendo o  $\lambda$
- 8 Extensões do LASSO



A inferência no LASSO ainda é um tópico de pesquisa aberto. O tópico de pesquisa que engloba esse tipo de problema se chama “inferência após seleção de variáveis”. A razão para o surgimento desta área é que ao selecionarmos as covariáveis do modelo estamos influenciando outras coisas, como o valor-P.

Por exemplo, suponha que temos uma variável resposta e cinco mil variáveis explicativas. Vamos ajustar todos os modelos lineares simples possíveis. Então encontramos uma variável explicativa cuja estimativa do parâmetro associado a ela tenha valor-P igual a 0,01. Será que de fato essa variável explicativa é realmente “importante” para o modelo? Ou há apenas uma associação espúria entre a variável explicativa e a variável resposta, já que é natural esperar que pelo menos uma das cinco mil variáveis terá forte correlação com a variável explicativa?



Existem propostas para a inferência após a aplicação do LASSO, para mais detalhes veja Tibshirani (2015) e Lee et al. (2016). Um pacote também foi desenvolvido para este fim, *selectiveInference*, veja Tibshirani et al. (2017) para mais detalhes.



Após escolher um modelo  $\widehat{M}$  baseado nos dados queremos testar uma hipótese  $\widehat{H}_0$ . Note que  $\widehat{H}_0$  é aleatória, pois  $\widehat{H}_0$  varia em função do modelo selecionado,  $\widehat{M}$ , que por sua vez varia em função da amostra.

Neste contexto temos que controlar o erro seletivo do tipo 1 que é dado por

$$P(T(\mathbf{Y}) \in R | \widehat{M}, \widehat{H}_0) \leq \alpha,$$

em que a região de rejeição é dada pelo evento  $\{T(\mathbf{Y}) \in R\}$ .



- 1 Introdução
- 2 Modelo Linear
- 3 LASSO Clássico
- 4 Escolhendo o Valor de  $\lambda$
- 5 Inferência no LASSO
- 6 LASSO no R**
- 7 Abordagem Bayesiana: BLASSO, the Bayesian LASSO
  - Amostrador de Gibbs
  - Escolhendo o  $\lambda$
- 8 Extensões do LASSO





Software *R* (R Core Team, 2018).

Pacotes:

- *glmnet*, ajusta MLGs com o LASSO (ridge regression também).
- *monomvn*, oferece ajuste com LASSO, BLASSO, Ridge reg., LAR, entre outros.
- *selectiveInference*, dá suporte para realizar inferência em modelos em que o LASSO foi aplicado.



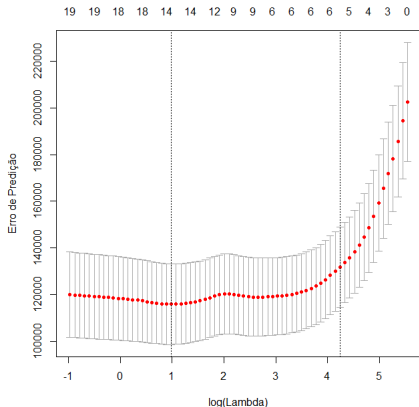
## Um exemplo no R

Para usar o LASSO no R usaremos o pacote *glmnet* que permite o ajuste de modelos lineares, modelos lineares generalizados, modelos de Cox (modelo semi-paramétrico) via LASSO ridge regression e outras penalizações. Nesta seção aplicaremos o LASSO no contexto de modelo linear.

Utilizaremos o conjunto de dados chamado Hitters do pacote *ISLR* do R. Esse conjunto de dados contém 263 observações sobre jogadores de baseball e 20 variáveis. Ajustaremos um modelo de regressão linear via LASSO aos dados, mas antes excluiremos observações com dados faltantes. O nosso objetivo é modelar a variável resposta salário anual dos jogadores por meio do LASSO. Neste caso temos um total de 19 variáveis explicativas.



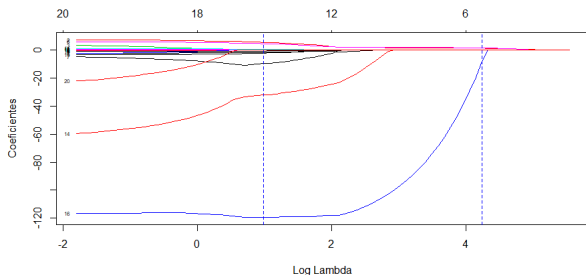
# Um exemplo no R



**Figura:** Gráfico da validação cruzada com  $k = 10$  para o conjunto de dados Hitters na escala logarítmica. Com retas verticais no ponto  $\lambda = 2,674375$  que minimiza o erro de predição e o ponto a um desvio padrão a mais de distância  $\lambda = 69,40069$ .



# Um exemplo no R



**Figura:** Gráfico da estimativa dos parâmetros via LASSO à medida que  $\lambda$  varia na escala logarítmica para o conjunto de dados Hitters. Com retas verticais no ponto  $\lambda = 2,674375$  que minimiza o erro de predição e o ponto a um desvio padrão a mais de distância  $\lambda = 69,40069$ .



# Um exemplo no R

**Tabela:** Estimativas dos parâmetros do modelo ajustado via LASSO com  $\lambda = 69,40069$ .

Parâmetros	Estimativas
Intercepto	127,96
AtBat	0,00
Hits	1,42
HmRun	0,00
Runs	0,00
RBI	0,00
Walks	1,58
Years	0,00
CAtBat	0,00
CHits	0,00
CHmRun	0,00
CRuns	0,16
CRBI	0,34
CWalks	0,00
LeagueA	0,00
LeagueN	0,00
DivisionW	-8,06
PutOuts	0,08
Assists	0,00
Errors	0,00
NewLeagueN	0,00



IME-USP

# Funções no R

```
library(ISLR)
summary(Hitters)

Hitters = na.omit(Hitters)
with(Hitters, sum(is.na(Salary)))

library(glmnet) #carregando pacote
x = model.matrix(Salary~.-1,data=Hitters) #definindo matriz X
y = Hitters$Salary #definindo variável resposta

modelo = glmnet(x,y) #ajustando o modelo via LASSO
valCruz = cv.glmnet(x,y, nfolds=10) #validação cruzada para escolher lambda

plot(valCruz) #gráfico da validação cruzada

coef(valCruz) #extraíndo coeficientes

plot(modelo)
#gráfico para checar o que acontece com as estimativas quando lambda aumenta
```



- 1 Introdução
- 2 Modelo Linear
- 3 LASSO Clássico
- 4 Escolhendo o Valor de  $\lambda$
- 5 Inferência no LASSO
- 6 LASSO no R
- 7 Abordagem Bayesiana: BLASSO, the Bayesian LASSO**
  - Amostrador de Gibbs
  - Escolhendo o  $\lambda$
- 8 Extensões do LASSO



$$\mathbf{y} = \mu \mathbf{1}_n + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

em que  $\mathbf{y}_{n \times 1}$  é o vetor com as variáveis resposta,  $\mathbf{X}_{n \times p}$  é a matriz com as variáveis explicativas (já padronizadas),  $\mu$  é a média geral,  $\boldsymbol{\beta}_{p \times 1}$  é o vetor de parâmetros angulares e  $\boldsymbol{\epsilon}_{n \times 1}$  é o vetor com as fontes de variação.

Sejam  $\hat{\boldsymbol{\beta}}$  o estimador de mínimos quadrados e  $\tilde{\mathbf{y}} = \mathbf{y} - \mu \mathbf{1}_n$ .

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^\top (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})$$





Tibshirani (1996) notou que a estimativa do LASSO pode ser vista como a moda a posteriori quando os  $\beta_j | \sigma^2$  têm distribuição Laplace( $0, \sqrt{\sigma^2}/\lambda$ ) independentes, isto é,

$$\pi(\boldsymbol{\beta} | \sigma^2) = \prod_{j=1}^p \frac{\lambda}{2\sqrt{\sigma^2}} \exp\left(-\frac{\lambda|\beta_j|}{\sqrt{\sigma^2}}\right) = \frac{\lambda^p}{2^p \sqrt{\sigma^2}^p} \exp\left(-\frac{\lambda}{\sqrt{\sigma^2}} \sum_{j=1}^p |\beta_j|\right).$$

Obs: condicionar a  $\sigma^2$  é importante, pois garante que a posteriori  $\pi(\boldsymbol{\beta}, \sigma^2 | \tilde{\mathbf{y}})$  será unimodal.

Se a posteriori não fosse unimodal qual seria a moda a posteriori?

Outra consequência da falta de unimodalidade é a lenta convergência do amostrador de Gibbs.



# Consequência de não condicionar a distribuição de $\beta$

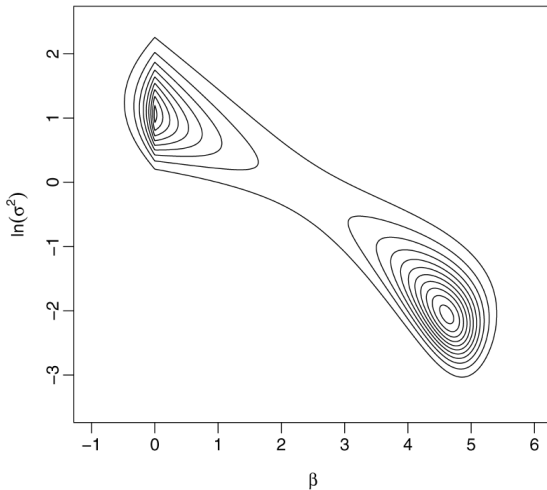


Figura: Fonte: Park e Casella (2008).



Park e Casella (2008) propuseram o amostrador de Gibbs partindo da representação da Laplace pela mistura na escala da normal (com densidade exponencial), ou seja,

$$\frac{a}{2}e^{-a|z|} = \int_0^{+\infty} \frac{1}{\sqrt{2\pi s}} e^{-z^2/(2s)} \frac{a^2}{2} e^{-a^2 s/2} ds, \quad a > 0. \quad (6)$$

Perceba que  $Z|S = s \sim N(0, s)$  e  $S \sim \text{Exp}(a^2/2)$ .



A estrutura hierárquica do modelo é dada por

$$\mathbf{y} | \mu, \mathbf{X}, \boldsymbol{\beta}, \sigma^2 \sim N_p \left( \mu \mathbf{1}_n + \mathbf{X} \boldsymbol{\beta}, \sigma^2 \mathbf{I}_n \right),$$

$$\boldsymbol{\beta} | \sigma^2, \tau_1^2, \dots, \tau_p^2 \sim N_p \left( \mathbf{0}_p, \sigma^2 \mathbf{D}_\tau \right), \mathbf{D}_\tau = \text{diag} \left( \tau_1^2, \dots, \tau_p^2 \right)$$

$$\sigma^2 \sim \pi \left( \sigma^2 \right) \perp \tau_j^2 \stackrel{iid}{\sim} \text{Exp} \left( \frac{\lambda^2}{2} \right)$$

Note que para obter a distribuição de  $\boldsymbol{\beta} | \sigma^2$  basta integrar o produto de  $N_p \left( \mathbf{0}_p, \sigma^2 \mathbf{D}_\tau \right)$  por  $\prod_{j=1}^p \text{Exp} \left( \frac{\lambda^2}{2} \right)$  em relação a  $\tau_1^2, \dots, \tau_p^2$ .



$$\pi(\beta|\sigma^2) = \int \prod_{j=1}^p \frac{1}{\sqrt{2\pi\sigma^2\tau_j^2}} \exp\left(-\frac{\beta_j^2}{2\sigma^2\tau_j^2}\right) \prod_{j=1}^p \frac{\lambda^2}{2} \exp\left(-\tau_j^2 \frac{\lambda^2}{2}\right) d\tau$$

$$\pi(\beta|\sigma^2) = \prod_{j=1}^p \int_0^\infty \frac{1}{\sqrt{2\pi\sigma^2\tau_j^2}} \exp\left(-\frac{\beta_j^2}{2\sigma^2\tau_j^2}\right) \frac{\lambda^2}{2} \exp\left(-\tau_j^2 \frac{\lambda^2}{2}\right) d\tau_j^2,$$

Perceba que em cada integral podemos usar (6) com  $\beta_j|\sigma^2, \tau_j^2 = \tau_j^* \sim N(0, \sigma^2\tau_j^2)$  e  $\tau_j^2 \sim \text{Exp}(\lambda^2/2)$ .

$$\pi(\beta|\sigma^2) = \prod_{j=1}^p \frac{\lambda}{2\sqrt{\sigma^2}} \exp\left(-\frac{\lambda}{\sqrt{\sigma^2}}|\beta_j|\right).$$



# Porque o BLASSO é obtido a partir da Moda a Posteriori?

Note que

$$\ln \left( \pi \left( \beta | \sigma^2, \tilde{\mathbf{y}} \right) \right) = \ln \left( f \left( \tilde{\mathbf{y}} | \beta, \sigma^2 \right) \right) + \ln \left( \pi \left( \beta | \sigma^2 \right) \right) + C$$

$$\ln \left( \pi \left( \beta | \sigma^2, \tilde{\mathbf{y}} \right) \right) = -\frac{1}{2\sigma^2} \sum_{i=1}^n \left( \tilde{\mathbf{y}} - \mathbf{x}_i^\top \beta \right)^2 - \frac{\lambda}{\sqrt{\sigma^2}} \sum_{j=1}^p |\beta_j| + C,$$

queremos maximizar a posteriori para encontrar a moda, logo

$$\beta_{BLASSO} = \arg \min_{\beta} \sum_{i=1}^n \left( \tilde{\mathbf{y}} - \mathbf{x}_i^\top \beta \right)^2 + 2\sqrt{\sigma^2} \lambda \sum_{j=1}^p |\beta_j|.$$



# Distribuição dos Parâmetros a Posteriori ( $\pi(\sigma^2) = 1/\sigma^2$ )

$$\beta | \tilde{\mathbf{y}} \sim N(\mathbf{A}^{-1} \mathbf{X}^T \tilde{\mathbf{y}}, \sigma^2 \mathbf{A}^{-1}), \quad \mathbf{A} = \mathbf{X}^T \mathbf{X} + \mathbf{D}_\tau^{-1}$$

$$\sigma^2 | \tilde{\mathbf{y}} \sim Gl(a, b), \quad a = \frac{n + p - 1}{2} \text{ e } b = \frac{1}{2} \sum_{i=1}^n (\tilde{y}_i - \mathbf{x}_i^T \beta)^2 + \frac{1}{2} \beta^T \mathbf{D}_\tau^{-1} \beta$$

$$\frac{1}{\tau_j^2} | \tilde{\mathbf{y}} \stackrel{ciid}{\sim} NI(\mu', \lambda'), \quad \mu' = \sqrt{\frac{\lambda^2 \sigma^2}{\beta_j^2}} \text{ e } \lambda' = \lambda^2,$$

cuja densidade é dada por

$$f(w) = \sqrt{\frac{\lambda'}{2\pi}} w^{-3/2} \exp\left[-\frac{\lambda'(w - \mu')^2}{2(\mu')^2 w}\right], \quad w > 0$$



Existem duas abordagens propostas para “escolher o  $\lambda$ ”

- Usar Bayes empírico via máxima verossimilhança marginal (algoritmo MCEM);
- Atribuir uma priori gama “não informativa” para  $\lambda^2$  (não  $\lambda$  !)  
É interessante atribuir priori gama para  $\lambda^2$  devido a conjugação resultante.

Obs: não usar priori imprópria como  $\pi(\lambda^2) = 1/\lambda^2$ , pois isto acarreta posteriori imprópria.





## Escolhendo o $\lambda$ :

Casella (2001) propôs um algoritmo MCEM que complementa o amostrador de Gibbs e dá estimativas de máxima verossimilhança dos hiperparâmetros.

$$\lambda^{(k)} = \sqrt{\frac{2p}{\sum_{j=1}^p E_{\lambda^{(k-1)}} [\tau_j^2 | \tilde{\mathbf{y}}]}},$$

em que o valor inicial (sugerido) é

$$\lambda^{(0)} = p \sqrt{\hat{\sigma}_{MQ}^2} / \sum_{j=1}^p |\hat{\beta}_j^{MQ}|$$

e a esperança condicional é substituída pela média do amostrador de Gibbs.



- Em geral, os resultados provenientes do LASSO e do BLASSO são muito similares.
- Embora o BLASSO seja computacionalmente intensivo ele é mais fácil de implementar e também tem a vantagem de gerar estimativas intervalares (e erro padrão) automaticamente durante o processo (afinal teremos as distribuições a posteriori dos parâmetro).
- Não há como calcular os erros padrão analiticamente via LASSO (exceto aproximações analíticas). Por isso para obter os erros padrão e assim poder construir estimativas intervalares é necessário utilizar Bootstrap.
- Além das formas de escolher o  $\lambda$  no caso do LASSO (validação cruzada, dentre outros) o BLASSO oferece mais duas formas de escolher o  $\lambda$ .



- 1 Introdução
- 2 Modelo Linear
- 3 LASSO Clássico
- 4 Escolhendo o Valor de  $\lambda$
- 5 Inferência no LASSO
- 6 LASSO no R
- 7 Abordagem Bayesiana: BLASSO, the Bayesian LASSO
  - Amostrador de Gibbs
  - Escolhendo o  $\lambda$
- 8 Extensões do LASSO



- Bridge Regression.

$$\hat{\beta}_{Bridge} = \arg \min_{\beta} \left[ (\mathbf{Y} - \mathbf{X}\beta)^{\top} (\mathbf{Y} - \mathbf{X}\beta) + \lambda \sum_{j=1}^p |\beta_j|^q \right], \lambda \geq 0 \text{ (fixo)},$$

priori adequada

$$\pi(\beta|\sigma^2) \propto \prod_{j=1}^p \exp \left[ -\lambda \left( |\beta_j|/\sqrt{\sigma^2} \right)^q \right],$$

isto é, exponenciais potência estão fazendo o papel da exponencial dupla neste caso.



- LASSO “Huberizado” (LASSO Robusto).

$$\hat{\beta}_H = \arg \min_{\beta} \left[ L(\mathbf{Y} - \mathbf{X}\beta) + \lambda \sum_{j=1}^p |\beta_j| \right], \quad \lambda \geq 0 \text{ (fixo),}$$

em que  $L(\cdot)$  é uma função de perda do tipo Huber que é quadrática na vizinhança de zero e linear crescente fora da vizinhança de zero.



# Usando os Pacotes monomvn e glmnet

```
# Carregando os pacotes
library(monomvn); library(lars); library(glmnet); library(miscTools)
data(diabetes); attach(diabetes)

# definindo o número de iterações descartadas (burn-in), número de amostras do mcmc e valores iniciais
burnin <- 500
iter <- 1000
initial.beta <- rep(-500, dim(x2)[2]) # atribuindo um valor inicial extremo para todos os betas
initial.lambda2 <- 10 # atribuindo um valor inicial extremo para lambda (parâmetro de penalização)
initial.variance <- 500 # atribuindo um valor inicial extremo para o parâmetro de variância

# Iniciando o amostrador de Gibbs
lasso <- blasso(X = x2, # matriz de covariáveis 442 x 64
               y = y, # variáveis resposta of 442
               T = iter, # número de iterações
               beta = initial.beta,
               lambda2 = initial.lambda2,
               s2 = initial.variance)
#rd = c(1, 1.78)) # hiperparâmetros sugeridos por Park e Casella (2008)

# extraindo valores de alguns parâmetros para visualização
coef.lasso <- as.data.frame(cbind(iter = seq(iter),
                                beta1 = lasso$beta[, "b.1"],
                                beta2 = lasso$beta[, "b.2"],
                                variance = lasso$s2,
                                lambda.square = lasso$lambda2))

colMedians(coef.lasso[-seq(burnin), -1])
#beta1      beta2      variance      lambda.square
#0.0000000 -172.3840906 2841.4410472 0.3031814
```



IME-USP

# Usando os Pacotes monomvn e glmnet

```
#####lasso clássico  
  
#Vamos comparar o LASSO (glmnet) com o BLASSO (monomvn)  
  
fit.glmnet <- glmnet(as.matrix(x2), y,  
                    lambda=cv.glmnet(as.matrix(x2), y)$lambda.1se)  
coef.glmnet <- coef(fit.glmnet)  
sum(coef.glmnet == 0)  
#53  
sum(colMedians(lasso$beta[-seq(burnin), ]) == 0)  
#56
```

O LASSO atribuiu o valor zero para 53 parâmetros. E o BLASSO atribuiu o valor zero para 56 parâmetros.





Casella, G. (2001).

Empirical Bayes gibbs sampling.

*Bioinformatics*, **20**, 3423-3430.



Friedman, J.; Hastie, T.; Tibshirani, R. (2010).

Regularization Paths for Generalized Linear Models via Coordinate Descent .

*Journal of Statistical Software*, **33**(1), 1-22.



Hastie, T.; Tibshirani, R.; Friedman, J. (2008).

*The Elements of Statistical Learning: Data Mining, Inference, and Prediction* 2ed .

New York: Springer-Verlag.



Hastie, T.; Tibshirani, R.; Wainwright, M. (2015).

*Statistical Learning with Sparsity: The Lasso and Generalizations*.

Chapman & Hall/CRC.







James, G.; Witten, D; Hastie, T.; Tibshirani, R. (2014).

*An Introduction to Statistical Learning: With Applications in R*, 2nd ed..  
pringer Publishing Company.



Lee, J. D.; Sun, D. L.; Sun, Y.; Taylor, J. E. (2016).

Exact post-selection inference, with application to the lasso  
*The Annals of Statistics*, **44**(3), 907-927.



Park, T.; Casella, G. (2008).

The Bayesian Lasso.

*Journal of the American Statistical Association*, **103**, 681-686.



R Core Team (2018).

R: A language and environment for statistical computing.

R Foundation for Statistical Computing, Vienna, Austria. URL  
<https://www.R-project.org/>.





Taylor, J.; Tibshirani, R. J. (2015).

Statistical learning and selective inference.

*Proceedings of the National Academy of Sciences*. **112**(25), 7629-7634.



Tibshirani, R. (1996).

Regression Shrinkage and Selection via the LASSO.

*Journal of the Royal Statistical Society. Series B (Methodological)*, **58**, 267-288.



Tibshirani, R.; Tibshirani, R.; Taylor, J.; Loftus, J.; Reid, S. (2017).

*selectiveInference: Tools for Post-Selection Inference* (R package version 1.2.4)

url = <https://CRAN.R-project.org/package=selectiveInference>



Obrigado pela atenção!!!



IME-USP