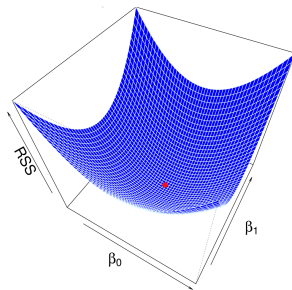
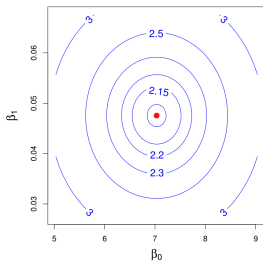


# Regularização

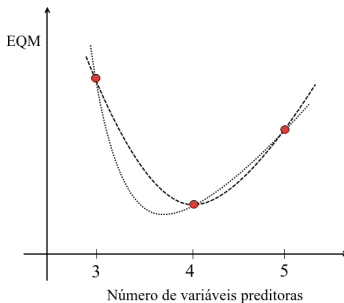
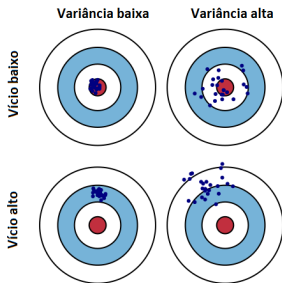
Eduardo Vargas Ferreira

- Quando estimamos os parâmetros da regressão por **mínimos quadrados ordinários**, estamos interessados na seguinte minimização:

$$\min \{J[y_i, h(x)]\} = \min \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\}$$



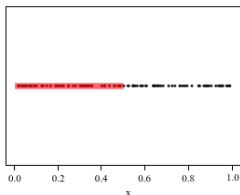
- Uma solução é flexibilizar o modelo admitindo certo vício das estimativas.



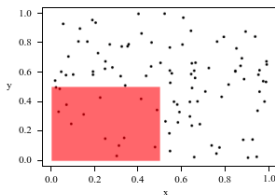
$$E[y_0 - h(x_0)]^2 = \text{Var}[h(x_0)] + [\text{Vício}(h(x_0))]^2 + \text{Var}(\varepsilon).$$

- Em situações com “*small n, large p*” a maioria dos métodos modernos de análise de dados falha, por diferentes razões, p. ex.:
  - ★ **Modelos Lineares Generalizados:** falham, pois a matriz do modelo não tem posto completo;
  - ★ **Random Forests** falham, pois a probabilidade de selecionar variáveis importantes diminui muito.
  - ★ **Análise de Clusters** e métodos baseados em distâncias no plano cartesiano falham devido à “**maldição da dimensionalidade**”.

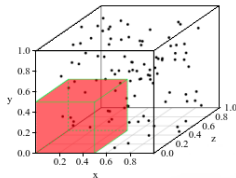
42% dos dados capturados



14% dos dados capturados



7% dos dados capturados



# Regularização

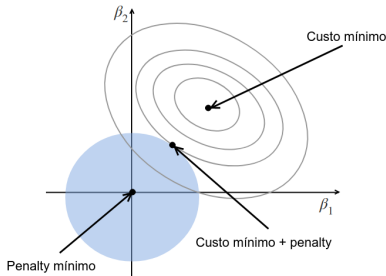
# Uma solução é a Regularização



- Regulamos o número de variáveis, impondo um custo ao algoritmo:

$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \quad \text{sujeito a } g(\beta) < t, \text{ com } t > 0$$

- $g(\beta)$  representa a **função penalty** (*shrinkage penalty*).



- Estamos diante de uma otimização com restrição:

$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \quad \text{sujeito a } g(\beta) < t, \text{ com } t > 0$$

- Fazemos isso aumentando a função objetivo:

$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda g(\beta), \text{ com } \lambda > 0$$

- $t$  e  $\lambda$  são inversamente proporcionais.

- Temos a função objetivo aumentada:

$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda g(\beta), \text{ com } \lambda > 0$$

- Utilizaremos a **família das potências** para penalizar o modelo:

$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|^q$$

- $\lambda$  é o **tuning parameter**, determinado separadamente. Ele controla o impacto do *penalty*,  $g(\beta)$ , nas estimativas dos parâmetros.



- Os coeficientes obtidos por mínimos quadrados ordinários são equivariantes por transformação de escala;
- Na regressão penalizada,  $\mathbf{X}_j \hat{\beta}_{j,\lambda}^{restrito}$  depende não somente de  $\lambda$ , mas da escala do  $j$ -ésimo preditor.

$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|^q$$

- Diante disso, deve-se padronizar os preditores na regressão penalizada:

$$\tilde{x}_{ij} = \frac{x_{ij}}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}}.$$

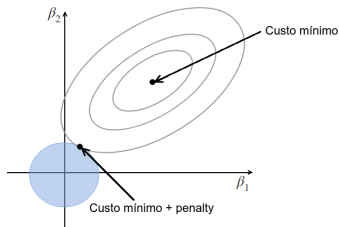
# Penalização Ridge

- Neste caso, o problema de otimização utiliza o *penalty*  $\ell^2$ :

$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

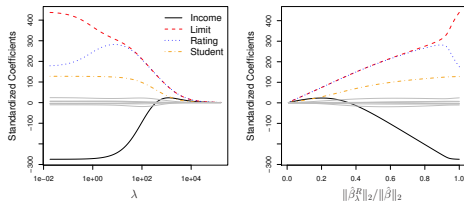
- No caso de vetores em  $\mathbf{X}$  ortonormais

$$\hat{\beta}_{\lambda}^R = \frac{\hat{\beta}^{OLS}}{1 + \lambda}.$$

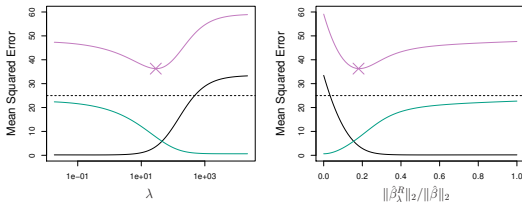


- Os  $\beta$ 's podem ser próximos de zero, mas não assumem esse valor.

- Inadimplência no cartão de crédito:** o objetivo é prever se um cliente será ou não inadimplente no próximo mês.



- Exemplo simulado:** foi gerado utilizando 45 preditores relacionados à resposta, em que nenhum dos verdadeiros  $\beta_1, \beta_2, \dots, \beta_{45}$  eram zero.

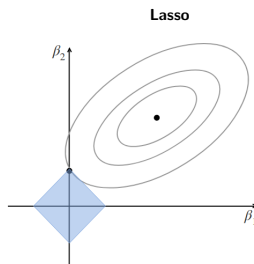
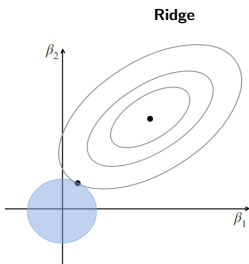


# Penalização Lasso

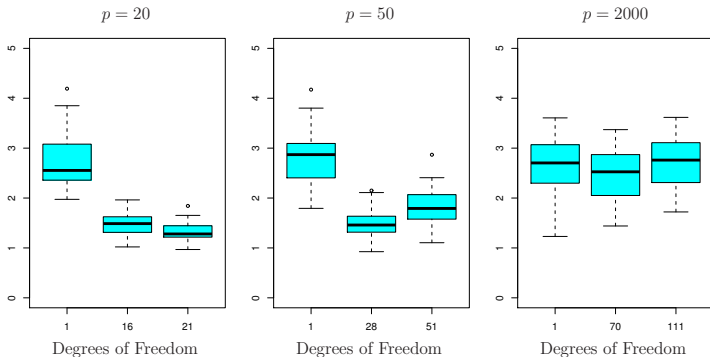
- Os coeficientes *Lasso*,  $\hat{\beta}_{\lambda}^L$ , minimizam a quantidade

$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

- O *penalty*  $\ell^1$  funciona também como um **selecionador de variáveis**.



- No exemplo abaixo temos 100 observações com  $p = 20, 50$  e  $2000$ , das quais apenas 20 são relacionadas com a resposta;

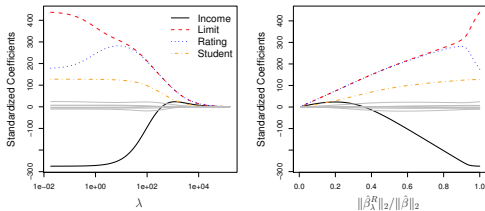


- Os graus de liberdade, substituem o  $\lambda$ , e representam o número de parâmetros estimados não nulos.

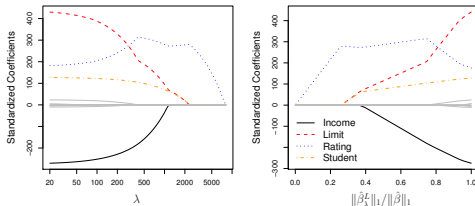
# Exemplo: Inadimplência no cartão de crédito



- Regressão Ridge



- Regressão Lasso

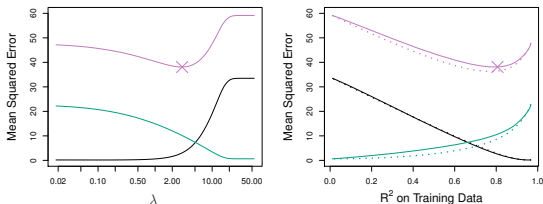




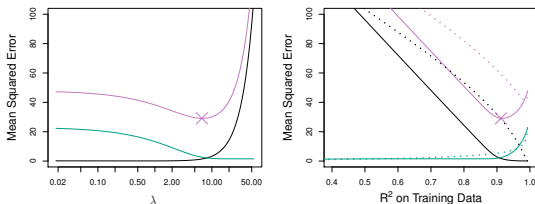
# Comparação entre Ridge e Lasso



- Este exemplo foi gerado utilizando 45 preditores relacionados à resposta, em que nenhum dos verdadeiros coeficientes  $\beta_1, \beta_2, \dots, \beta_{45}$  eram zero;



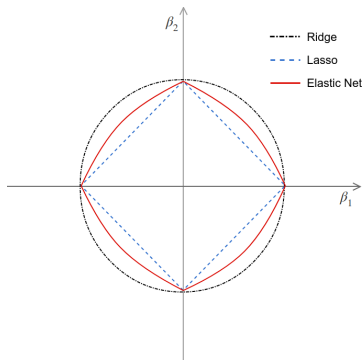
- No exemplo seguinte, a resposta é função apenas de 2 dos 45 preditores. I.e., 43 dos verdadeiros coeficientes  $\beta_1, \beta_2, \dots, \beta_{45}$  eram zero.



# Penalização Elastic net

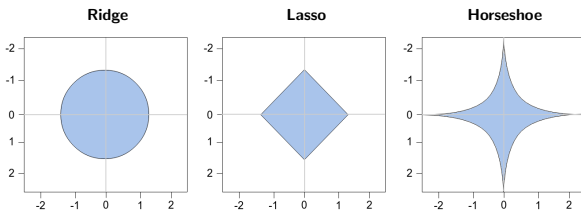
- **Elastic net** é um compromisso entre a regressão *Ridge* e *Lasso*. Os coeficientes elastic net,  $\hat{\beta}_{\lambda}^E$ , minimizam a quantidade

$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \left( \alpha |\beta_j| + (1 - \alpha) \beta_j^2 \right)$$



## **Penalização horseshoe**

- Ela favorece mais ainda a presença de 0's (maior esparsidade);
- Ou seja, tende a encontrar as elipses geradas pelos mínimos quadrados em cima dos eixos com mais frequência que *Ridge* e *Lasso*;



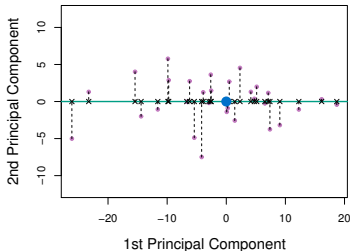
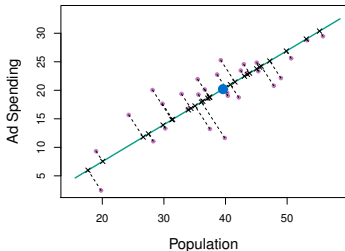
- E quando  $q = 0$  voltamos ao **Best subset selection**.

## **Regressão com Componentes Principais**

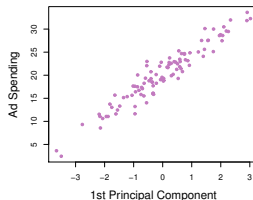
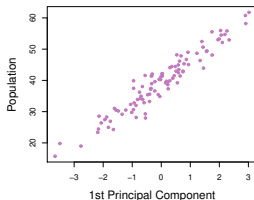
- Suponha que temos a informação sobre a população (**pop**) em 100 cidades dos EUA e gasto com propaganda de determinada empresa (**ad**);
- Podemos resumir essas variáveis em apenas uma,  $Z_1$ , da seguinte forma:

$$\begin{aligned} Z_1 &= \phi_{11} (\text{pop} - \overline{\text{pop}}) + \phi_{21} (\text{ad} - \overline{\text{ad}}) \\ &= 0.839 \times (\text{pop} - \overline{\text{pop}}) + 0.544 \times (\text{ad} - \overline{\text{ad}}) . \end{aligned}$$

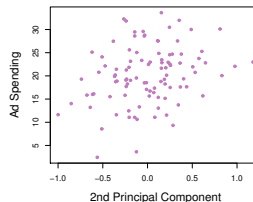
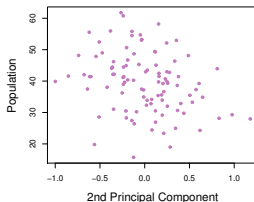
- Tal que  $\text{Var}(Z_1)$  seja o máximo possível.



$$Z_1 = 0.839 \times (\text{pop} - \overline{\text{pop}}) + 0.544 \times (\text{ad} - \overline{\text{ad}})$$



$$Z_2 = 0.839 \times (\text{pop} - \overline{\text{pop}}) + 0.544 \times (\text{ad} - \overline{\text{ad}})$$





- Considere que  $Z_1, Z_2, \dots, Z_M$  ( $M < p$ ) represente as combinações lineares das variáveis originais  $X_j$ ,  $j = 1, \dots, p$ :

$$Z_m = \sum_{j=1}^p \phi_{jm} X_j, \text{ com } m = 1, \dots, M.$$

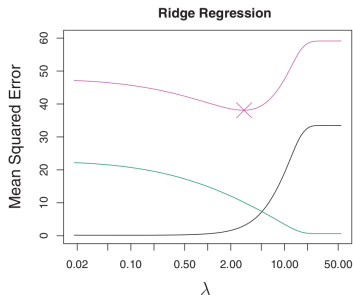
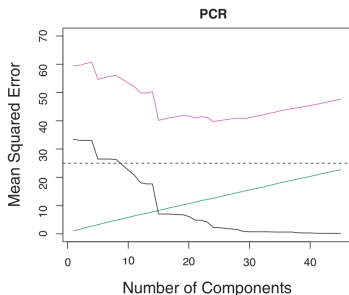
- Assumindo que direção de máxima variabilidade é a mesma associada ao  $Y$ , podemos construir uma regressão utilizando as C.P's, da forma:

$$y_i = \theta_0 + \sum_{m=1}^M \theta_m z_{im} + \varepsilon_i, \text{ com } i = 1, \dots, n.$$

- Mas, note que

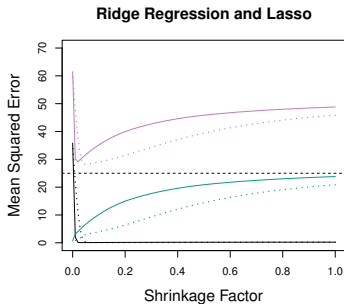
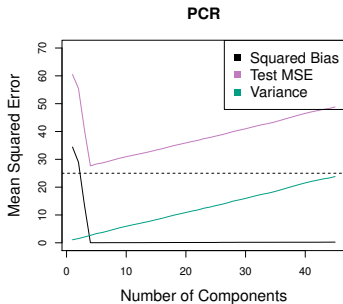
$$\sum_{m=1}^M \theta_m z_{im} = \sum_{m=1}^M \theta_m \sum_{j=1}^p \phi_{jm} x_{ij} = \sum_{j=1}^p \sum_{m=1}^M \theta_m \phi_{jm} x_{ij} = \sum_{j=1}^p \beta_j x_{ij}.$$

- Foi gerado utilizando 45 preditores relacionados à resposta, em que nenhum dos verdadeiros  $\beta_1, \beta_2, \dots, \beta_{45}$  eram zero.



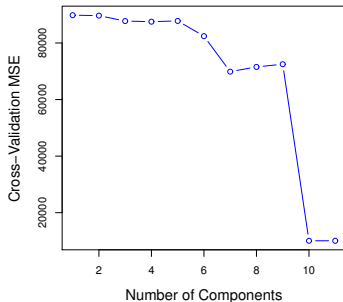
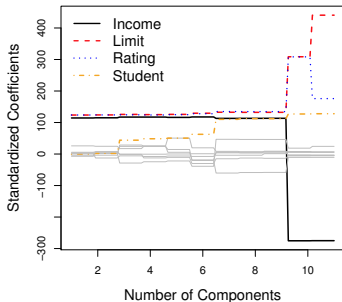
- Os resultados via PCR e Ridge são muito próximos. Pode-se dizer até que Ridge é uma versão contínua de PCR.

- Os dados foram gerados, tal que a resposta dependesse apenas das 5 primeiras componentes principais.



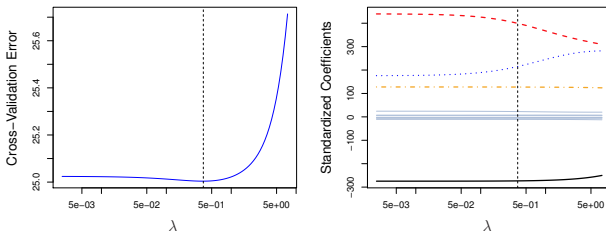
- A PCR se sai melhor quando poucas componentes são necessárias para explicar a maior parte da variabilidade dos dados.

- Nestes dados o menor erro de validação foi obtido quando  $M = 10$ . Equivalente aos mínimos quadrados.

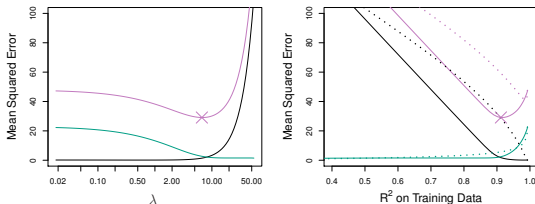


## **Seleccionando o tuning parameter**

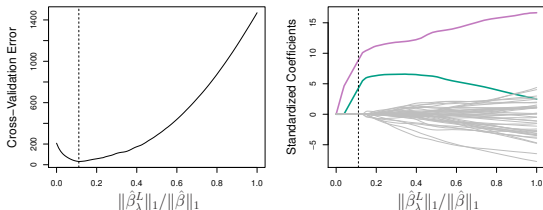
- **Validação cruzada** fornece uma maneira simples de resolver este problema:
  - (a) A partir de uma grade de valores de  $\lambda$ , calculamos a taxa de erro de validação (para cada  $\lambda$ );
  - (b) Escolhemos o valor de  $\lambda$  que fornece a menor taxa de erro;
  - (c) Ajustamos novamente o modelo, utilizando todas as observações disponíveis, com o valor de  $\lambda$  encontrado anteriormente.



- Este exemplo foi gerado utilizando 45 preditores, em que 43 dos verdadeiros coeficientes  $\beta_1, \beta_2, \dots, \beta_{45}$  eram zero.



- Aplicamos validação cruzada 10 – *fold* para selecionar o melhor  $\lambda$ .



- James, G., Witten, D., Hastie, T. e Tibshirani, An Introduction to Statistical Learning, 2013;
- Hastie, T., Tibshirani, R. e Friedman, J., The Elements of Statistical Learning, 2009;
- Lantz, B., Machine Learning with R, Packt Publishing, 2013;
- Tan, Steinbach, and Kumar, Introduction to Data Mining, Addison-Wesley, 2005;
- Some of the figures in this presentation are taken from "An Introduction to Statistical Learning, with applications in R" (Springer, 2013) with permission from the authors: G. James, D. Witten, T. Hastie and R. Tibshirani