

Lecture 2: Sparse Regression

24/03/2023

Lecturer: André Ramos

Scribe: Mateus Waga

1 Introdução

Nesta aula foi apresentado o capítulo 3 do livro [1], baseado nos artigos [2, 3, 4].

Dado um conjunto de dados com n observações e p variáveis preditoras, denotadas por $\mathbf{X} = [x_1, x_2, \dots, x_p]$ e uma variável de destino $\mathbf{Y} = [y_1, y_2, \dots, y_n]$, o modelo de regressão linear com regularização é definido como

$$\begin{aligned} \min_{\boldsymbol{\beta}} \quad & \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \frac{1}{2\gamma} \|\boldsymbol{\beta}\|_2^2 \\ \text{s.t.} \quad & \|\boldsymbol{\beta}\|_0 \leq k \end{aligned}$$

onde $\boldsymbol{\beta} = [\beta_1, \beta_2, \dots, \beta_p]$ é um vetor de coeficientes de regressão e γ é um parâmetro de regularização que controla a força da penalidade de regularização e $\|\boldsymbol{\beta}\|_0$ é o número de componentes não zero.

Se $n < p$ O problema sem a restrição de norma 0 possui infinitas , logo o problema não está bem definido. Se $n \geq p$ o problema é bem definido mas a restrição adicional torna o problema mais simples e mais interpretável.

A seguir, descreveremos alguns métodos para solucionar o problema.

2 Método Primal

Trata-se de uma reformulação do problema inicial como um MIP quadrático. É introduzida a variável binária s_j que é igual a 1 se a variável x_j é selecionada e 0 caso contrário.

$$\begin{aligned} \min \quad & \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \frac{1}{2\gamma} \|\boldsymbol{\beta}\|_2^2 \\ \text{s.t.} \quad & \boldsymbol{\beta} \in \mathbb{R}^p, \mathbf{s} \in \{0, 1\} \\ & \mathbf{1}'\mathbf{s} \leq k \\ & M^- s_j \leq \beta_j \leq M^+ s_j, \quad \forall j \in [p] \end{aligned}$$

A escolha do M afeta o desempenho dos métodos de solução. Assim para casos onde $n \geq p$ podemos utilizar a seguinte formulação para encontrar valores de M mais "apertados":

$$M^- = \min_{\beta} \beta_j$$

$$\text{s.t. } \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \frac{1}{2\gamma} \|\beta\|_2^2 \leq \text{UB}$$

$$M^+ = \max_{\beta} \beta_j$$

$$\text{s.t. } \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \frac{1}{2\gamma} \|\beta\|_2^2 \leq \text{UB}$$

Onde UB é uma solução viável do problema original (escolher a priori k variáveis não zero).

Uma alternativa proposta no livro, para não utilizar o M é considerar:

$$g(\beta) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \frac{1}{2\gamma} \|\beta\|_2^2,$$

com

$$\nabla g(\beta) = -\mathbf{X}^T(\mathbf{y} - \mathbf{X}\beta) + \frac{1}{\gamma}\beta.$$

uma vez que $g(\beta) \geq 0$ e possui um gradiente Lipschitz contínuo, temos:

$$\|\nabla g(\beta) - \nabla g(\tilde{\beta})\| \leq \ell \|\beta - \tilde{\beta}\|.$$

Assim, considere:

$$\hat{\beta} \in \arg \min_{\|\beta\|_0 \leq k} \|\beta - \mathbf{c}\|_2^2$$

e a função $H_k(c)$

$$H_k(c) = \hat{\beta}_i = \begin{cases} c_i, & \text{if } i \in \{(1), \dots, (k)\} \\ 0, & \text{otherwise.} \end{cases}$$

Abaixo o algoritmo proposto no livro, ele considera como input: $g(\beta)$, L , ϵ e uma solução inicial $\beta_1 \in \mathbb{R}^p$ tal que $\|\beta_1\|_0 \leq k$. Output: uma solução estacionária de primeira ordem β^* .

- 1: $m \leftarrow 1$
- 2: repeat
- 3: $m \leftarrow m + 1$
- 4: $\eta_m \leftarrow \mathbf{H}_k(\beta_{m-1} - \frac{1}{L} \nabla g(\beta_{m-1}))$
- 5: $\lambda_m \leftarrow \arg \min_{\lambda} g(\lambda \eta_m + (1 - \lambda) \beta_{m-1})$
- 6: $\beta_m \leftarrow \lambda_m \eta_m + (1 - \lambda_m) \beta_{m-1}$
- 7: until $g(\beta_m) - g(\beta_{m-1}) \leq \epsilon$
- 8: return β_m

3 Método Dual

A ideia deste método é utilizar uma solução dual que encontra o vetor \mathbf{s} e "plugar" estes valores no modelo original. O modelo utilizado para encontrar o \mathbf{s} é:

$$\begin{aligned} \min \quad & \frac{1}{2} \mathbf{y}^T \left(\mathbf{I}_n + \gamma \sum_{j \in [p]} s_j \mathbf{K}_j \right)^{-1} \mathbf{y} \\ \text{s.t.} \quad & \mathbf{s} \in \mathbf{S}_k^p, \end{aligned}$$

onde \mathbf{K}_j em \mathbf{S}_+^n é definido como:

$$\mathbf{K}_j = \mathbf{X}_j \mathbf{X}_j^T$$

References

- [1] D. Bertsimas and J. Dunn. *Machine learning under a modern optimization lens*. Dynamic Ideas LLC Charlestown, MA, 2019.
- [2] D. Bertsimas and A. King. Or forum—an algorithmic approach to linear regression. *Operations Research*, 64(1):2–16, 2016.
- [3] D. Bertsimas, A. King, and R. Mazumder. Best subset selection via a modern optimization lens. 2016.
- [4] D. BERTSIMAS and B. VAN PARYS. Sparse high-dimensional regression: Exact scalable algorithms and phase transitions. *The Annals of Statistics*, 48(1):300–323, 2020.