Regularización: Lasso

Big Data y Machine Learning para Economía Aplicada

Ignacio Sarmiento-Barbieri

Universidad de los Andes

Agenda

- 1 Recap
- 2 Lasso
- 3 Familia de regresiones penalizadas
- 4 k > n
- 5 Elastic Net

Agenda

- 1 Recap
- 2 Lasso
- 3 Familia de regresiones penalizadas
- 4 k > n
- 5 Elastic Net

Regularización: Motivación

- Las técnicas econometricas estándar no están optimizadas para la predicción porque se enfocan en la insesgadez.
- ▶ OLS por ejemplo es el mejor estimador lineal *insesgado*
- lacktriangle OLS minimiza el error "dentro de muestra", eligiendo eta de forma tal que

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_{i} - \beta_{0} - x_{i1}\beta_{1} - \dots - x_{ip}\beta_{p})^{2}$$
(1)

- pero para predicción, no estamos interesados en hacer un buen trabajo dentro de muestra
- Queremos hacer un buen trabajo, fuera de muestra

- (ロ) (部) (注) (注) 注 り(()

Ridge

- Asegurar cero sesgo dentro de muestra crea problemas fuera de muestra: trade-off Sesgo-Varianza
- Las técnicas de machine learning fueron desarrolladas para hacer este trade-off de forma empírica.
- ▶ Vamos a proponer modelos del estilo

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - \beta_0 - x_{i1}\beta_1 - \dots - x_{ip}\beta_p)^2 + \lambda \sum_{j=1}^{p} R(\beta_j)$$
 (2)

- ▶ donde *R* es un regularizador que penaliza funciones que crean varianza
- ► Explícitamente en la minimización incluimos un termino de sesgo y un termino de varianza.

Ridge

Para un $\lambda \geq 0$ dado, consideremos ahora el siguiente problema de optimización

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_{i} - \beta_{0} - x_{i1}\beta_{1} - \dots - x_{ip}\beta_{p})^{2} + \lambda \sum_{j=1}^{p} (\beta_{j})^{2}$$

$$\sum_{i=1}^{n+p} (y_{i} - \beta_{0} - x_{i1}\beta_{1} - \dots - x_{ip}\beta_{p})^{2} + \lambda \sum_{j=1}^{p} (\beta_{j})^{2}$$

$$\sum_{i=1}^{n+p} (y_{i} - \beta_{0} - x_{i1}\beta_{1} - \dots - x_{ip}\beta_{p})^{2} + \lambda \sum_{j=1}^{p} (\beta_{j})^{2}$$

$$\sum_{i=1}^{n+p} (y_{i} - \beta_{0} - x_{i1}\beta_{1} - \dots - x_{ip}\beta_{p})^{2} + \lambda \sum_{j=1}^{p} (\beta_{j})^{2}$$

$$\sum_{i=1}^{n+p} (y_{i} - \beta_{0} - x_{i1}\beta_{1} - \dots - x_{ip}\beta_{p})^{2} + \lambda \sum_{j=1}^{p} (\beta_{j})^{2}$$

$$\sum_{i=1}^{n+p} (y_{i} - \beta_{0} - x_{i1}\beta_{1} - \dots - x_{ip}\beta_{p})^{2} + \lambda \sum_{j=1}^{p} (\beta_{j})^{2}$$

$$\sum_{i=1}^{n+p} (y_{i} - \beta_{0} - x_{i1}\beta_{1} - \dots - x_{ip}\beta_{p})^{2} + \lambda \sum_{j=1}^{p} (\beta_{j})^{2}$$

$$\sum_{i=1}^{n+p} (y_{i} - \beta_{0} - x_{i1}\beta_{1} - \dots - x_{ip}\beta_{p})^{2} + \lambda \sum_{j=1}^{p} (\beta_{j})^{2}$$

Sarmiento-Barbieri (Uniandes)

Agenda

- 1 Recap
- 2 Lasso
- 3 Familia de regresiones penalizadas
- 4 k > n
- 5 Elastic Net

Lasso



lacktriangle Para un $\lambda \geq 0$ dado, consideremos el siguiente problema de optimización

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - \beta_0 - x_{i1}\beta_1 - \dots - x_{ip}\beta_p)^2 + \lambda \sum_{j=1}^{p} |\beta_j|$$

$$|\beta_j - \emptyset|$$
(4)

Lasso

lacktriangle Para un $\lambda \geq 0$ dado, consideremos el siguiente problema de optimización

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - \beta_0 - x_{i1}\beta_1 - \dots - x_{ip}\beta_p)^2 + \lambda \sum_{j=1}^{p} |\beta_j|$$
 (4)

- ► "LASSO's free lunch": selecciona automáticamente los predictores que van en el modelo $(\beta_j \neq 0)$ y los que no $(\beta_j = 0)$
- Por qué? Los coeficientes que no van son soluciones de esquina
- $ightharpoonup L(\beta)$ es no differentiable



Sarmiento-Barbieri (Uniandes)

Lasso Intuición en 1 Dimension

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - x_i\beta)^2 + \lambda|\beta|$$
 (5)

- ► Un solo predictor, un solo coeficiente, estondor who < = (x:-x)! = 2 x:=1
- ightharpoonup Si $\lambda = 0$

$$\min_{\beta} E(\beta) = \sum_{i=1}^{n} (y_i - x_i \beta)^2$$
 (6)

y la solución es

$$\hat{\beta}_{OLS} = \frac{\sum (x_i - \bar{x}) (y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \overline{\sum} x_i y_i(7)$$

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - x_i\beta)^2 + \lambda|\beta|$$
(8)

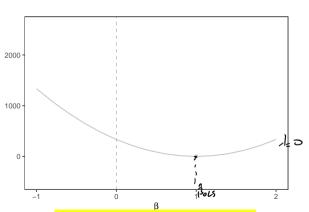
$$E(\beta) = \begin{cases} \frac{2}{2} (y_1 - x_2 \beta)^2 + \lambda \beta \beta > 0 \\ \frac{2}{2} (y_1 - x_2 \beta)^2 - \lambda \beta \beta < 0 \end{cases}$$

Cradiaticos y Afrenciosos pora BED

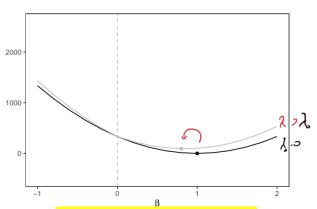
Fina pio y b Sien

Intuición en 1 Dimensión $\hat{\beta} > 0$ per distinfes colores de Co perdi lad λ

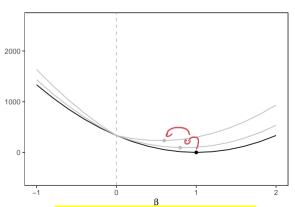
$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - x_i\beta)^2 + \beta$$
(9)



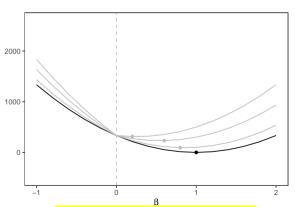
$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - x_i\beta)^2 + \lambda\beta$$
 (10)



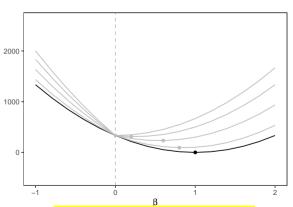
$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - x_i\beta)^2 + \lambda\beta$$
 (11)



$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - x_i\beta)^2 + \lambda\beta$$
 (12)

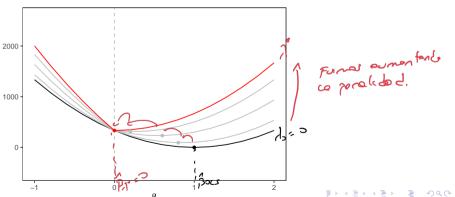


$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - x_i\beta)^2 + \lambda\beta$$
 (13)



 $\hat{\beta} > 0$

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - x_i\beta)^2 + \lambda\beta$$
 (14)

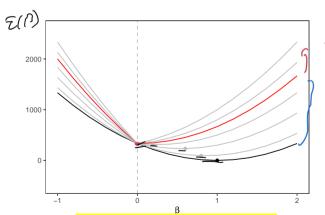


Sarmiento-Barbieri (Uniandes)

Regularización: Lasso

$$\hat{\beta} > 0$$

$$\min_{\beta} E(\beta) = \sum_{i=1}^{n} (y_i - x_i \beta)^2 + \lambda \beta$$
(15)



so to be solved

solucion intrany

Ilustración en R



photo from https://www.dailydot.com/parsec/batman-1966-labels-tumblr-twitter-vine/

$$min_{eta}E(eta) = \sum_{i=1}^{n}(y_i - x_ieta)^2$$

$$\min_{\beta} E(\beta) = \sum_{i=1}^{n} (y_i - x_i \beta)^2 + \lambda |\beta|$$

$$= -2 \sum_{\beta} y_i y_i^{\beta} + 2 \sum_{\beta} x_i^{\beta} y_i^{\beta} + \lambda = 0$$

-> Z (y: - x: B) + 13 (3)=

(16)

Solución analitica

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - x_i\beta)^2 + \lambda |\underline{\beta}|$$
 (16)

la solución analítica es

$$\hat{\beta}_{lasso} = \begin{cases} 0 & \text{si } \lambda \ge \lambda^* \\ \hat{\beta}_{OLS} - \frac{\lambda}{2} & \text{si } \lambda < \underline{\lambda}^* \end{cases}$$
 (17)

Intuición en 2 Dimensiones (Lasso)

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - x_{i1}\beta_1 - x_{i1}\beta_2)^2 \text{ s.a.} (|\beta_1| + |\beta_2|) \le c$$

$$|\beta_i - \delta_i| + |\beta_2| = 0$$
OLS solution
$$|\beta_i| = b_i$$
OLS solution

(18)

Example



photo from https://www.dailydot.com/parsec/batman-1966-labels-tumblr-twitter-vine/

Resumen

 Ridge y Lasso son sesgados, pero las disminuciones en varianza pueden compensar estoy v llevar a un MSE menor

Lasso encoje a cero, Ridge no tanto

13; -0/~losso

- Importante para aplicación:
 - Estandarizar los datos
 - \triangleright Como elegimos λ ?

- > glownt (coret) to home nor lefact.



Resumen

- ► Ridge y Lasso son sesgados, pero las disminuciones en varianza pueden compensar estoy y llevar a un MSE menor
- Lasso encoje a cero, Ridge no tanto
- ► Importante para aplicación:
 - Estandarizar los datos
 - ightharpoonup Como elegimos λ ? ightharpoonup Validación cruzada

Agenda

- 1 Recap
- 2 Lasso
- 3 Familia de regresiones penalizadas
- 4 k > n
- 5 Elastic Net

Family of penalized regressions

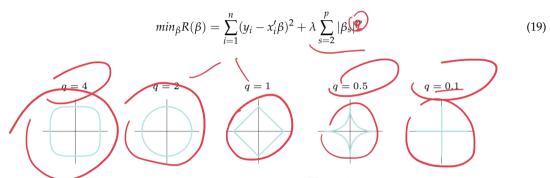


FIGURE 3.12. Contours of constant value of $\sum_{j} |\beta_{j}|^{q}$ for given values of q.

Agenda

- 1 Recap
- 2 Lasso
- 3 Familia de regresiones penalizadas
- 4 k > n
- 5 Elastic Net

More predictors than observations (k > n)

- ► Objective 1: Accuracy
 - ightharpoonup Minimize prediction error (in one step) ightharpoonup Ridge, Lasso
- ► Objective 2: Dimensionality
 - ► Reduce the predictor space → Lasso's free lunch

Ly Pudge no encoje a 0 no nor seleccións.

- ▶ What happens when we have more predictors than observations (k > n)?
 - ► OLS fails _
 - Ridge augments data
 - ▶ and Lasso?

Lasso when k > n

$$V(\vec{\beta}) = \sqrt{(\kappa)(\sqrt{2})}$$

- Lasso works fine in this case
- ▶ However, there are some issues to keep in mind
 - ightharpoonup When k > n chooses at most n variables
 - ▶ When we have a group of highly correlated variables,
 - Lasso chooses only one. Makes it unstable for prediction. (Doesn't happen to Ridge)
 - ▶ Ridge shrinks the coefficients of correlated variables toward each other. This makes Ridge "work" better than Lasso. "Work" in terms of prediction error

Agenda

- 1 Recap
- 2 Lasso
- 3 Familia de regresiones penalizadas
- 4 k > n
- 5 Elastic Net

Elastic net

$$min_{\beta}EN(\beta) = \sum_{i=1}^{n} (y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j)^2 + \lambda \left(\alpha \sum_{j=1}^{p} |\beta_j| + \frac{(1-\alpha)}{2} \sum_{j=1}^{p} (\beta_j)^2\right)$$
(20)

- ightharpoonup Si $\alpha = 1$ Lasso
- ► Si $\alpha = 0$ Ridge

Elastic Net

- ► Elastic net: happy medium.
 - ► Good job at prediction and selecting variables

$$min_{\beta}EN(\beta) = \sum_{i=1}^{n} (y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j)^2 + (\lambda) \left(\alpha \sum_{j=1}^{p} |\beta_j| + \frac{(1-\alpha)}{2} \sum_{j=1}^{p} (\beta_j)^2 \right)$$
(23)

- Mixes Ridge and Lasso
- ► Lasso selects predictors
- ► Strict convexity part of the penalty (ridge) solves the grouping instability problem
- ▶ How to choose (λ, α) ? → Bidimensional Crossvalidation
- ► Recomended lecture: Zou, H. & Hastie, T. (2005)
- ▶ H.W.: $\beta_{OLS} > 0$ one predictor standarized

$$\hat{\beta}_{EN} = \frac{\left(\hat{\beta}_{OLS} - \frac{\lambda_1}{2}\right)_+}{1 + \lambda_-} \tag{22}$$



Sarmiento-Barbieri (Uniandes)

Example



photo from https://www.dailydot.com/parsec/batman-1966-labels-tumblr-twitter-vine/