

# 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 econométricas estándar no están optimizadas para la predicción porque se enfocan en la insesgadez.
- ▶ OLS por ejemplo es el mejor estimador lineal *insesgado*
- ▶ OLS minimiza el error “*dentro de muestra*”, eligiendo  $\beta$  de forma tal que

$$\min_{\beta} E(\beta) = \sum_{i=1}^n (y_i - \beta_0 - x_{i1}\beta_1 - \cdots - x_{ip}\beta_p)^2 \quad (1)$$

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

# 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

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

- ▶ 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

$$\min_{\beta} E(\beta) = \sum_{i=1}^n (y_i - \beta_0 - x_{i1}\beta_1 - \cdots - x_{ip}\beta_p)^2 + \lambda \sum_{j=1}^p (\beta_j)^2 \quad (3)$$

# Agenda

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

# Lasso

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

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



# Lasso

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

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

- “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
- $L(\beta)$  es no differentiable

# Lasso Intuición en 1 Dimension

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

- Un solo predictor, un solo coeficiente
- Si  $\lambda = 0$

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

- y la solución es

$$\hat{\beta}_{OLS} \quad (7)$$

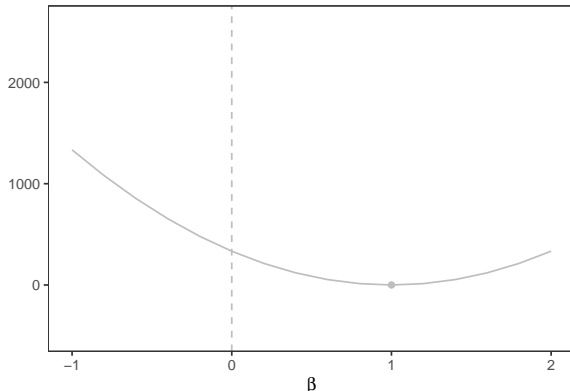
# Intuición en 1 Dimension

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

# Intuición en 1 Dimensión

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

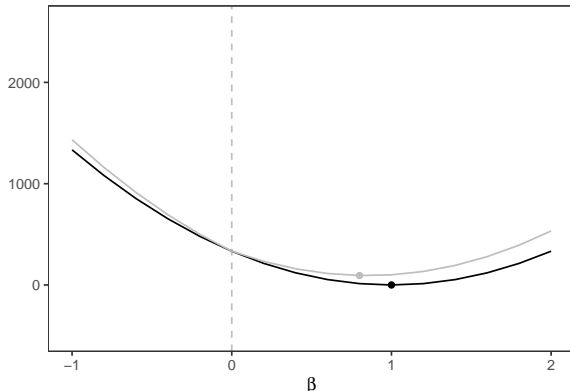
$$\min_{\beta} E(\beta) = \sum_{i=1}^n (y_i - x_i \beta)^2 + \lambda \beta \quad (9)$$



# Intuición en 1 Dimensión

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

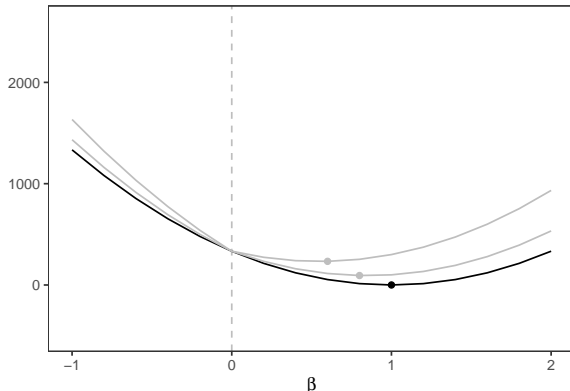
$$\min_{\beta} E(\beta) = \sum_{i=1}^n (y_i - x_i \beta)^2 + \lambda \beta \quad (10)$$



# Intuición en 1 Dimensión

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

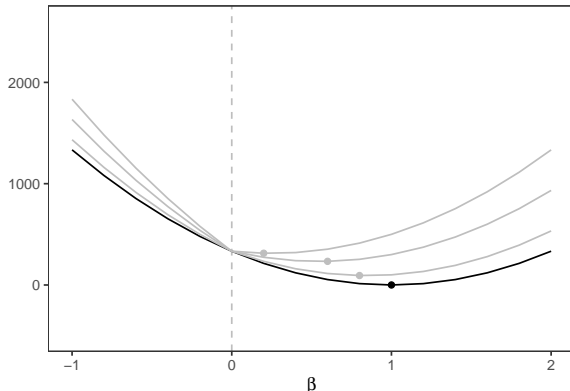
$$\min_{\beta} E(\beta) = \sum_{i=1}^n (y_i - x_i \beta)^2 + \lambda \beta \quad (11)$$



# Intuición en 1 Dimensión

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

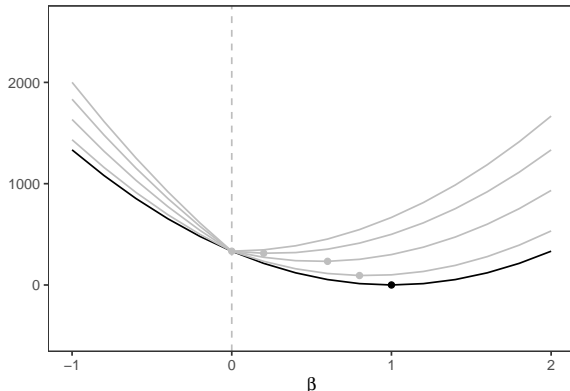
$$\min_{\beta} E(\beta) = \sum_{i=1}^n (y_i - x_i \beta)^2 + \lambda \beta \quad (12)$$



# Intuición en 1 Dimensión

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

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

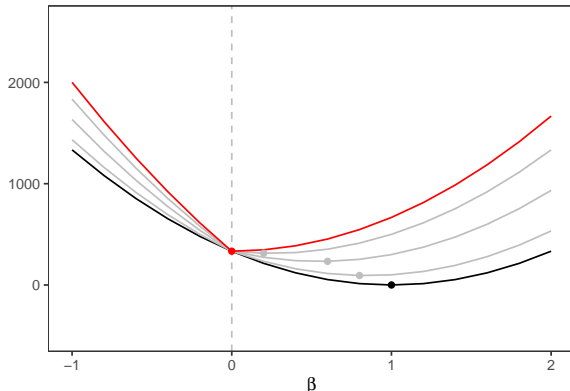




# Intuición en 1 Dimensión

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

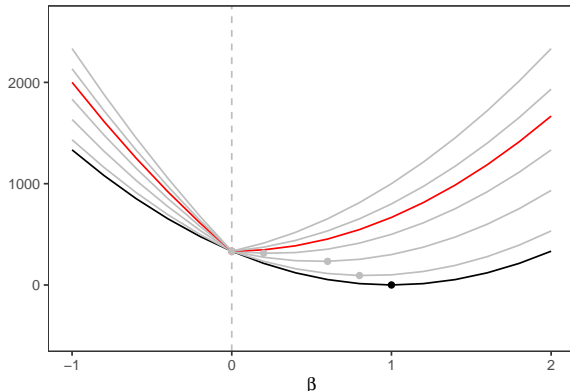
$$\min_{\beta} E(\beta) = \sum_{i=1}^n (y_i - x_i \beta)^2 + \lambda \beta \quad (14)$$



# Intuición en 1 Dimensión

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

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



# Ilustración en R



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

# Intuición en 1 Dimension

Solución analítica

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

# Intuición en 1 Dimension

## Solución analítica

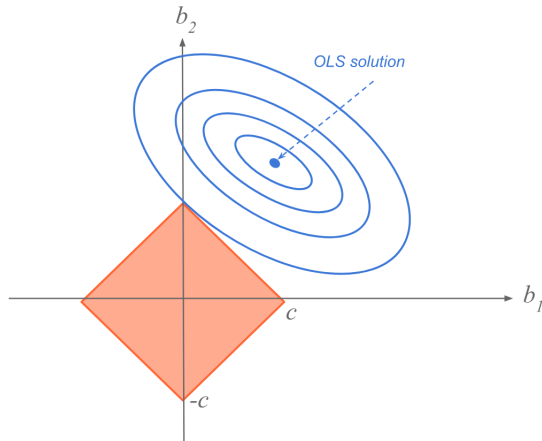
$$\min_{\beta} E(\beta) = \sum_{i=1}^n (y_i - x_i \beta)^2 + \lambda |\beta| \quad (16)$$

► la solución analítica es

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

# Intuición en 2 Dimensiones (Lasso)

$$\min_{\beta} E(\beta) = \sum_{i=1}^n (y_i - x_{i1}\beta_1 - x_{i2}\beta_2)^2 \text{ s.a. } (|\beta_1| + |\beta_2|) \leq c \quad (18)$$



# 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 esto y llevar a un MSE menor
- ▶ Lasso encoje a cero, Ridge no tanto
- ▶ Importante para aplicación:
  - ▶ Estandarizar los datos
  - ▶ Como elegimos  $\lambda$ ?



# Resumen

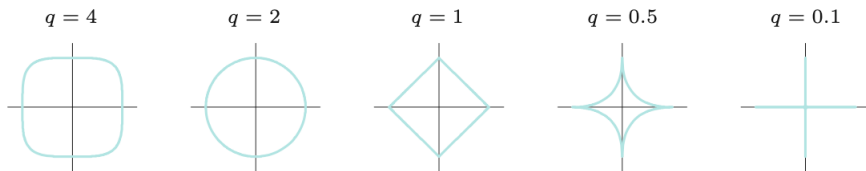
- ▶ Ridge y Lasso son sesgados, pero las disminuciones en varianza pueden compensar esto y llevar a un MSE menor
- ▶ Lasso encoje a cero, Ridge no tanto
- ▶ Importante para aplicación:
  - ▶ Estandarizar los datos
  - ▶ Como elegimos  $\lambda$ ? → Validación cruzada

# Agenda

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

# Family of penalized regressions

$$\min_{\beta} R(\beta) = \sum_{i=1}^n (y_i - x_i' \beta)^2 + \lambda \sum_{s=2}^p |\beta_s|^q \quad (19)$$



**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
  - ▶ Minimize prediction error (in one step)  $\rightarrow$  Ridge, Lasso
- ▶ Objective 2: Dimensionality
  - ▶ Reduce the predictor space  $\rightarrow$  Lasso's free lunch
- ▶ What happens when we have more predictors than observations ( $k > n$ )?
  - ▶ OLS fails
  - ▶ Ridge augments data
  - ▶ and Lasso?

# Lasso when $k > n$

- ▶ Lasso works fine in this case
- ▶ However, there are some issues to keep in mind
  - ▶ 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

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

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



# Elastic Net

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

$$\min_{\beta} EN(\beta) = \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \left( \alpha \sum_{j=1}^p |\beta_j| + \frac{(1-\alpha)}{2} \sum_{j=1}^p (\beta_j)^2 \right) \quad (21)$$

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

$$\hat{\beta}_{EN} = \frac{\left( \hat{\beta}_{OLS} - \frac{\lambda_1}{2} \right)_+}{1 + \lambda_2} \quad (22)$$

# Example



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