Selección de Modelos y Regularización Big Data y Machine Learning para Economía Aplicada

Ignacio Sarmiento-Barbieri

Universidad de los Andes

- 1 Ridge
 - Trade-off Sesgo-Varianza
 - Escala de las variables /
 - More predictors than observations /
 - Selección de $\widehat{\lambda}$
- 2 Lasso
- 3 Recap

- 1 Ridge
 - Trade-off Sesgo-Varianza
 - Escala de las variables
 - More predictors than observations
 - \circ Selección de λ
- 2 Lasso
- 3 Recap

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 > (a evis northo)
- ightharpoonup OLS minimiza el error "dentro de muestra", eligiendo β de forma tal que

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - \widehat{\beta_0 - x_{i1}\beta_1 - \dots - x_{lp}\beta_p})^2$$
 (1)

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



- 1 Ridge
 - Trade-off Sesgo-Varianza
 - Escala de las variables
 - More predictors than observations
 - Selección de λ
- 2 Lasso
- 3 Recap

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

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

$$(3)$$

Formalmente

- Las Xs estan estandarizadas (x_i con media 0 ($\bar{x}=0$) y varianza 1 ($\sum x_i^2=1$))
- Regresión: $y = \beta x + u$
- ► OLS

$$\hat{\beta}_{ols} = \sum x_i y_i$$

▶ Ridge

$$\hat{\beta}_{ridge} = \frac{\sum x_i y_i}{(1+\lambda)} = \frac{\hat{\beta}_{ols}}{(1+\lambda)} \qquad \begin{array}{c} -\gamma & \bigcirc \\ \gamma & \rightarrow \infty \end{array}$$

4□ ▶ 4回 ▶ 4 豆 ▶ 4 豆 ・ 夕 ♀ ○

Formalmente

- 290 #
- -> No esso rada mesora ► En regresión multiple (X es una matriz $n \times \underline{k}$)
- Regresión: $y = X\beta + u$
- ► OLS

- (XKX) XAKK) KILL H.W. Como lo Ce Ji eston Lorizono, Cos X's
- $\hat{\beta}_{ols} = (X'X)^{-1}X'y$
- Ridge

 $\hat{\beta}_{ridge} = (X'X + \lambda I)^{-1}X'y$

Ridge vs OLS



- ► Ridge es sesgado $E(\hat{\beta}_{ridge}) \neq \beta$
- ▶ Pero la varianza es menor que la de OLS
- lacktriangleq Para ciertos valores del parámetro λ $MSE_{OLS} > MSE_{ridge}$
- ► Mostremos esto para el caso de 1 variable

Ridge vs OLS

222

 $\sum x_i^2 = 1$

- - ightharpoonup Varianza $V(\hat{\beta}_{ols}) = \sqrt{2}$
 - $MSE(\hat{\beta}_{ols}) = \mathcal{B}_{oos}^{2} + \mathcal{V} = 6^{2}$

Sesgo $E(\hat{\beta}_{ols}) - \beta = \bigcirc \square$ In lest $g \circ \omega$

y = x3+4

7 es 20 170 do (1+2)

► Ridge:

- Sesgo $E(\hat{\beta}_{ridge}) \beta = \begin{bmatrix} \beta \\ \beta \end{bmatrix} + \begin{bmatrix} \beta \\ \beta \end{bmatrix}$ Varianza $V(\hat{\beta}_{ridge}) =$

$$\blacktriangleright MSE(\hat{\beta}_{ridge}) = \begin{pmatrix} \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix} \end{pmatrix}$$

1 varas6

- V (Bridge) = D

Ridge vs OLS

$$MSE(\hat{\beta}_{ols}) - MSE(\hat{\beta}_{ridge}) =$$

$$(4)$$

- 1 Ridge
 - Trade-off Sesgo-Varianza
 - Escala de las variables
 - More predictors than observations
 - Selección de λ
- 2 Lasso
- 3 Recap

- ▶ La escala de las variables importa en Ridge, mientras que en OLS no.
- ► Tiene consecuencias
 - ► En la solución $(\hat{\beta})$
 - ► En la predicción (ŷ)

Ridge no es invariante a las escala

- ► Supongamos z = c * x
- ► Vamos a mostrar que $\hat{y}_i^z = \hat{y}_i^x$
- ▶ Partamos del modelo

$$y_i = \beta_0^z + \beta_1^z z_i + u \tag{5}$$

$$\hat{\beta}_{1}^{z} = \frac{\sum (z_{i} - \bar{z})(y_{i} - \bar{y})}{\sum (z_{i} - \bar{z})^{2}}$$
(6)

Ridge no es invariante a las escala

Continuando

$$\hat{\beta}_{1}^{z} = \frac{\sum (z_{i} - \bar{z})(y_{i} - \bar{y})}{\sum (z_{i} - \bar{z})^{2}}$$
(7)

ightharpoonup Pero z = c * x

$$\hat{\beta}_{1}^{z} = \frac{\sum (cx_{i} - c\bar{x})(y_{i} - \bar{y})}{\sum (cx_{i} - c\bar{x})^{2}}$$

$$= \frac{1}{c} \frac{\sum (x_{i} - \bar{x})(y_{i} - \bar{y})}{\sum (x_{i} - \bar{x})^{2}}$$

$$= \frac{1}{c} \hat{\beta}_{1}^{x}$$

► En Ridge?

(8)

(9)

Ridge no es invariante a las escala

▶ Para un $\lambda \ge 0$ dado, el problema de optimización

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - \beta_0^z - \beta_1^z z_i)^2 + \lambda(\beta_1^z)^2$$
 (11)

▶ Demo: baticomputer, math: Homework



Ridge no es invariante a las escala

► En la predicción

$$\hat{\beta}_1^z z_i = \hat{\beta}_1^z c x_i \tag{12}$$

$$=\frac{1}{c}\hat{\beta}_1^x cx_i \tag{13}$$

$$=\hat{\beta}_1^x x_i \tag{14}$$

Ridge no es invariante a las escala

ightharpoonup En términos generales, si Z = cX

$$\begin{split} \hat{\beta}_{OLS}^{Z} &= (Z'Z)^{-1}Z'y \\ &= ((cX)'(cX))^{-1}(cX)'y \\ &= \frac{c}{c^2}(X'X)^{-1}X'y \\ &= \frac{1}{c}(X'X)^{-1}X'y \\ &= \frac{1}{c}\hat{\beta}_{OLS}^{X} \end{split}$$

Ridge no es invariante a las escala

Entonces

$$\hat{\beta}_{OLS}^{Z} Z = \frac{1}{c} \hat{\beta}_{OLS}^{X} cX$$
$$= \hat{\beta}_{OLS}^{X} X$$

Con Ridge esto no funciona

$$\hat{\beta}_{Ridge}^{Z}Z \neq \hat{\beta}_{Ridge}^{X}X$$

► Es importante estandarizar las variables (la mayoría de los softwares lo hace automáticamente)



- 1 Ridge
 - Trade-off Sesgo-Varianza
 - Escala de las variables
 - More predictors than observations
 - Selección de λ
- 2 Lasso
- 3 Recap

More predictors than observations (k > n)

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

OLS when k > n

- ► Rank? Max number of rows or columns that are linearly independent
 - ▶ Implies $rank(X_{k \times n}) \le min(k, n)$
- ▶ MCO we need $rank(X_{k \times n}) = k \implies k \le n$
- ▶ If $rank(X_{k \times n}) = k$ then rank(X'X) = k
- ▶ If k > n, then $rank(X'X) \le n < k$ then (X'X) cannot be inverted
- ▶ Ridge and Lasso work when $k \ge n$

Ridge when k > n

$$min_{\beta}R(\beta) = \sum_{i=1}^{n} (y - x\beta)^2 + \lambda(\beta)^2$$
(15)

- ▶ Solution → data augmentation
- ► Intuition: Ridge "adds" *k* additional points.
- ▶ Allows us to "deal" with $k \ge n$

Ridge when k > n

$$min_{\beta}R(\beta) = \sum_{i=1}^{n} (y_i - x_i'\beta)^2 + \lambda(\beta_s)^2$$
 (16)

Ridge when k > n

$$min_{\beta}R(\beta) = \sum_{i=1}^{n} (y_i - \sum_{j=1}^{k} x'_{ij}\beta_j)^2 + \lambda (\sum_{j=1}^{k} \beta_j)^2$$
(17)

- 1 Ridge
 - Trade-off Sesgo-Varianza
 - Escala de las variables
 - More predictors than observations
 - ullet Selección de λ
- 2 Lasso
- 3 Recap

Selección de λ

- Asegurar cero sesgo dentro de muestra crea problemas fuera de muestra: trade-off Sesgo-Varianza
- Ridge hace este trade-off de forma empírica.

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - \beta_0 - x_{i1}\beta_1 - \dots - x_{ip}\beta_p)^2 + \lambda \sum_{j=1}^{p} R(\beta_j)$$
 (18)

- $ightharpoonup \lambda$ es el precio al que hacemos este trade off
- Como elegimos λ?



Selección de λ

- lacktriangledown λ es un hiper-parámetro y lo elegimos usando validación cruzada
 - ▶ Partimos la muestra de entrenamiento en K Partes: $MUESTRA = M_{fold \, 1} \cup M_{fold \, 2} \cdots \cup M_{fold \, K}$
 - ► Cada conjunto $M_{fold \, K}$ va a jugar el rol de una muestra de evaluación $M_{eval \, k}$.
 - Entonces para cada muestra
 - $ightharpoonup M_{train-1} = M_{train} M_{fold 1}$
 - •
 - $ightharpoonup M_{train-k} = M_{train} M_{fold\,k}$

Selección de λ

- Luego hacemos el siguiente loop
 - Para $i = 0, 0.001, 0.002, \dots, \lambda_{max}$ {
 - Para k = 1, ..., K {
 - Ajustar el modelo $m_{i,k}$ con λ_i en $M_{train-k}$
 - Calcular y guardar el $MSE(m_{i,k})$ usando M_{eval-k}
 - } # fin para k
 - Calcular y guardar $MSE_i = \frac{1}{K}MSE(m_{i,k})$
 - $\}$ # fin para λ
- ► Encontramos el menor MSE_i y usar ese $\lambda_i = \lambda^*$





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

- 1 Ridge
 - Trade-off Sesgo-Varianza
 - Escala de las variables
 - More predictors than observations
 - \circ Selección de λ
- 2 Lasso
- 3 Recap



Lasso

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

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - \beta_0 - x_{i1}\beta_1 - \dots - x_{ip}\beta_p)^2 + \lambda \sum_{i=1}^{p} |\beta_i|$$
 (19)

Lasso

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

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - \beta_0 - x_{i1}\beta_1 - \dots - x_{ip}\beta_p)^2 + \lambda \sum_{j=1}^{p} |\beta_j|$$
 (19)

- 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



- 1 Ridge
 - Trade-off Sesgo-Varianza
 - Escala de las variables
 - More predictors than observations
 - Selección de λ
- 2 Lasso
- 3 Recap

Recap

- ► El objetivo es predecir bien fuera de muestra, donde nos enfrentamos al trade-off Sesgo-Varianza
- ► Propusimos modelos

$$min_{\beta}E(\beta) = \sum_{i=1}^{n} (y_i - \beta_0 - x_{i1}\beta_1 - \dots - x_{ip}\beta_p)^2 + \lambda \sum_{j=1}^{p} R(\beta_j)$$
 (20)

- ▶ 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
 - ► Lasso
 - Elastic Net
- Próxima clase: detalles de Lasso y EN

