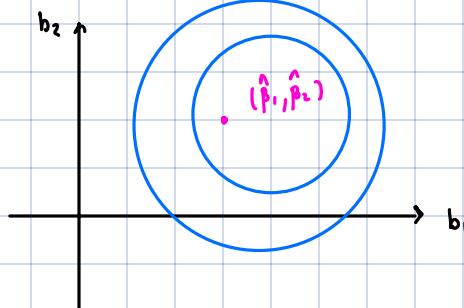


Recordemos nuestro modelo lineal

$$Y_i = X_i' \beta + u_i, \quad E[u_i | X_i] = 0.$$

Tenemos que  $\hat{\beta}_n$  es la solución a

$$\hat{\beta}_n = \underset{b}{\operatorname{argmin}} \sum_{i=1}^n \frac{(Y_i - X_i' b)^2}{2n}$$



Las curvas de nivel se minimizan en  $(\hat{\beta}_1, \hat{\beta}_2)$

### 1. Ridge

La regresión Ridge propone estimar

$$\hat{\beta}_n^R = \underset{b}{\operatorname{argmin}} \sum_{i=1}^n \frac{(Y_i - X_i' b)^2}{2n} + \lambda b^T b / 2,$$

$$\text{Penalizamos } \|b\|_2^2 := \sum_{j=1}^p b_j^2$$

donde  $\lambda$  es un parámetro externo que elegimos (no estimamos directamente). A esto se le conoce como un hiperparámetro.

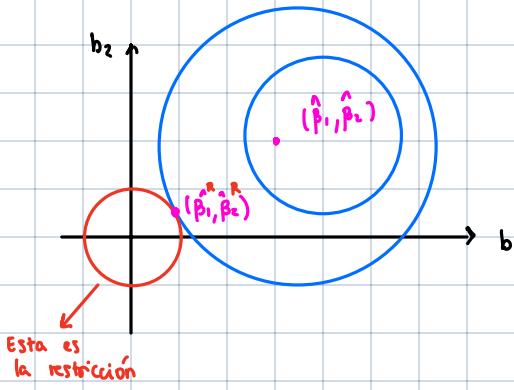
$$\text{Derivando } Q_n(b; \lambda) = \frac{1}{2n} (Y - Xb)^T (Y - Xb) + \frac{\lambda}{2} b^T b$$

$$- X^T (Y - X\hat{\beta}_n^R) + \lambda \hat{\beta}_n^R = 0$$

$$\Rightarrow X^T X \hat{\beta}_n^R + \lambda \hat{\beta}_n^R = X^T Y$$

$$\Rightarrow \hat{\beta}_n^R = (X^T X + \lambda I_p)^{-1} X^T Y$$

Incluso cuando  $X^T X$  es cercano a ser singular, al sumar este último término hacen que una solución siempre exista.



\* La regresión Ridge acerca los coeficientes hacia cero (shrinkage)

Sin embargo, noten que hemos introducido sesgo incluso en muestras finitas.  
Supongamos  $E[U_i | X_i] = 0$ .

$$\hat{\beta}_n^R = (X^T X + \lambda I_n)^{-1} X^T X \beta + (X^T X + \lambda I_n)^{-1} X^T U$$

$$\Rightarrow E[\hat{\beta}_n^R | X] = (X^T X + \lambda I_n)^{-1} X^T X \beta \neq \beta.$$

Es decir, Ridge introduce sesgo a costa de reducir la varianza de las predicciones. Noten como estabiliza la varianza en el caso  $k=1$  y homocedástico :

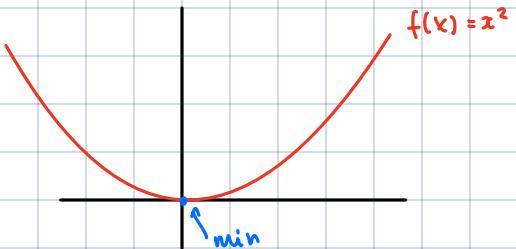
$$\hat{\beta}_{1n} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n (x_i^2 + \lambda)}, \quad \text{Var}(\hat{\beta}_1 | X) = \frac{\sigma^2}{\sum_{i=1}^n (x_i^2 + \lambda)}$$

Este denominador está más lejos de 0.

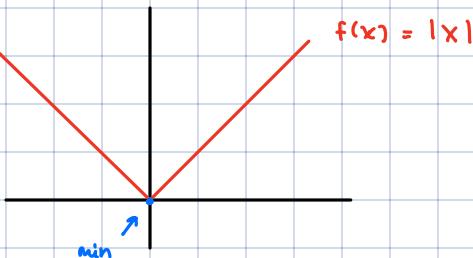
Finalmente, ya que Ridge **NO** impone una solución de esquina l.i.e. que algún  $\hat{\beta}_j$  sea 0, no sirve para hacer selección de regresores. Por ello, estudiaremos un siguiente estimador.

## 2. Minimización convexa

Existen casos donde queremos minimizar una función objetivo que es convexa. Por tanto, definitivamente podemos encontrar un óptimo. Sin embargo, no siempre son objetos diferenciables.



Se halla como  $\partial f(x) = 2x = 0$



¿Cómo se halla?

### 2.1. Subgradiente

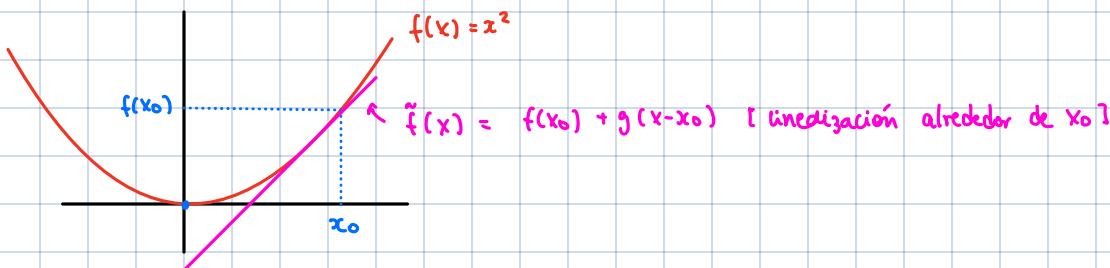
Def.:  $g$  es subgradiente de  $f$  en  $x_0$  si  $\forall x \in \mathbb{R}^n$  tenemos

$$f(x) - f(x_0) \geq g^T(x - x_0)$$

Para ganar intuición consideremos  $n=1$ . La definición implica que

$$\begin{aligned} f(x) &\geq f(x_0) - g^T(x_0) + g^T x \\ &= f(x_0) + g^T(x - x_0) \end{aligned}$$

Ejemplo:



Es decir, la subgradiente es una de las posibles pendientes de las líneas que pasan por el punto  $(x_0, f(x_0))$ , y que siempre están por debajo de  $f(x)$ . Vemos que en el caso de  $f(x) = x^2$  solo hay una posible línea.

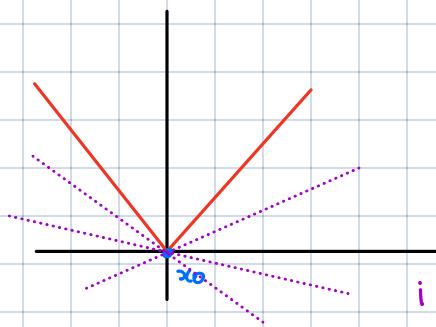
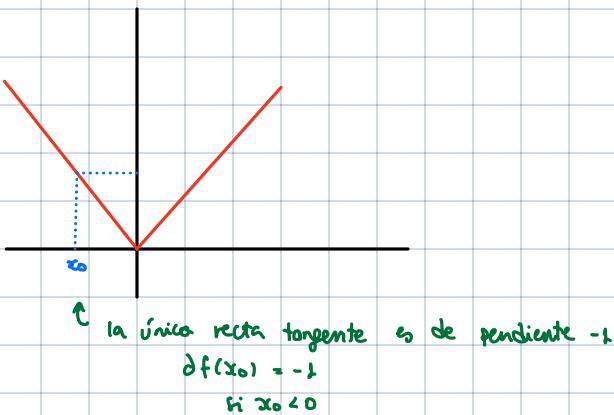
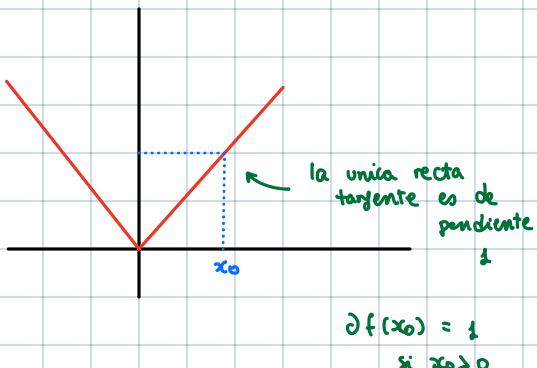
Def.: El set de todos los subgradientes en  $x_0$  se conoce como subdiferencial en  $x_0$ .

$$\partial f(x_0) = \{ g \in \mathbb{R}^n : f(x) - f(x_0) \geq g^T(x - x_0) \quad \forall x \in \mathbb{R}^n \}$$

Decimos que una función es diferenciable en  $x_0$  si su subdiferencial en  $x_0$  **solo incluye un único elemento**.

$$\partial f(x_0) = \{g\} \Leftrightarrow g = \nabla f(x).$$

Veamos ahora el valor absoluto:



i Hay varias posibles tangentes!

$$\partial f(x_0) = [-1, 1] \text{ si } x_0 = 0$$

Es decir, cualquier pendiente entre -1 y 1 es un subgradiente.

Por ello, decimos que:

$$\partial|x| = \begin{cases} -1 & , x < 0 \\ [-1, 1] & , x = 0 \\ 1 & , x > 0 . \end{cases}$$

Proposición.- Sea  $M$  el set de puntos mínimos de una función convexa  $f$ :

$$M = \{x^* \in \mathbb{R}^k : f(x^*) \leq f(x) \quad \forall x \in \mathbb{R}^k\}$$

Entonces

$$x^* \in M \Leftrightarrow 0 \in \partial f(x^*).$$

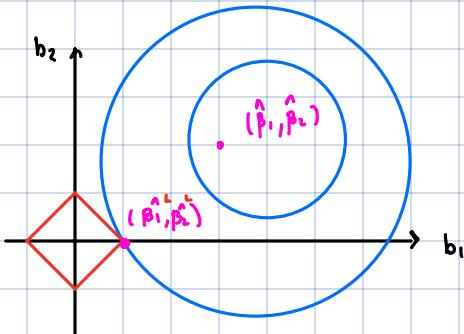
Por ello, podemos decir que  $x_0 = 0$  minimiza  $f(x) = \|x\|$ :

$$0 \in \partial \|x\| \text{ cuando } x = x_0 \Rightarrow x^* \in M.$$

### 3. LASSO

La regresión Ridge propone estimar

$$\hat{\beta}_n^L = \underset{b}{\operatorname{argmin}} \frac{\sum_{i=1}^n (y_i - x_i' b)^2}{2n} + \lambda \sum_{j=1}^k |b_j| = \frac{\|y - x_b\|_2^2}{2n} + \lambda \|b\|_1$$



Está diseñado para hallar soluciones de eje en donde algún  $\hat{\beta}_j$  será igual a 0. Es decir, selecciona regresores.

Importantemente, este es un problema de minimización convexa.

#### 3.1. Solución Analítica (caso especial)

Consideremos

$$\sum_{i=1}^n \frac{x_{ij} x_{il}}{n} = \begin{cases} 0 & j \neq l \\ 1 & j = l \end{cases} \quad \left. \begin{array}{l} \text{regresores ortogonales} \\ \text{segundo momento igual a 1.} \end{array} \right\}$$

En ese caso

$$\begin{aligned} \frac{x' x}{n} &= I_n \\ \Rightarrow \hat{\beta}_{OLS} &= \frac{x' y}{n}, \end{aligned}$$

por lo que

$$\hat{\beta}_j^{OLS} = \sum_{i=1}^n \frac{x_i y_i}{n}.$$

Vamos a definir  $(x)^+ = \max\{x, 0\}$

$$\operatorname{sign}(x) = \begin{cases} -1 & , x < 0 \\ 0 & , x = 0 \\ 1 & , x > 0 \end{cases}$$

Proposición : Suponer  $X'X/n = I_K$ . Entonces el estimador LASSO resuelve la minimización de

$$Q_{n,\lambda}(b) = \left\{ \frac{1}{2n} \|y - Xb\|_2^2 + \lambda \|b\|_1 \right\}$$

y satisface

$$\hat{\beta}_j = \text{sign}(\hat{\beta}_j^{\text{OLS}}) (|\hat{\beta}_j^{\text{OLS}}| - \lambda)^+$$

Achica los coeficientes en una magnitud  $\lambda$ .

Si  $\lambda > |\hat{\beta}_j^{\text{OLS}}|$  lo reemplaza  $\hat{\beta}_j^{\text{OLS}}$  por  $\hat{\beta}_j = 0$ .

Shrinkage + Selection

proof :

Definimos

$$\underbrace{v}_{\text{vector}} + c \underbrace{S}_{\text{set de vectores}} = \{ v + cg : g \in S \}.$$

Usando esa notación definimos

$$\begin{aligned} \partial Q_{n,\lambda}(b) &= -\frac{X'(y - Xb)}{n} + \lambda \partial \|b\|_1 \\ &= -\frac{X'y}{n} + \frac{X'X}{n} b + \lambda \partial \|b\|_1 \\ &= -\hat{\beta}^{\text{OLS}} + I_K b + \lambda \partial \|b\|_1 \\ &= -(\hat{\beta}^{\text{OLS}} - b) + \lambda \partial \|b\|_1 \\ &= -\begin{pmatrix} \hat{\beta}_1^{\text{OLS}} - b_1 & -\lambda \partial |b_1| \\ \vdots & \\ \hat{\beta}_K^{\text{OLS}} - b_K & -\lambda \partial |b_K| \end{pmatrix}. \end{aligned}$$

Noten que podemos resolver elemento por elemento. En concreto, el estimador LASSO satisface

$$0 \in \hat{\beta}_j^{\text{OLS}} - \hat{\beta}_j - \lambda \partial |\hat{\beta}_j|$$

Entonces  $\hat{\beta}_j = 0$  si y solo si

$$0 \in \hat{\beta}_j^{\text{OLS}} - 0 - \lambda \partial |0| = \hat{\beta}_j^{\text{OLS}} - \lambda [-1, 1] = [\hat{\beta}_j^{\text{OLS}} - \lambda, \hat{\beta}_j^{\text{OLS}} + \lambda]$$

Esto es equivalente a decir

$$\hat{\beta}^{OLS} - \lambda \leq 0 \leq \hat{\beta}^{OLS} + \lambda$$

$$\Rightarrow -\lambda \leq \hat{\beta}^{OLS} \leq \lambda$$

$$\Rightarrow |\hat{\beta}^{OLS}| < \lambda$$

$$\Rightarrow (|\hat{\beta}^{OLS}| - \lambda)^+ = 0.$$

El estimador LASSO  $\hat{\beta}_j$  sera  $\neq 0$  si y solo si  $|\hat{\beta}^{OLS}| > \lambda$ , que ocurre si

$$\hat{\beta}^{OLS} > \lambda \quad \text{o} \quad \hat{\beta}^{OLS} < -\lambda.$$

En esos casos tenemos:

$$\underbrace{\hat{\beta}_j > \lambda}_{\hat{\beta}_j > 0} \Rightarrow \hat{\beta}_j^{OLS} - \hat{\beta}_j - \lambda \operatorname{d}|\hat{\beta}_j| = \hat{\beta}_j^{OLS} - \hat{\beta}_j - \lambda_{\times 1} = 0$$

$$\Rightarrow \hat{\beta}_j = \hat{\beta}_j^{OLS} - \lambda = + (|\hat{\beta}_j^{OLS}| - \lambda)^+$$

$$\underbrace{\hat{\beta}_j < -\lambda}_{\hat{\beta}_j < 0} \Rightarrow \hat{\beta}_j^{OLS} - \hat{\beta}_j - \lambda \operatorname{d}|\hat{\beta}_j| = \hat{\beta}_j^{OLS} - \hat{\beta}_j - \lambda_{\times (-1)} = 0$$

$$\Rightarrow \hat{\beta}_j = \hat{\beta}_j^{OLS} + \lambda = - (|\hat{\beta}_j^{OLS}| - \lambda)^+$$

Ahora, recorden que  $\hat{\beta}_j^{OLS} = \beta_j + O_p(n^{-1/2})$ , que para  $\beta_j=0$  implica  $\hat{\beta}_j^{OLS} = O_p(n^{-1/2})$ .

Para que LASSO seleccione correctamente,  $\lambda_n$  debe dominar ese componente de ruido.

Sin embargo  $\lambda_n$  tampoco puede ser muy grande o enviará todo a 0. Esto empezará a ocurrir si  $\lambda > \min_{j \neq 0} |\beta_j|$ .

Supongamos que  $\min_j |\beta_j| \geq \delta > 0$ ,  $\lambda_n = \sqrt{\frac{\log n}{n}}$ .

$$\begin{aligned} \Pr(\hat{\beta}_j = 0 \mid \beta_j = 0) &= \Pr(|\hat{\beta}_j^{OLS}| < \lambda_n \mid \beta_j = 0) \\ &= \Pr(O_p(1/\sqrt{n}) < \sqrt{\frac{\log n}{n}}) \rightarrow 1. \end{aligned}$$

$$\begin{aligned} \Pr(\hat{\beta}_j \neq 0 \mid \beta_j \neq 0) &= \Pr(|\hat{\beta}_j^{OLS}| > \lambda_n \mid \beta_j \neq 0) \\ &= \Pr(|\beta_j + O_p(n^{-1/2})| > \sqrt{\frac{\log n}{n}} \mid \beta_j \neq 0) \\ &= \Pr(\delta + O_p(n^{-1/2}) > \sqrt{\frac{\log n}{n}}) \rightarrow 1. \end{aligned}$$

### 3.2. Caso General

No existe una fórmula cerrada de la solución, pero la vamos a caracterizar. En concreto LASSO resuelve

$$-\frac{1}{n} \sum_{i=1}^n x_i (y_i - x_i' \hat{\beta}) + \lambda_n \hat{g} = 0 ,$$

donde  $\hat{g} \in \|\hat{\beta}\|_1$  :  $\hat{g}_j = \begin{cases} \text{sign}(\hat{\beta}_j) & \hat{\beta}_j \neq 0 \\ \in [-1, 1] & \hat{\beta}_j = 0 \end{cases}$

Denotamos  $A_n^*$  como los regresores seleccionados por Lasso

$$A_n^* = \{j : \hat{\beta}_j \neq 0\}$$

Podemos escribir la condición de primer orden así

$$\frac{1}{n} \sum_{i=1}^n x_{i,A_n^*} (y_i - x_{i,A_n^*}' \hat{\beta}_{A_n^*}) - \lambda_n \underbrace{\text{sign}(\hat{\beta}_{A_n^*})}_{\text{sign}(.) \text{ de cada uno de sus elementos}} = 0$$

Tenemos que Lasso excluye regresores ( $j \in A_n^c$  o  $\hat{\beta}_j = 0$ ) cuando

$$0 = -\frac{1}{n} \sum_{i=1}^n x_{ij} (y_i - x_{i,A_n^*}' \hat{\beta}_{A_n^*}) + \lambda_n \cdot \hat{g}_j , \text{ donde } |\hat{g}_j| \leq 1$$

o equivalentemente

$$\left| \frac{1}{n} \sum_{i=1}^n x_{ij} (y_i - x_{i,A_n^*}' \hat{\beta}_{A_n^*}) \right| \leq \lambda_n$$

Denotamos  $A_0 = \{j : \beta_j \neq 0\}$  como el set de regresores relevantes, tal que el verdadero modelo es

$$y_i = x_{i,A_0}' \beta_{A_0} + u_i$$

Fíjense que  $\text{sign}(\hat{\beta}) = \text{sign}(\beta)$  implica que seleccionamos los regresores correctos  $\hat{\beta}_j = 0 \Rightarrow \hat{\beta}_j = 0$  (ocurre si  $j \notin A_0^c$ )  $\Rightarrow A_n^* = A_0$   
 $\hat{\beta}_j \neq 0 \Rightarrow \hat{\beta}_j \neq 0$  (ocurre si  $j \in A_0$ )

Vamos a demostrar cuando ocurre esto.

Proposición - Supongamos que  $\mathbf{X}_{i,A_0} \mathbf{x}_{i,A_0}'$  es p.d. Entonces  $\text{sign}(\hat{\beta}) = \text{sign}(\beta)$  si y solo si

$$\textcircled{1} \quad \text{sign}(\beta_{A_0}) = \text{sign} \left( \beta_{A_0} + \left( \frac{\sum x_{i,A_0} x_{i,A_0}'}{n} \right)^{-1} \left( \frac{1}{n} \sum x_{i,A_0} u_i - \lambda_n \text{sign}(\beta_{A_0}) \right) \right)$$

y para todo  $j \in A_0^c$

$$\textcircled{2} \quad \left| \frac{1}{n} \sum_{i=1}^n x_{ij} x_{i,A_0}' \left( \frac{\sum x_{i,A_0} x_{i,A_0}'}{n} \right)^{-1} \left( \frac{\sum x_{i,A_0} u_i - \lambda_n \text{sign}(\beta_{A_0})}{n} \right) - \frac{\sum x_{ij} u_i}{n} \right| \leq \lambda_n$$

proof: ( $\Rightarrow$ ) Supongamos  $\text{sign}(\hat{\beta}) = \text{sign}(\beta)$ . Entonces se cumple

$$0 = \frac{1}{n} \sum x_{i,A_0} (y_i - x_{i,A_0}' \hat{\beta}_{A_0}) - \lambda_n \text{sign}(\hat{\beta}_{A_0})$$

$$\Rightarrow 0 = \frac{1}{n} \sum x_{i,A_0} (u_i - x_{i,A_0}' (\hat{\beta}_{A_0} - \beta_{A_0})) - \lambda_n \text{sign}(\hat{\beta}_{A_0})$$

$$\Rightarrow \hat{\beta}_{A_0} = \beta_{A_0} + \left( \frac{\sum x_{i,A_0} x_{i,A_0}'}{n} \right)^{-1} \left( \frac{1}{n} \sum x_{i,A_0} u_i - \lambda_n \text{sign}(\beta_{A_0}) \right)$$

Lo que significa que

$$\begin{aligned} \text{sign}(\beta_{A_0}) &= \text{sign}(\hat{\beta}_{A_0}) \\ &= \text{sign} \left( \beta_{A_0} + \left( \frac{\sum x_{i,A_0} x_{i,A_0}'}{n} \right)^{-1} \left( \frac{1}{n} \sum x_{i,A_0} u_i - \lambda_n \text{sign}(\beta_{A_0}) \right) \right). \end{aligned}$$

(obtenemos  $\textcircled{1}$ )

La ecuación  $\textcircled{2}$  viene de la cond. de primer orden, ya que  $A_0^c = A_0$

excluye a  $j$  si:

$$\left| \frac{1}{n} \sum_{i=1}^n x_{ij} (y_i - x_{i,A_0}' \hat{\beta}_{A_0}) \right| \leq \lambda_n$$

↑ reemplazamos la solución de arriba y nos da  $\textcircled{2}$

( $\Leftarrow$ ) revisarla en las notas, pero sigue técnicas que ya usamos.

Podemos re-escribir la proporción

$$\textcircled{1} \quad \text{sign}(\beta_{A_0}) = \text{sign} \left( \beta_{A_0} + \left( \sum \frac{x_{i,A_0} x_{i,A_0}'}{n} \right)^{-1} \left( \frac{1}{n} \sum x_{i,A_0} u_i - \lambda_n \text{sign}(\beta_{A_0}) \right) \right)$$

para todo  $j \in A_0^c$

$$\textcircled{2} \quad \left| \frac{1}{n} \sum_{i=1}^n x_{ij} x_{i,A_0}' \left( \sum \frac{x_{i,A_0} x_{i,A_0}'}{n} \right)^{-1} \left( \frac{\sum x_{i,A_0} u_i}{n} - \lambda_n \text{sign}(\beta_{A_0}) \right) - \frac{\sum x_{ij} u_i}{n} \right| \leq \lambda_n$$

Como

$$\left| \beta_j + \left[ \left( \sum \frac{x_{i,A_0} x_{i,A_0}'}{n} \right)^{-1} \left( \frac{1}{n} \sum x_{i,A_0} u_i - \lambda_n \text{sign}(\beta_{A_0}) \right) \right]_j \right| > 0, \quad j \in A_0^c$$

$$\left| \frac{1}{n} \sum_{i=1}^n x_{ij} x_{i,A_0}' \left( \sum \frac{x_{i,A_0} x_{i,A_0}'}{n} \right)^{-1} \left( \frac{\sum x_{i,A_0} u_i}{n} - \lambda_n \text{sign}(\beta_{A_0}) \right) - \frac{\sum x_{ij} u_i}{n} \right| \leq \lambda_n, \quad j \in A_0^c$$

Supongamos que

$$\lambda_n = \sqrt{\frac{\log n}{n}}.$$

Entonces, para todo  $j \in A_0$

$$\left| \beta_j + \left[ \left( \sum \frac{x_{i,A_0} x_{i,A_0}'}{n} \right)^{-1} \left( \frac{1}{n} \sum x_{i,A_0} u_i - \lambda_n \text{sign}(\beta_{A_0}) \right) \right]_j \right|$$

$$= \left| \beta_j + O_p(1) \left( O_p\left(\frac{1}{\sqrt{n}}\right) - O\left(\sqrt{\frac{\log n}{n}}\right) \right) \right|$$

$$= \left| \beta_j + o_p(1) \right| > 0 \quad \text{con probabilidad aprox. a 1 siempre que } \min_{j \in A_0} |\beta_j| \geq \delta > 0.$$

Ahora, para todo  $j \in A_0^c$

$$\left| \frac{1}{n} \sum_{i=1}^n x_{ij} x_{i,A_0}' \left( \sum \frac{x_{i,A_0} x_{i,A_0}'}{n} \right)^{-1} \left( \frac{\sum x_{i,A_0} u_i}{n} - \lambda_n \text{sign}(\beta_{A_0}) \right) - \frac{\sum x_{ij} u_i}{n} \right| \leq \lambda_n$$

$\underbrace{O_p(1/\sqrt{n})}_{\text{Op}(1/\sqrt{n})}$        $\underbrace{O_p(1/\sqrt{n})}_{\text{Op}(1/\sqrt{n})}$

$$\Rightarrow \left| O_p(1/\sqrt{n}) + \lambda_n \sum_{i=1}^n x_{ij} x_{i,A_0}' \left( \sum \frac{x_{i,A_0} x_{i,A_0}'}{n} \right)^{-1} \text{sign}(\beta_{A_0}) \right| \leq \lambda_n$$

Fríjante que debe cumplirse que en nuestras grandes

$$\left| E X_{ij} X_{iA_0}' (E X_{iA_0} X_{iA_0}')^{-1} \text{sign}(\beta_{A_0}) \right| \leq 1$$

para  $j \in A_0^c$ . Esta condición se llama **irrepresentabilidad**. Lo que dice es que los regresores irrelevantes no deben estar fuertemente correlacionados con los relevantes. De lo contrario, Lasso mantendrá regresores irrelevantes.