

## 1. Modelo de Regresión Lineal

Estamos proponiendo estimar el siguiente modelo:

$$Y_i = X_i' \beta + u_i$$
$$= \sum_{j=1}^k \beta_j x_{ij} + u_i, \quad (1)$$

$\underbrace{\hspace{10em}}$   
 $k$  potenciales regresores

donde estos regresores son exógenos en el sentido de  $E[X_{ij} u_i] = 0$ ,  $\forall j=1, \dots, k$ .  
Además, por el momento supondremos que  $k$  es fijo y no depende de  $n$ . Noten que esto implica que  $k/n \rightarrow 0$  cuando  $n \rightarrow \infty$ .

Denotamos  $A$  como la lista de regresores relevantes (coeficientes distintos de cero)

$$A = \{j : \beta_j \neq 0\}.$$

Ejemplo:  $A = \{1, 4\}$  dice que solo  $X_{i1}$  y  $X_{i4}$  son relevantes.

denotamos  $A_0$  como el verdadero set de regresores relevantes

$\uparrow$  Verdadero modelo!

$$\Rightarrow Y_i = \sum_{j \in A_0} \beta_j X_{ij} + u_i$$

Nuestra meta es estimar  $A_0$  con data  $\{(Y_i, X_i')' : i=1, \dots, n\}$ . Vamos a estimar el set de regresores relevantes con un proceso de selección.  
Decimos que es un proceso de selección consistente si

$$(2) \quad P(\hat{A}_n = A_0) \rightarrow 1 \quad \text{cuando } n \rightarrow \infty.$$

Asimismo  $\beta_A$  es el subvector de  $\beta$  que solo incluye los coeficientes de  $A$ :

$$\beta_A = (\beta_j : j \in A).$$

$\hookrightarrow |A|$ : # de elementos en  $A$ .

Con el procedimiento de selección  $\hat{A}_n$  que produce  $\hat{\beta}_n$  donde  $\hat{\beta}_{nj} = 0$  para  $j \notin \hat{A}_n$ . Decimos que el proceso es  $\text{oráculo}$  si

$$\sqrt{n}(\hat{\beta}_{A_0} - \beta_{A_0}) \xrightarrow{d} N(0, V(A_0)),$$

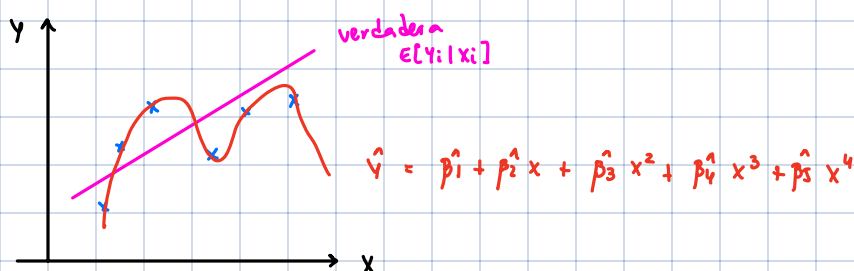
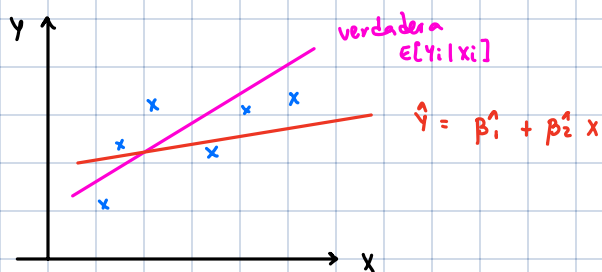
donde  $V(A_0)$  es la varianza asintótica cuando  $A_0$  es conocido.

\* Oráculo significa que es como si en nuestras grandes hubieramos usado el verdadero modelo  $A_0$ .

## 1.1. Penalización

Recorden que  $\hat{\beta}_n$  resuelve  $\min_b \text{SSR} := \sum_{i=1}^n (y_i - x_i' b)^2$  y siempre que aumentamos regresores podemos disminuirlo.

Sin embargo, caemos en el problema de que hacemos **overfitting**



Consideremos que el verdadero proceso generador de datos es:

$$\begin{aligned} y_i &= \sum_{j \in A_0} \beta_j x_{ij} + u_i \\ &= x_{i,A_0}' \beta_{A_0} + u_i \end{aligned}$$

Nosotros proponemos un modelo  $A$  tal que estimamos

$$\begin{aligned} \hat{\beta}_{n,A}(A) &= \left( \sum_{i=1}^n x_{i,A} x_{i,A}' \right)^{-1} \sum_{i=1}^n x_{i,A} y_i, \\ \hat{\beta}_{n,A^c}(A) &= 0 \end{aligned}$$

Ejemplo:  $A = \{1, 3\} \Rightarrow A^c = \{2, 4\} \Rightarrow \begin{aligned} \hat{\beta}_1 &\neq 0, \hat{\beta}_3 \neq 0 \\ \hat{\beta}_2 &= 0, \hat{\beta}_4 = 0 \end{aligned}$   
 $u = 4$

Además, bajo  $A$  definimos

$$\text{SSR}_n(A) = \sum_{i=1}^n (y_i - x_{i,A}' \hat{\beta}_{n,A}(A))^2$$

Una forma razonable de penalizar es usando el número de regresores como un proxy de complejidad

$$BIC_n(A) := SSR_n(A) + \underbrace{|A| \log n}_{\text{penalización}}$$

↑  
 $\approx$  lo que llamamos Bayesian Information Criterion

Entonces

$$\hat{A}_n^{BIC} = \underset{A}{\operatorname{argmin}} BIC_n(A)$$

Teorema :- Si  $EX_i X_i'$  tiene rango completo y  $EU_i^2 X_i X_i' = O(1)$  y es p.d.  
 Entonces

$$P(\hat{A}_n = A_0) \rightarrow 1 \text{ cuando } n \rightarrow \infty.$$

\* Es decir, la selección es consistente.

proof :

Debemos mostrar que para todo  $A \neq A_0$

$$P(\underbrace{SSR_n(A) + |A| \log n}_{BIC_n(A)} > \underbrace{SSR_n(A_0) + |A_0| \log n}_{BIC_n(A_0)}) \rightarrow 1$$

(Esto significa que  $A_0$  en muestras grandes será lo que minimiza  $BIC_n$  con probabilidad alta)

Empezamos mostrando lo siguiente:

$$\begin{aligned} \bullet \frac{SSR_n(A_0)}{n} &= \frac{1}{n} \sum_{i=1}^n (y_i - x_{i,A_0}' \hat{\beta}_{n,A_0}(A_0))^2 \\ &= \frac{1}{n} \sum_{i=1}^n (y_i - x_{i,A_0}' \hat{\beta}_{n,A_0}(A_0) \pm x_{i,A_0}' \beta_{A_0})^2 \\ &= \frac{1}{n} \sum_{i=1}^n (u_i - x_{i,A_0}' (\hat{\beta}_{n,A_0}(A_0) - \beta_{A_0}))^2 \\ &= \underbrace{\frac{1}{n} \sum_{i=1}^n u_i^2}_{EU_i^2 + o_p(1)} - 2 \underbrace{\frac{1}{n} \sum_{i=1}^n u_i x_{i,A_0}' (\hat{\beta}_{n,A_0}(A_0) - \beta_{A_0})}_{\leq E \|u_i X_i\| \cdot \|\hat{\beta}_{n,A_0}(A_0) - \beta_{A_0}\| = 0 \times o_p(1)} + \underbrace{\frac{1}{n} \sum_{i=1}^n [x_{i,A_0}' (\hat{\beta}_{n,A_0}(A_0) - \beta_{A_0})]^2}_{\leq E \|x_i X_i'\| \cdot \|\hat{\beta}_{n,A_0}(A_0) - \beta_{A_0}\|^2 = O(1) \times o_p(1) = o_p(1)} \\ &= EU_i^2 + o_p(1) \end{aligned}$$

$$\bullet \quad \frac{1}{n} \text{SSR}_n(A) = \frac{1}{n} \sum_{i=1}^n (y_i - x_{i,A}' \hat{\beta}_{n,A}(A))^2$$

(s.a.)  
 $(A \cap A_0) \neq A_0$   
 (omitimos regresores relevantes)

$$= \frac{1}{n} \sum_{i=1}^n u_i^2 + \frac{1}{n} \sum_{i=1}^n [x_{i,A}' (\hat{\beta}_{n,A}(A) - \beta_A)]^2 - 2 \frac{1}{n} \sum_{i=1}^n u_i x_{i,A}' (\hat{\beta}_{n,A}(A) - \beta_A)$$

\* Excluir regresores relevantes implica  $\hat{\beta}_n(A) - \beta \rightarrow \delta \neq 0$

$$= E u_i^2 + \delta' E x_i x_i' \delta + o_p(1).$$

$$\bullet \quad \text{SSR}_n(A) = \sum_{i=1}^n u_i^2 + \sum_{i=1}^n [x_{i,A}' (\hat{\beta}_{n,A}(A) - \beta_A)]^2 - 2 \sum_{i=1}^n u_i x_{i,A}' (\hat{\beta}_{n,A}(A) - \beta_A)$$

Nota: que no vale  $1/n$   
 $A_0 \subset A$   
 (contiene regresores irrelevantes)

$$= \sum_{i=1}^n u_i^2 - 2 \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n x_{i,A} u_i \right)' \sqrt{n} (\hat{\beta}_{n,A} - \beta_A)$$

$$+ \sqrt{n} (\hat{\beta}_{n,A} - \beta_A)' \left( \frac{1}{n} \sum_{i=1}^n x_{i,A} x_{i,A}' \right) \sqrt{n} (\hat{\beta}_{n,A} - \beta_A)$$

$$= \underbrace{\sum_{i=1}^n u_i^2}_{\geq 0} + O_p(1) + o_p(1)$$

Entonces

Caso 1:  $P(\text{BIC}_n(A) > \text{BIC}_n(A_0)) = P\left(\frac{\text{SSR}_n(A)}{n} + |A| \frac{\log n}{n} > \frac{\text{SSR}_n(A_0)}{n} + |A_0| \frac{\log n}{n}\right)$

(s.a.)  
 $A \cap A_0 \neq A_0$

$$= P\left(E u_i^2 + \delta' E x_i x_i' \delta + o_p(1) + |A| \frac{\log n}{n} > E u_i^2 + |A_0| \frac{\log n}{n} + o_p(1)\right)$$

$$= P\left(\underbrace{(|A| - |A_0|)}_{\geq 0} \frac{\log n}{n} + o_p(1) + \underbrace{\delta' E x_i x_i' \delta}_{> 0} > 0\right)$$

$\rightarrow 1.$

Caso 2:  $P(\text{BIC}_n(A) > \text{BIC}_n(A_0)) = P(\text{SSR}_n(A) + |A| \log n > \text{SSR}_n(A_0) + |A_0| \log n)$

(s.a.)  
 $A_0 \subset A$

$$= P\left(\sum_{i=1}^n u_i^2 + O_p(1) + o_p(1) + |A| \log n > \sum_{i=1}^n u_i^2 + O_p(1) + o_p(1) + |A_0| \log n\right)$$

$$= P\left(O_p(1) + o_p(1) > \underbrace{(|A_0| - |A|)}_{\leq 0} \log n\right) \rightarrow 1.$$

¿Qué pasa si la penalización fuese algo que no crece con  $n$ ?

$$AIC_n(A) := SSR_n(A) + 2|A|$$

Entonces

Caso 1 :  $P( AIC_n(A) > AIC_n(A_0) ) = P( \underbrace{SSR_n(A)}_{s.a.} + |A| \frac{2}{n} > \underbrace{SSR_n(A_0)}_{s.a.} + |A_0| \frac{2}{n} )$   
 $A \cap A_0 \neq A_0$

$$= P( EU_i^2 + \delta' EX_i X_i' \delta + o_p(1) + |A| \frac{2}{n} > EU_i^2 + |A_0| \frac{2}{n} + o_p(1) )$$

$$= P( \underbrace{(|A| - |A_0|)}_{\geq 0} \underbrace{\frac{2}{n}}_{=o(1)} + \underbrace{\delta' EX_i X_i' \delta}_{>0} > 0 )$$

$\rightarrow 1$ .

Caso 2 :  $P( BIC_n(A) > BIC_n(A_0) ) = P( SSR_n(A) + |A| \log n > SSR_n(A_0) + |A_0| \log n )$   
 $s.a.$   
 $A_0 \subset A$

$$= P( \sum_{i=1}^n u_i^2 + O_p(1) + o_p(1) + |A| \log n > \sum_{i=1}^n u_i^2 + O_p(1) + o_p(1) + |A_0| \log n )$$

$$= P( O_p(1) + o_p(1) > \underbrace{(|A_0| - |A|)}_{\leq 0} \log n ) \rightarrow 1.$$

En conclusión AIC puede seleccionar correctamente los regresores relevantes (caso 1) pero en muestras grandes puede seguir seleccionando regresores irrelevantes (más conservador).

## 1.2. Inferencia Post Selección

Si tenemos un proceso de selección consistente

$$P( \hat{A}_n = A_0 ) \rightarrow 1, \quad (\text{por ejemplo con BIC})$$

entonces

$$P( \sqrt{n} ( \hat{\beta}_{n,j}(\hat{A}_n) - \beta_j ) \leq u ) = P( \sqrt{n} ( \hat{\beta}_{n,j}(A_0) - \beta_j ) \leq u ) + o(1)$$

proof:

$$\begin{aligned} P(\sqrt{n}(\hat{\beta}_{n,j}(\hat{A}_n) - \beta_j) \leq u) &= P(\sqrt{n}(\hat{\beta}_{n,j}(\hat{A}_n) - \beta_j) \leq u, \hat{A}_n = A_0) \\ &\quad + \\ &\quad P(\sqrt{n}(\hat{\beta}_{n,j}(\hat{A}_n) - \beta_j) \leq u, \hat{A}_n \neq A_0) \\ &\leq P(\sqrt{n}(\hat{\beta}_{n,j}(\hat{A}_n) - \beta_j) \leq u, \hat{A}_n = A_0) \\ &\quad + P(\hat{A}_n \neq A_0) \\ &= P(\sqrt{n}(\hat{\beta}_{n,j}(\hat{A}_n) - \beta_j) \leq u, \hat{A}_n = A_0) + o(1) \\ &= P(\sqrt{n}(\hat{\beta}_{n,j}(\hat{A}_n) - \beta_j) \leq u \mid \hat{A}_n = A_0) P(\hat{A}_n = A_0) + o(1) \\ &= P(\sqrt{n}(\hat{\beta}_{n,j}(A_0) - \beta_j) \leq u) (1 + o(1)) + o(1) \\ &= P(\sqrt{n}(\hat{\beta}_{n,j}(A_0) - \beta_j) \leq u) + o(1). \quad \blacksquare \end{aligned}$$

Por lo tanto, todo procedimiento que seleccione consistentemente los regresores es también uno que cumple la propiedad de ser oráculo.

### \* Paréntesis: Delta Method

La prueba anterior ha mostrado una herramienta útil de como probar distintos teoremas cuando hay consistencia.

Teorema - (Método Delta) Sea  $\hat{\beta}_n$  consistente y  $h(\cdot)$  una función continua en  $\beta_0$  (verdadero valor). Entonces, si  $\text{Var}(X_i u_i) = O(1)$  tenemos

$$\sqrt{n}(\hat{\beta}_n - \beta_0) \xrightarrow{d} \mathcal{N}\left(0, \underbrace{(E X_i X_i')^{-1} E[u_i^2 X_i X_i'] (E X_i X_i')^{-1}}_V\right)$$

$$\sqrt{n}(h(\hat{\beta}_n) - h(\beta_0)) \xrightarrow{d} \mathcal{N}\left(0, \underbrace{\nabla h(\beta_0)^T}_K \underbrace{V}_K \underbrace{\nabla h(\beta_0)}_{K^T}\right)$$

proof:

$$h(\hat{\beta}_n) = h(\beta_0) + \nabla h(\beta_n^*) (\hat{\beta}_n - \beta_0) \quad * \text{ Mean value Expansion}$$

$\beta_0 \leq \beta_n^* \leq \hat{\beta}_n$

Tenemos que

$$\begin{aligned} P( \|\nabla h(\beta_n^*) - \nabla h(\beta_0)\| > \varepsilon ) &= P( \|\nabla h(\beta_n^*) - \nabla h(\beta_0)\| > \varepsilon, \|\beta_n^* - \beta_0\| > \delta ) \\ &\quad + P( \|\nabla h(\beta_n^*) - \nabla h(\beta_0)\| > \varepsilon, \|\beta_n^* - \beta_0\| \leq \delta ) \\ &\leq P( \|\beta_n^* - \beta_0\| > \delta ) \\ &\quad + P( \|\nabla h(\beta_n^*) - \nabla h(\beta_0)\| > \varepsilon, \|\beta_n^* - \beta_0\| \leq \delta ) \\ &\leq P( \underbrace{\|\beta_n^* - \beta_0\|}_{o(1)} > \delta ) \\ &\quad + P( \|\nabla h(\beta_n^*) - \nabla h(\beta_0)\| > \varepsilon, \|\beta_n^* - \beta_0\| \leq \delta ) \\ &\stackrel{\text{def de continuidad}}{=} P( \underbrace{\frac{\varepsilon}{2}}_{> \varepsilon} > \varepsilon ) + o(1) \quad \text{si elegimos un } \delta \text{ suficientemente pequeño.} \\ &= o(1). \end{aligned}$$

Entonces,

$$\begin{aligned} h(\beta_n^*) &= h(\beta_0) + \nabla h(\beta_n^*) (\beta_n^* - \beta_0) \\ &= h(\beta_0) + [\nabla h(\beta_0) + o_p(1)] (\beta_n^* - \beta_0) \end{aligned}$$

Reordenando y multiplicando por  $\sqrt{n}$

$$\begin{aligned} \sqrt{n}(h(\beta_n^*) - h(\beta_0)) &= [\nabla h(\beta_0) + o_p(1)] \sqrt{n}(\beta_n^* - \beta_0) \\ &= \nabla h(\beta_0) \sqrt{n}(\beta_n^* - \beta_0) + o_p(1) O_p(1) \\ &= \nabla h(\beta_0) \sqrt{n}(\beta_n^* - \beta_0) + o_p(1) \\ &\xrightarrow{d} \mathcal{N}(0, \underbrace{\nabla h(\beta_0)^T}_{1 \times k} \underbrace{V}_{k \times k} \underbrace{\nabla h(\beta_0)}_{k \times 1}). \quad \blacksquare \end{aligned}$$