

Transparencias para Métodos Cuantitativos I

Marcos Bujosa

29 de mayo de 2023

Transparencias para Métodos Cuantitativos I

Marcos Bujosa

29 de mayo de 2023

¿Por qué modelizar?

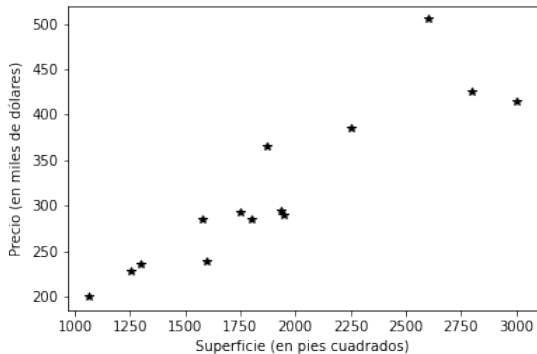
Modelado consiste en intentar ajustar un modelo matemático (estadístico) a un conjunto de datos (“la muestra”).

Un modelo es útil cuando (pese a ser *simple*) *capta las características* de los datos que consideramos más interesantes.

Ejemplos de objetivos por los que construir modelos

- ▶ Estimación: sensibilidad de un valor financiero a movimientos de un índice de referencia (evaluación de exposición al riesgo y cobertura con derivados sobre el índice).
- ▶ Previsión: probabilidad de impago de préstamos (función de las características de la operación y del solicitante).
- ▶ Simulación: rendimiento de una cartera de valores en diferentes escenarios.
- ▶ Control: (*bancos centrales*) intervención de tipos para controlar la inflación.

¿Hay relación entre tamaño y precio de una vivienda?



$$Precio_n = a + b(Superficie_n) + OtrasCosas_n$$

Función de consumo

Supongamos que consumo (*con*) y renta disponible (*rd*) de las familias siguen la relación:

$$con = \beta_1 + \beta_2 rd + otrascosas$$

donde *otrascosas* son otros aspectos distintos de la renta (activos financieros, estado de ánimo, edad, lugar de residencia, etc.).

Si disponemos datos de *consumo* y *renta disponible* de N familias como vectores de \mathbb{R}^N , *podemos construir una aproximación* (\widetilde{con}) *del consumo* mediante una combinación lineal de la renta disponible (*rd*) y de un término cte. (1) *ignorando las otrascosas*:

$$\widetilde{con} = \widetilde{\beta}_1 \mathbf{1} + \widetilde{\beta}_2 rd = \begin{bmatrix} \mathbf{1}; & rd; \end{bmatrix} \begin{pmatrix} \widetilde{\beta}_1 \\ \widetilde{\beta}_2 \end{pmatrix}.$$

Nomenclatura y notación

- ▶ *regresando*: vector de datos de *consumo* (*con*)
- ▶ *regresores*: vector de unos (**1**) y de rentas disp. (*rd*):

$$\mathbf{X} = [\mathbf{1}; \mathbf{rd};], \quad \text{donde} \quad \mathbf{X}_{|1} = \mathbf{1} \quad \text{y} \quad \mathbf{X}_{|2} = \mathbf{rd}.$$

- ▶ *vector de parámetros*: $\tilde{\boldsymbol{\beta}} = \begin{pmatrix} \widetilde{\beta_1} \\ \widetilde{\beta_2} \end{pmatrix}$

Otro ejemplo: Un modelo para los salarios

$$\textit{salario} = \beta_1 + \beta_2 \textit{educ} + \beta_3 \textit{exper} + \beta_4 \textit{IQ} + \textit{otras cosas}$$

(disponiendo de datos de N trabajadores) el **ajuste** es

$$\widetilde{\textit{salario}} = \widetilde{\beta_1} \mathbf{1} + \widetilde{\beta_2} \textit{educ} + \widetilde{\beta_3} \textit{exper} + \widetilde{\beta_4} \textit{iq}$$

Ajuste MCO. Función lineal en los parámetros

- La aproximación o ajuste $\tilde{\mathbf{y}}$ es una combinación lineal de los regresores $\mathbf{X}_{|j}$:

$$\begin{pmatrix} \tilde{y}_1 \\ \vdots \\ \tilde{y}_N \end{pmatrix} = \tilde{\beta}_1 \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} + \tilde{\beta}_2 \begin{pmatrix} x_{12} \\ \vdots \\ x_{N2} \end{pmatrix} + \cdots + \tilde{\beta}_k \begin{pmatrix} x_{1k} \\ \vdots \\ x_{Nk} \end{pmatrix}$$

ó

$$\begin{aligned} \tilde{\mathbf{y}} &= \tilde{\beta}_1 \mathbf{1} + \tilde{\beta}_2 \mathbf{X}_{|2} + \tilde{\beta}_3 \mathbf{X}_{|3} + \cdots + \tilde{\beta}_k \mathbf{X}_{|k} \\ &= \begin{bmatrix} \mathbf{1}; \mathbf{X}_{|2}; \dots \mathbf{X}_{|k}; \end{bmatrix} \begin{pmatrix} \tilde{\beta}_1 \\ \vdots \\ \tilde{\beta}_k \end{pmatrix} = \mathbf{X} \tilde{\boldsymbol{\beta}}. \end{aligned}$$

Así los valores ajustados son: $\tilde{\mathbf{y}} = \mathbf{X} \tilde{\boldsymbol{\beta}} \in \mathbb{R}^N$

Datos del ejemplo del precio de las viviendas

Precio viviendas (miles de \$) y superficie útil (pies al cuadrado)

14 casas unifamiliares en *University City*. San Diego, California. Año 1990.

	price (y)	sqft (x)	price (\tilde{y})
0	199.9	1065	
1	228	1254	
2	235	1300	
3	285	1577	
4	239	1600	
5	293	1750	
6	285	1800	
7	365	1870	
8	295	1935	
9	290	1948	
10	385	2254	
11	505	2600	
12	425	2800	
13	415	3000	

Si asumimos que el precio y se relaciona con la superficie x del siguiente modo:

$$y_n = a + b x_n + \text{otras cosas}_n,$$

podemos “*aproximar*” el vector de precios, \mathbf{y} , con una **combinación lineal de los regresores**:

$$\tilde{\mathbf{y}} = \tilde{\beta}_1 \mathbf{1} + \tilde{\beta}_2 \mathbf{x} = \begin{bmatrix} 1; & \mathbf{x}; \end{bmatrix} \begin{pmatrix} \tilde{\beta}_1 \\ \tilde{\beta}_2 \end{pmatrix} = \mathbf{X} \tilde{\boldsymbol{\beta}}.$$

Los precios ajustados como combinación lineal de los regresores

$$\tilde{\mathbf{y}} = (\mathbf{x}_{|1})\tilde{\beta}_1 + (\mathbf{x}_{|2})\tilde{\beta}_2 = \begin{pmatrix} \tilde{y}_1 \\ \tilde{y}_2 \\ \tilde{y}_3 \\ \tilde{y}_4 \\ \tilde{y}_5 \\ \tilde{y}_6 \\ \tilde{y}_7 \\ \tilde{y}_8 \\ \tilde{y}_9 \\ \tilde{y}_{10} \\ \tilde{y}_{11} \\ \tilde{y}_{12} \\ \tilde{y}_{13} \\ \tilde{y}_{14} \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} \tilde{\beta}_1 + \begin{pmatrix} 1065 \\ 1254 \\ 1300 \\ 1577 \\ 1600 \\ 1750 \\ 1800 \\ 1870 \\ 1935 \\ 1948 \\ 2254 \\ 2600 \\ 2800 \\ 3000 \end{pmatrix} \tilde{\beta}_2 = \begin{pmatrix} 1 & 1065 \\ 1 & 1254 \\ 1 & 1300 \\ 1 & 1577 \\ 1 & 1600 \\ 1 & 1750 \\ 1 & 1800 \\ 1 & 1870 \\ 1 & 1935 \\ 1 & 1948 \\ 1 & 2254 \\ 1 & 2600 \\ 1 & 2800 \\ 1 & 3000 \end{pmatrix} \begin{pmatrix} \tilde{\beta}_1 \\ \tilde{\beta}_2 \end{pmatrix} = \mathbf{X}\tilde{\beta};$$

El precio ajustado para el séptimo piso de la muestra será

$$\tilde{y}_7 = {}_{7|}\tilde{\mathbf{y}} = {}_{7|}\mathbf{X}\tilde{\beta} = (1, 1800,) \begin{pmatrix} \tilde{\beta}_1 \\ \tilde{\beta}_2 \end{pmatrix} = (1)\tilde{\beta}_1 + (1800)\tilde{\beta}_2 \neq y_7.$$

Pero *¿qué criterio empleamos para elegir $\tilde{\beta}_1$ y $\tilde{\beta}_2$ en el ajuste $\tilde{\mathbf{y}} = \mathbf{X}\tilde{\beta}$?*

Dados \mathbf{X} e \mathbf{y} , el “*error de ajuste*” empleando $\tilde{\beta}$ es

$$\tilde{\mathbf{e}} = \mathbf{y} - \mathbf{X}\tilde{\beta} = \mathbf{y} - \tilde{\mathbf{y}};$$

Así, descomponemos los datos observados \mathbf{y} en: $\mathbf{y} = \tilde{\mathbf{y}} + \tilde{\mathbf{e}}$.

Llamamos “Suma de los Residuos al Cuadrado” (*SRC*) del ajuste $\tilde{\mathbf{y}}$ a

$$SRC(\tilde{\beta}) \equiv \sum_{n=1}^N (\tilde{e}_n)^2 = \tilde{\mathbf{e}} \cdot \tilde{\mathbf{e}} = \|\tilde{\mathbf{e}}\|^2$$

es decir, al cuadrado de la longitud del vector $\tilde{\mathbf{e}} = (\mathbf{y} - \tilde{\mathbf{y}})$.

- **Objetivo.** Encontrar una combinación lineal de los regresores

$$\text{Sean } \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \text{ y } \mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}. \quad \text{Buscamos una } \mathbf{X}\tilde{\boldsymbol{\beta}}.$$

- **Criterio** de búsqueda (ajuste MCO)

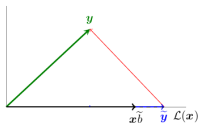
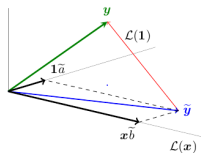
Elegir el vector $\tilde{\boldsymbol{\beta}}$ tal que $\mathbf{X}\tilde{\boldsymbol{\beta}}$ esté *lo más cerca* de \mathbf{y} ; es decir, tal que la componente $\tilde{\mathbf{e}}$ sea lo más pequeña posible en la descomposición:

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\tilde{\boldsymbol{\beta}} + \tilde{\mathbf{e}} \\ &= \tilde{\mathbf{y}} + \tilde{\mathbf{e}}. \end{aligned}$$

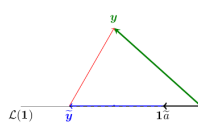
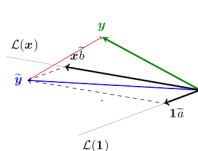
Geometría de un mal ajuste lineal

Un \tilde{a} demasiado pequeño y un \tilde{b} demasiado grande.

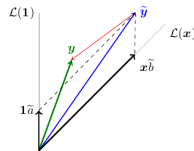
Desde un lado



Desde el otro



Desde arriba



$$\mathbf{X} = \begin{bmatrix} 1 & x; \end{bmatrix}; \quad \tilde{\beta} = \begin{pmatrix} \tilde{a} \\ \tilde{b} \end{pmatrix}; \quad \boxed{\tilde{y} = \mathbf{X}\tilde{\beta}}; \quad \mathbf{y} = \tilde{y} + \tilde{e}; \quad \tilde{e} = \mathbf{y} - \tilde{y}$$

Como el vector $\hat{e} = (y - \hat{y})$ es **mínimo** si es *perpendicular* a cada regresor. Es decir, si:

$$\hat{e} \perp \mathbf{X}_{|j} \Leftrightarrow \mathbf{0} = \mathbf{X}^T \hat{e} = \mathbf{X}^T (y - \hat{y}).$$

Tenemos que,

$$\hat{y} = \mathbf{X} \hat{\beta} \Leftrightarrow \mathbf{X}^T (y - \mathbf{X} \hat{\beta}) = \mathbf{0} \Leftrightarrow \mathbf{X}^T y - \mathbf{X}^T \mathbf{X} \hat{\beta} = \mathbf{0}$$

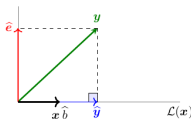
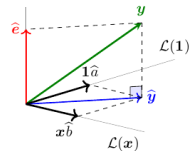
$$\hat{y} = \mathbf{X} \hat{\beta} \quad \text{si} \quad \boxed{(\mathbf{X}^T \mathbf{X}) \hat{\beta} = \mathbf{X}^T y} \quad (1)$$

Las soluciones $\hat{\beta}$ son los parámetros del ajuste MCO: $\hat{y} = \mathbf{X} \hat{\beta}$

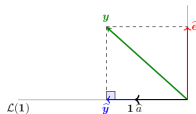
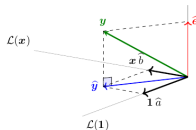
Ajuste MCO (Geometría de la proyección ortogonal)

$$\hat{\mathbf{e}} \perp \mathbf{X} \iff \hat{\boldsymbol{\beta}} \text{ es tal que } (\mathbf{X}^T \mathbf{X}) \hat{\boldsymbol{\beta}} = \mathbf{X}^T \mathbf{y}$$

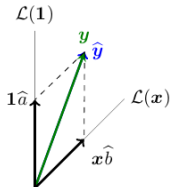
Desde un lado



Desde el otro



Desde arriba



$$\mathbf{X} = \begin{bmatrix} 1 & x; \end{bmatrix}; \quad \hat{\boldsymbol{\beta}} = \begin{pmatrix} \hat{a} \\ \hat{b} \end{pmatrix}; \quad \boxed{\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}}}; \quad \mathbf{y} = \hat{\mathbf{y}} + \hat{\mathbf{e}}; \quad \hat{\mathbf{e}} = \mathbf{y} - \hat{\mathbf{y}}$$

¿Es el sistema de ecuaciones normales determinado?

Puesto que

$$\mathbf{X}\hat{\boldsymbol{\beta}} = \hat{\mathbf{y}} \quad \Longleftrightarrow \quad (\mathbf{X}^T\mathbf{X})\hat{\boldsymbol{\beta}} = \mathbf{X}^T\mathbf{y}, \quad \text{donde } \mathbf{X} ;$$

$N \times k$

ambos sistemas tendrán *solución única si y sólo* si sus matrices de coeficientes son de *rango* k .

En tal caso, multiplicando ambos lados de las ecuaciones normales por $(\mathbf{X}^T\mathbf{X})^{-1}$ tenemos que

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} \quad (2)$$

es *la única solución*.

Ecuación de salarios:

Supongamos el siguiente modelo

$$Salar_n = e^{\beta_1 + \beta_2(educ_n) + \beta_3(antig_n) + \beta_4(exper_n) + otrascosas_n};$$

Tomando logaritmos tenemos una nueva variable $\ln(Salar_n)$ que podremos ajustar con un modelo lineal en los parámetros pues,

$$\ln(Salar_n) = \beta_1 + \beta_2(educ_n) + \beta_3(antig_n) + \beta_4(exper_n) + otrascosas_n.$$

Pero ¿qué pasará si jamás ningún trabajador cambió de empresa?

La matriz $\mathbf{X}^T \mathbf{X}$ será singular. De hecho, como *experiencia* y *antigüedad* coinciden, como mucho sólo podemos calcular su **efecto conjunto**:

$$\ln(Salar_n) = \beta_1 + \beta_2(educ_n) + (\beta_3 + \beta_4)exper_n + otrascosas_n,$$