

# Regresión lineal: mínimos cuadrados

Matías Carrasco

4 de septiembre de 2023

## Índice

|   |    |
|---|----|
| 1. Modelo lineal univariado                       | 1  |
| 2. Cálculo de los coeficientes                    | 2  |
| 3. Modelo lineal univariado en notación matricial | 4  |
| 4. Ambos cálculos son equivalentes                | 5  |
| 5. Interpretación de la correlación               | 6  |
| 6. Modelo lineal multivariado                     | 7  |
| 7. Interpretación de $Z^\top u$                   | 8  |
| 8. Interpretación de $Z^\top Z$                   | 9  |
| 9. Efecto de la multicolinealidad                 | 9  |
| 10. Regresión lineal con funciones base           | 11 |
| 11. Regularización Ridge                          | 11 |
| 12. Multicolinealidad en regresión polinomial     | 12 |

## 1. Modelo lineal univariado

Consideremos el siguiente setting:

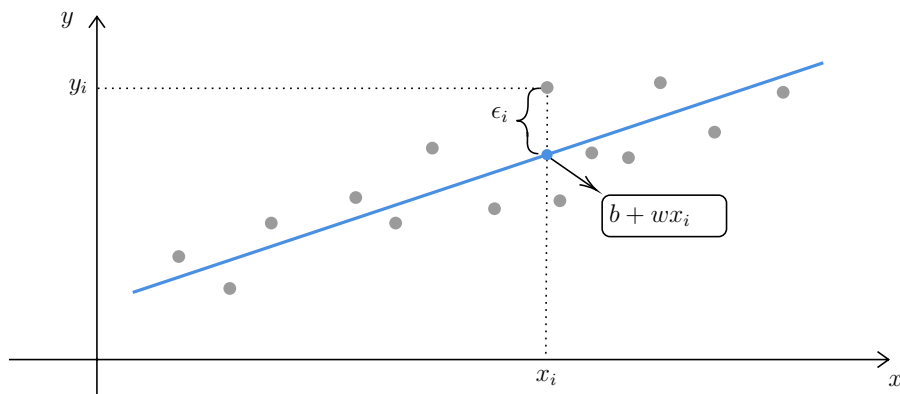
- Espacio de atributos  $\mathcal{X} \subset \mathbb{R}$ .

- Espacio de etiquetas  $\mathcal{Y} \subset \mathbb{R}$ .
- Espacio de observaciones  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ .
- Conjunto de datos  $S = \{(x_i, y_i)\}_{i=1}^N$ .
- Función de pérdida cuadrática:  $\text{Loss}(y, y') = (y - y')^2$

El modelo lineal consiste en aproximar la relación entre  $x$  e  $y$  mediante una recta

$$y = f(x; \boldsymbol{\theta}) = b + wx$$

en donde el vector de parámetros es  $\boldsymbol{\theta} = [b, w]^\top \in \mathbb{R}^2$ .



La pérdida promedio en todo el conjunto de datos  $S$  se llama MSE (Mean Square Error) y la usaremos para elegir los mejores parámetros:

$$L(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i; \boldsymbol{\theta}))^2 = \frac{1}{N} \sum_{i=1}^N (y_i - b - wx_i)^2$$

de modo que

$$\hat{b}, \hat{w} = \arg \min_{b, w} \{L(b, w)\}.$$

## 2. Cálculo de los coeficientes

Observación: dados números reales  $A_1, \dots, A_N$

$$\arg \min_a \left\{ \sum_{i=1}^N (A_i - a)^2 \right\} = \bar{A}$$

en donde denotamos  $\bar{A}$  el promedio  $\frac{1}{N} \sum_i A_i$ .

Por la observación, si fijamos  $w$ , el mejor  $b$  es

$$\arg \min_b \left\{ \frac{1}{N} \sum_{i=1}^N \underbrace{(y_i - wx_i - b)}_{A_i}^2 \right\} = \bar{y} - w\bar{x}$$

Luego basta encontrar  $\hat{w}$

$$\begin{aligned} L(w) &= L(\bar{y} - w\bar{x}, w) \\ &= \frac{1}{N} \sum_{i=1}^N (y_i - wx_i - \bar{y} + w\bar{x})^2 \\ &= \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y} - w(x_i - \bar{x}))^2 \\ &= \frac{1}{N} \sum_{i=1}^N \underbrace{(y_i - \bar{y})^2}_{s_y^2} - 2w \frac{1}{N} \sum_{i=1}^N \underbrace{(y_i - \bar{y})(x_i - \bar{x})}_{s_{xy}} + w^2 \frac{1}{N} \sum_{i=1}^N \underbrace{(x_i - \bar{x})^2}_{s_x^2} \\ &= s_y^2 - 2s_{xy}w + s_x^2w^2 \quad (\text{polinomio de grado 2 en } w) \end{aligned}$$

Sus raíces son

$$\frac{2s_{xy} \pm \sqrt{4s_{xy}^2 - 4s_x^2s_y^2}}{2s_x^2}$$

pero como  $L(w) \geq 0$  a lo sumo tiene una raíz doble, de donde obtenemos la condición:

$$\frac{s_{xy}^2}{s_x^2s_y^2} \leq 1 \Leftrightarrow \left| \frac{s_{xy}}{s_x s_y} \right| \leq 1$$

El número  $r = \frac{s_{xy}}{s_x s_y}$  se llama coeficiente de correlación y acabamos de probar que  $r \in [-1, 1]$ . Más aún, probamos que  $r = \pm 1$  si, y sólo si,  $y_i = b + wx_i$  para todo  $i$ .

Volviendo a  $L(w)$  y completando el cuadrado tenemos:

$$\begin{aligned} L(w) &= s_y^2 - 2s_{xy}w + s_x^2w^2 \\ &= s_y^2 - 2rs_y s_x w + s_x^2w^2 \\ &= s_y^2 - r^2s_y^2 + r^2s_y^2 - 2rs_y s_x w + s_x^2w^2 \\ &= s_y^2(1 - r^2) + (rs_y - s_x w)^2 \end{aligned}$$

De aquí vemos que el mínimo se alcanza en  $\hat{w} = r \frac{s_y}{s_x}$  y vale  $L(\hat{w}) = s_y(1 - r^2)$ .

Si llamamos  $\epsilon_i = y_i - (b + wx_i)$  llegamos a la siguiente fórmula para la MSE mínima:

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N \epsilon_i^2 = s_y^2(1 - r^2)$$

En estadística se suele usar  $r^2$  como medida de ajuste en lugar de la MSE.

### 3. Modelo lineal univariado en notación matricial

La función lineal puede escribirse como un producto escalar de vectores

$$b + wx_i = [1, x_i] \cdot \begin{bmatrix} b \\ w \end{bmatrix} = \mathbf{x}_i^\top \boldsymbol{\theta} = \boldsymbol{\theta}^\top \mathbf{x}_i$$

Considerando la matriz de diseño

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}$$

Podemos juntar todos los valores de  $\{b + wx_i\}_{i=1}^N$  en un solo producto matricial

$$\underbrace{\begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}}_{\mathbf{X}} \cdot \underbrace{\begin{bmatrix} b \\ w \end{bmatrix}}_{\boldsymbol{\theta}} = \underbrace{\begin{bmatrix} b + wx_1 \\ b + wx_2 \\ \vdots \\ b + wx_N \end{bmatrix}}_{\mathbf{X}\boldsymbol{\theta}}$$

y compararlo con el vector de targets o etiquetas

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}.$$

De este modo la función de pérdida queda

$$L(\boldsymbol{\theta}) = \frac{1}{N} \underbrace{\sum_{i=1}^N (y_i - b - wx_i)^2}_{\text{norma}} = \frac{1}{N} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|^2$$

$$= \frac{1}{N} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})$$

Haciendo las distributivas obtenemos

$$L(\boldsymbol{\theta}) = \frac{1}{N} (\mathbf{y}^\top - \boldsymbol{\theta}^\top \mathbf{X}^\top) (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})$$

$$= \frac{1}{N} \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\boldsymbol{\theta} - \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{y} + \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta}$$

Derivando respecto al vector  $\boldsymbol{\theta}$  tenemos

$$\frac{\partial L}{\partial \boldsymbol{\theta}} = -\mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} = -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} = \mathbf{0}$$

de donde obtenemos

$$\boxed{\hat{\boldsymbol{\theta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}}$$

## 4. Ambos cálculos son equivalentes

Para obtener  $\hat{\boldsymbol{\theta}}$  a partir del producto matricial seguimos estos pasos:

1. Calculamos el producto transpuesto de la matriz de diseño  $\mathbf{X}^\top \mathbf{X}$ :

$$\mathbf{X}^\top \mathbf{X} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_N \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix} = \begin{bmatrix} N & \sum_{i=1}^N x_i \\ \sum_{i=1}^N x_i & \sum_{i=1}^N x_i^2 \end{bmatrix}$$

2. Calculamos la inversa de la matriz resultante:

$$(\mathbf{X}^\top \mathbf{X})^{-1} = \frac{1}{N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i\right)^2} \begin{bmatrix} \sum_{i=1}^N x_i^2 & -\sum_{i=1}^N x_i \\ -\sum_{i=1}^N x_i & N \end{bmatrix}$$

3. Calculamos el producto de la matriz de diseño transpuesta  $\mathbf{X}^\top$  y el vector de targets  $\mathbf{y}$ :

$$\mathbf{X}^\top \mathbf{y} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_N \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N y_i \\ \sum_{i=1}^N x_i y_i \end{bmatrix}$$

4. Finalmente, multiplicamos el resultado del paso 2 con el resultado del paso 3 para obtener los coeficientes  $\hat{\boldsymbol{\theta}}$ :

$$\begin{aligned} \hat{\boldsymbol{\theta}} &= \frac{1}{N \sum_{i=1}^N x_i^2 - \left( \sum_{i=1}^N x_i \right)^2} \begin{bmatrix} \sum_{i=1}^N x_i^2 & - \sum_{i=1}^N x_i \\ - \sum_{i=1}^N x_i & N \end{bmatrix} \begin{bmatrix} \sum_{i=1}^N y_i \\ \sum_{i=1}^N x_i y_i \end{bmatrix} \\ &= \frac{1}{N^2 \overline{x^2} - N^2 \bar{x}^2} \begin{bmatrix} N \overline{x^2} & -N \bar{x} \\ -N \bar{x} & N \end{bmatrix} \begin{bmatrix} N \bar{y} \\ N \overline{xy} \end{bmatrix} \\ &= \frac{1}{\overline{x^2} - \bar{x}^2} \begin{bmatrix} \overline{x^2} & -\bar{x} \\ -\bar{x} & 1 \end{bmatrix} \begin{bmatrix} \bar{y} \\ \overline{xy} \end{bmatrix} \\ &= \frac{1}{\overline{x^2} - \bar{x}^2} \begin{bmatrix} \overline{x^2 \bar{y}} - \bar{x} \overline{xy} \\ -\bar{x} \bar{y} + \overline{xy} \end{bmatrix} \\ &= \frac{1}{s_x^2} \begin{bmatrix} (s_x^2 + \bar{x}^2) \bar{y} - \bar{x} (s_{xy} + \bar{x} \bar{y}) \\ s_{xy} \end{bmatrix} = \begin{bmatrix} \bar{y} - \bar{x} r \frac{s_y}{s_x} \\ r \frac{s_y}{s_x} \end{bmatrix} = \begin{bmatrix} \hat{b} \\ \hat{w} \end{bmatrix} \end{aligned}$$

## 5. Interpretación de la correlación

Una forma directa de interpretar la correlación es como  $r^2 = \text{cor}(\mathbf{y}, \hat{\mathbf{y}})$  siendo  $\hat{\mathbf{y}} = \mathbf{X}\boldsymbol{\theta}$  el vector de predicciones. Esto es así pues la correlación es invariante por transformaciones lineales.

En estadística aplicada se suele interpretar  $r^2$  como porcentaje de la variabilidad de  $\mathbf{y}$  explicada por  $\mathbf{x}$ . De hecho, la varianza de  $\hat{\mathbf{y}}$  es

$$s_{\hat{\mathbf{y}}}^2 = \text{var}(\hat{\mathbf{y}}) = \text{var}(\hat{b} + \hat{w}\mathbf{x}) = \hat{w}^2 \text{var}(\mathbf{x}) = r^2 \frac{s_y^2}{s_x^2} s_x^2 = r^2 s_y^2$$

de donde

$$\% \text{ de variación explicada} = \frac{\text{var}(\hat{\mathbf{y}})}{\text{var}(\mathbf{y})} \times 100 = \frac{r^2 s_y^2}{s_y^2} \times 100 = r^2 \times 100$$

## 6. Modelo lineal multivariado

El setting para el modelo multivariado es análogo salvo que

- Espacio de atributos  $\mathcal{X} \subset \mathbb{R}^D$ .

pues disponemos de  $D$  atributos.

Ahora el modelo es

$$y = f(\mathbf{x}; \boldsymbol{\theta}) = \boldsymbol{\theta}^\top \mathbf{x}$$

en donde los vectores  $\boldsymbol{\theta}$  y  $\mathbf{x}$  están dados por

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_D \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} 1 \\ x^{(1)} \\ \vdots \\ x^{(D)} \end{bmatrix}$$

y el producto representa

$$y = [b \quad w_1 \quad \dots \quad w_D] \begin{bmatrix} 1 \\ x^{(1)} \\ \vdots \\ x^{(D)} \end{bmatrix} = b + \sum_{j=1}^D w_j x^{(j)}$$

Considerando la matriz de diseño y el vector de targets

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_1^{(D)} \\ 1 & x_2^{(1)} & \dots & x_2^{(D)} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_N^{(1)} & \dots & x_N^{(D)} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

los cálculos son idénticos y la solución es la misma  $\hat{\boldsymbol{\theta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ .

Al mirar en detalle la función de pérdida

$$L(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \left( \underbrace{y_i - \sum_{j=1}^D w_j x_i^{(j)}}_{A_i} - b \right)^2$$

vemos que si fijamos  $\mathbf{w}$  el mejor  $b$  es  $\bar{y} - \sum_{j=1}^D w_j \bar{x}^{(j)}$ . Es decir, nuevamente  $b$  cumple el rol de ajustar promedios.

Si escribimos la función de pérdida usando solamente  $\mathbf{w}$  obtenemos

$$L(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \left( (y_i - \bar{y}) - \sum_{j=1}^D w_j (x_i^{(j)} - \bar{x}^{(j)}) \right)^2$$

Si denotamos  $u_i = y_i - \bar{y}$  y  $z_i^{(j)} = x_i^{(j)} - \bar{x}^{(j)}$  podemos reescribir la pérdida

$$L(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \left( u_i - \sum_{j=1}^D w_j z_i^{(j)} \right)^2 = \frac{1}{N} (\mathbf{u} - \mathbf{Z}\mathbf{w})^\top (\mathbf{u} - \mathbf{Z}\mathbf{w})$$

en donde  $\mathbf{u}$  es el vector de targets centrado y  $\mathbf{Z}$  es la matriz de diseño cuyas columnas están centradas

$$\mathbf{Z} = \begin{bmatrix} \mathbf{z}_1^\top \\ \mathbf{z}_2^\top \\ \vdots \\ \mathbf{z}_N^\top \end{bmatrix} = \begin{bmatrix} z_1^{(1)} & \cdots & z_1^{(D)} \\ z_2^{(1)} & \cdots & z_2^{(D)} \\ \vdots & \vdots & \vdots \\ z_N^{(1)} & \cdots & z_N^{(D)} \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_N \end{bmatrix}$$

La solución óptima está dada por  $\hat{\mathbf{w}} = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{u}$ .

## 7. Interpretación de $\mathbf{Z}^\top \mathbf{u}$

El producto es

$$\mathbf{Z}^\top \mathbf{u} = \begin{bmatrix} z_1^{(1)} & | & z_2^{(1)} & | & \cdots & | & z_N^{(1)} \\ \vdots & | & \vdots & | & \ddots & | & \vdots \\ z_1^{(D)} & | & z_2^{(D)} & | & \cdots & | & z_N^{(D)} \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_N \end{bmatrix} = \sum_{i=1}^N \mathbf{z}_i u_i$$

Esto nos dará como resultado un vector columna de dimensiones  $D \times 1$  donde las entradas son las sumas ponderadas de los  $u_i$  por las respectivas coordenadas de los  $\mathbf{z}_i$ :

$$\mathbf{Z}^\top \mathbf{u} = \begin{bmatrix} \sum_{i=1}^N z_i^{(1)} u_i \\ \vdots \\ \sum_{i=1}^N z_i^{(D)} u_i \end{bmatrix} = N \text{cov}([\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(D)}], \mathbf{y})$$



## 8. Interpretación de $\mathbf{Z}^\top \mathbf{Z}$

Por otro lado el producto matricial  $\mathbf{Z}^\top \mathbf{Z}$  es:

$$\mathbf{Z}^\top \mathbf{Z} = \begin{bmatrix} \mathbf{z}_1 & \mathbf{z}_2 & \cdots & \mathbf{z}_N \end{bmatrix} \begin{bmatrix} \mathbf{z}_1^\top \\ \mathbf{z}_2^\top \\ \vdots \\ \mathbf{z}_N^\top \end{bmatrix}$$

Es una matriz cuadrada de tamaño  $D \times D$ , donde cada elemento  $(j, k)$  es el resultado del producto escalar entre  $\mathbf{z}^{(j)}$  y  $\mathbf{z}^{(k)}$ :

$$\mathbf{Z}^\top \mathbf{Z} = \begin{bmatrix} \mathbf{z}^{(1)} \cdot \mathbf{z}^{(1)} & \cdots & \mathbf{z}^{(1)} \cdot \mathbf{z}^{(D)} \\ \vdots & \ddots & \vdots \\ \mathbf{z}^{(D)} \cdot \mathbf{z}^{(1)} & \cdots & \mathbf{z}^{(D)} \cdot \mathbf{z}^{(D)} \end{bmatrix}$$

Es decir,  $\mathbf{Z}^\top \mathbf{Z} = N \text{cov}([\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(D)}])$ .

En resumen el vector de pesos óptimo se puede escribir como

$$\hat{\mathbf{w}} = \text{cov}([\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(D)}])^{-1} \text{cov}([\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(D)}], \mathbf{y})$$

## 9. Efecto de la multicolinealidad

Supongamos que el atributo (columna de  $\mathbf{Z}$ )

$$\mathbf{x}^{(D)} = \sum_{j=1}^{D-1} \alpha_j \mathbf{x}^{(j)}$$

es combinación lineal de los otros. Esto afectara las columnas de la matriz  $\mathbf{Z}^\top \mathbf{Z}$  pues

$$\begin{aligned} \mathbf{Z}^\top \mathbf{Z} &= \begin{bmatrix} \mathbf{z}^{(1)} \cdot \mathbf{z}^{(1)} & \cdots & \mathbf{z}^{(1)} \cdot \mathbf{z}^{(D)} \\ \vdots & \ddots & \vdots \\ \mathbf{z}^{(D)} \cdot \mathbf{z}^{(1)} & \cdots & \mathbf{z}^{(D)} \cdot \mathbf{z}^{(D)} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{z}^{(1)} \cdot \mathbf{z}^{(1)} & \cdots & \sum_{j=1}^{D-1} \alpha_j \mathbf{z}^{(1)} \cdot \mathbf{z}^{(j)} \\ \vdots & \ddots & \vdots \\ \mathbf{z}^{(D)} \cdot \mathbf{z}^{(1)} & \cdots & \sum_{j=1}^{D-1} \alpha_j \mathbf{z}^{(D)} \cdot \mathbf{z}^{(j)} \end{bmatrix} \end{aligned}$$

en donde también obtenemos una columna como combinación lineal de las demás. En particular el determinante  $\det(\mathbf{Z}^\top \mathbf{Z}) = 0$  y la matriz no es invertible. Esto hace que el vector de pesos  $\mathbf{w}$  no quede unívocamente determinado y dificulte su aproximación.

Veamos con más detalle el caso bi-variado, es decir con  $D = 2$  atributos. Supongamos que  $\mathbf{x}^{(1)}$  y  $\mathbf{x}^{(2)}$  han sido estandarizados y denotemos  $r = \text{cor}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)})$  su coeficiente de correlación. Si la  $y$  está estandarizada también, vemos que  $\hat{b} = 0$ , y los pesos están dados por:

$$\hat{\mathbf{w}} = \text{cor}([\mathbf{x}^{(1)}, \mathbf{x}^{(2)}])^{-1} \text{cor}([\mathbf{x}^{(1)}, \mathbf{x}^{(2)}], \mathbf{y})$$

La matriz de correlación entre  $\mathbf{x}^{(1)}$  y  $\mathbf{x}^{(2)}$  es

$$\begin{bmatrix} 1 & r \\ r & 1 \end{bmatrix}$$

por lo que su inversa es

$$\text{cor}([\mathbf{x}^{(1)}, \mathbf{x}^{(2)}])^{-1} = \frac{1}{1-r^2} \begin{bmatrix} 1 & -r \\ -r & 1 \end{bmatrix}$$

Si denotamos el vector de correlaciones entre  $\mathbf{x}^{(1)}$  y  $\mathbf{x}^{(2)}$  con  $\mathbf{y}$  como

$$\text{cor}([\mathbf{x}^{(1)}, \mathbf{x}^{(2)}], \mathbf{y}) = \begin{bmatrix} r_1 \\ r_2 \end{bmatrix}$$

tenemos que

$$\hat{\mathbf{w}} = \frac{1}{1-r^2} \begin{bmatrix} 1 & -r \\ -r & 1 \end{bmatrix} \begin{bmatrix} r_1 \\ r_2 \end{bmatrix} = \frac{1}{1-r^2} \begin{bmatrix} r_1 - rr_2 \\ r_2 - rr_1 \end{bmatrix}$$

Las correlaciones  $r$ ,  $r_1$ , y  $r_2$  no pueden elegirse arbitrariamente pero basta que satisfagan la relación

$$1 - r_1^2 - r_2^2 + 2r_1r_2r - r^2 \geq 0$$

Esto equivale a que  $r$  pertenezca al intervalo

$$\frac{2r_1r_2 \pm \sqrt{4r_1^2r_2^2 + 4(1-r_1^2-r_2^2)}}{2} = r_1r_2 \pm \sqrt{(1-r_1^2)(1-r_2^2)}$$

Se puede ver que es válido elegir  $r_1 = \rho$ ,  $r_2 = -\rho$  y  $r = 1 - 2\rho^2$ . En ese caso los pesos son

$$w_1 = \frac{1}{2\rho} \quad w_2 = -\frac{1}{2\rho}$$

y haciendo  $\rho \rightarrow 0$  vemos que pueden ser arbitrariamente grandes. Este ejemplo muestra cómo la multicolinealidad puede dar origen a pesos muy grandes y motiva la consideración de la siguiente sección.

## 10. Regresión lineal con funciones base

La regresión con funciones de base es una técnica que permite modelar relaciones no lineales entre el target  $y$  y uno o más atributos, supongamos en este caso, un sólo atributo  $x$ .

Supongamos que deseamos modelar  $y$  no simplemente como una función lineal de  $x$ , sino utilizando una serie de funciones  $h_1(x), h_2(x), \dots, h_D(x)$ . Entonces, nuestro modelo se ve así:

$$y = f(x; \boldsymbol{\theta}) = b + w_1 h_1(x) + w_2 h_2(x) + \dots + w_D h_D(x)$$

Aquí  $h_j(x)$  son las funciones de base y pueden ser cualquier cosa. Casos particulares son:

- Regresión polinomial: cuando  $h_j(x) = x^j$ .
- Regresión trigonométrica: cuando  $h_j(x) = \cos(jx)$  o  $h_j(x) = \sin(jx)$ .

Lo que es “lineal” acerca de este modelo es la relación entre  $y$  y los coeficientes  $\boldsymbol{\theta} = (b, \mathbf{w})$ . Por eso aún se llama regresión lineal, a pesar de que la relación entre  $y$  y  $x$  puede ser altamente no lineal.

Considerando la matriz de diseño

$$\mathbf{X} = \begin{bmatrix} 1 & h_1(x_1) & \dots & h_D(x_1) \\ \vdots & \vdots & \vdots & \vdots \\ 1 & h_1(x_N) & \dots & h_D(x_N) \end{bmatrix}$$

el procedimiento es exactamente igual que en el caso multivariado.

## 11. Regularización Ridge

La regularización Ridge (inglés significa cresta) consiste en agregar un término a la función de pérdida que de algún modo penalice el tamaño del vector de pesos  $\mathbf{w}$ :

$$L_\alpha(\mathbf{w}) = L(\mathbf{w}) + \alpha \|\mathbf{w}\|^2 = L(\mathbf{w}) + \alpha \sum_{j=1}^D w_j^2$$

El objetivo es evitar los problemas de la multicolinealidad.

De hecho, podemos calcular explícitamente el valor del parámetro  $\hat{w}_\alpha$  del mismo modo que hicimos antes. Escribiendo la función de pérdida como

$$L_\alpha(\mathbf{w}) = \frac{1}{N} \mathbf{w}^\top \mathbf{Z}^\top \mathbf{Z} \mathbf{w} - \frac{2}{N} \mathbf{w}^\top \mathbf{Z}^\top \mathbf{u} + \frac{1}{N} \mathbf{u}^\top \mathbf{u} + \alpha \mathbf{w}^\top \mathbf{w}$$

derivamos e igualamos a cero:

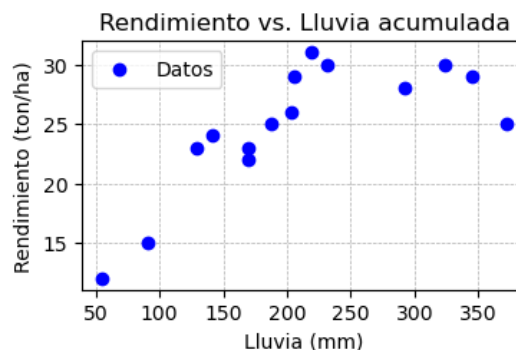
$$\begin{aligned}\frac{\partial L_\alpha}{\partial \mathbf{w}} &= \frac{2}{N} \mathbf{Z}^\top \mathbf{Z} \mathbf{w} - \frac{2}{N} \mathbf{Z}^\top \mathbf{u} + 2\alpha \mathbf{w} = 0 \\ \left( \frac{1}{N} \mathbf{Z}^\top \mathbf{Z} + 2\alpha I \right) \mathbf{w} &= \frac{1}{N} \mathbf{Z}^\top \mathbf{u} \\ \hat{\mathbf{w}}_\alpha &= (\text{cov}([\mathbf{x}^{(j)}]) + \alpha I)^{-1} \text{cov}([\mathbf{x}^{(j)}], \mathbf{y})\end{aligned}$$

Notar el efecto de  $\alpha$  en acercar proporcionalmente la matriz de covarianzas a la identidad, lo cual mejora su invertibilidad.

## 12. Multicolinealidad en regresión polinomial

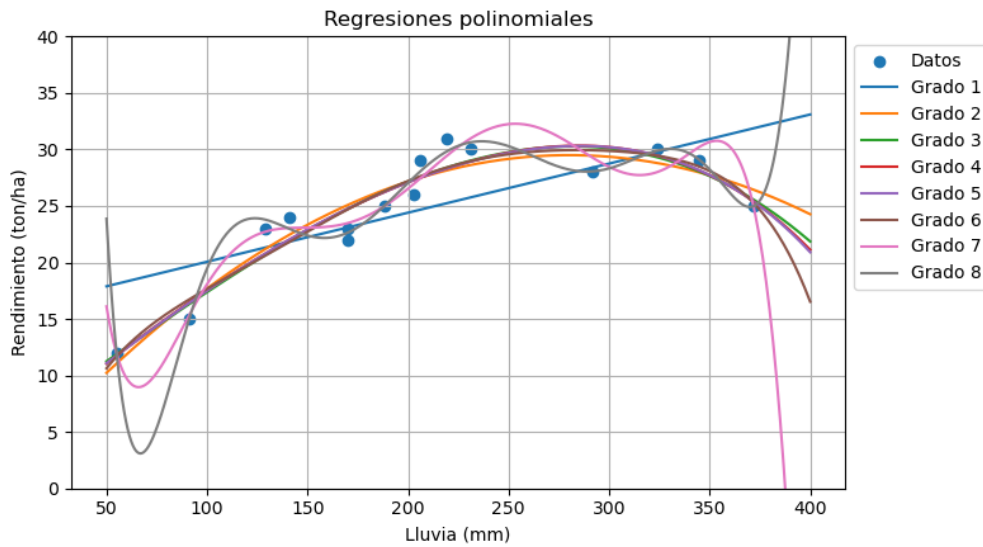
A modo de ilustración, considerar el siguiente conjunto de datos del rendimiento de un cultivo de papas en función de la lluvia acumulada:

| $x = \text{Lluvia (mm)}$ | $y = \text{Rendimiento (ton/ha)}$ |
|--------------------------|-----------------------------------|
| 206                      | 29                                |
| 188                      | 25                                |
| 219                      | 31                                |
| 372                      | 25                                |
| 345                      | 29                                |
| 231                      | 30                                |
| 203                      | 26                                |
| 170                      | 23                                |
| 55                       | 12                                |
| 91                       | 15                                |
| 292                      | 28                                |
| 141                      | 24                                |
| 129                      | 23                                |
| 170                      | 22                                |
| 324                      | 30                                |



Al mirar el gráfico vemos que la relación entre  $x$  e  $y$  no es lineal y parece razonable intentar con una regresión polinomial. Inmediatamente surge la pregunta de elegir el grado del polinomio.

El siguiente gráfico muestra varios polinomios, con grados que van desde 1 a 8, ajustados a estos datos:



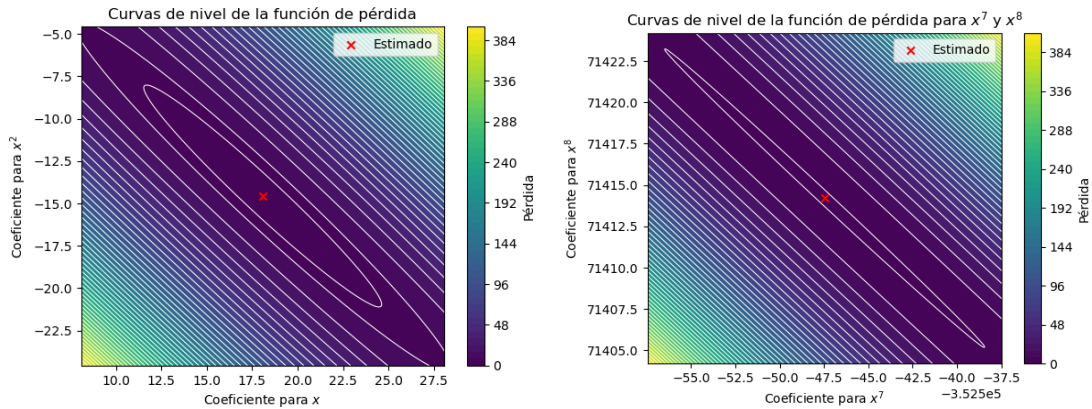
Para ajustar dichos polinomios siempre es conveniente estandarizar la matriz de diseño.

Observar la matriz de correlaciones para grado 8:

| Cor   | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ | $x^7$ | $x^8$ |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| $x^1$ | 1     | 0.977 | 0.936 | 0.897 | 0.863 | 0.834 | 0.808 | 0.784 |
| $x^2$ | 0.977 | 1     | 0.989 | 0.969 | 0.946 | 0.924 | 0.902 | 0.881 |
| $x^3$ | 0.936 | 0.989 | 1     | 0.994 | 0.982 | 0.967 | 0.951 | 0.934 |
| $x^4$ | 0.897 | 0.969 | 0.994 | 1     | 0.996 | 0.988 | 0.977 | 0.964 |
| $x^5$ | 0.863 | 0.946 | 0.982 | 0.996 | 1     | 0.998 | 0.991 | 0.982 |
| $x^6$ | 0.834 | 0.924 | 0.967 | 0.988 | 0.998 | 1     | 0.998 | 0.993 |
| $x^7$ | 0.808 | 0.902 | 0.951 | 0.977 | 0.991 | 0.998 | 1     | 0.998 |
| $x^8$ | 0.784 | 0.881 | 0.934 | 0.964 | 0.982 | 0.993 | 0.998 | 1     |

Más aún, el determinante de dicha matriz es  $5,14 \times 10^{-39}$ , y por lo tanto estamos ante la presencia de una marcada multicolinealidad, y por ende a riesgo de coeficientes grandes.

La multicolinealidad se puede visualizar al graficar las curvas de nivel de la función de pérdida. A modo de ejemplo:



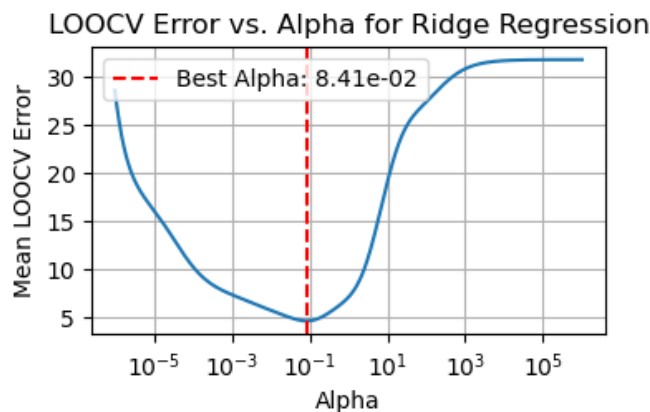
Las elipses estiradas (con los datos estandarizados) indican alta correlación entre los atributos.

Y de hecho los coeficientes para los distintos grados tienden a crecer rápidamente:

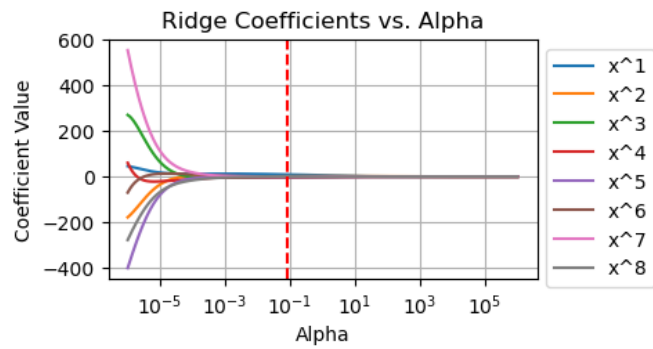
| Coefficiente | Grado 1 | Grado 2 | Grado 3 | Grado 4 | Grado 5 | Grado 6  | Grado 7    | Grado 8    |
|--------------|---------|---------|---------|---------|---------|----------|------------|------------|
| $w_1$        | 3.85    | 18.08   | 10.79   | 15.02   | 12.71   | 76.59    | -1533.75   | -5491.59   |
| $w_2$        |         | -14.56  | 3.56    | -13.02  | -0.23   | -453.11  | 13365.95   | 52990.47   |
| $w_3$        |         |         | -11.13  | 10.75   | -16.30  | 1318.39  | -50622.38  | -231805.88 |
| $w_4$        |         |         |         | -9.49   | 15.63   | -1950.95 | 102868.06  | 571182.95  |
| $w_5$        |         |         |         |         | -8.56   | 1419.35  | -116379.46 | -841240.02 |
| $w_6$        |         |         |         |         |         | -406.93  | 69088.87   | 735492.68  |
| $w_7$        |         |         |         |         |         |          | -16787.07  | -352547.47 |
| $w_8$        |         |         |         |         |         |          |            | 71414.21   |

Para controlar el tamaño de los coeficientes podemos correr una regresión polinomial de grado 8 regularizada. El problema se traslada ahora en elegir el valor de  $\alpha$ .

El gráfico a continuación muestra la curva de error para varios valores de  $\alpha$ , el error calculado usando la técnica de Leave One Out Cross Validation (LOOCV):



También podemos ver la evolución de los coeficientes en función de  $\alpha$ :



Por último, podemos ver el efecto de la regularización en los coeficientes:

| Coef  | Original     | Ridge        |
|-------|--------------|--------------|
| $x^1$ | -5491.591405 | 9.290364281  |
| $x^2$ | 52990.4687   | 1.239116016  |
| $x^3$ | -231805.8798 | -2.013422183 |
| $x^4$ | 571182.9466  | -2.683297203 |
| $x^5$ | -841240.0155 | -2.179430752 |
| $x^6$ | 735492.6833  | -1.218635966 |
| $x^7$ | -352547.4682 | -0.16433541  |
| $x^8$ | 71414.20573  | 0.79483848   |

Aquí tenemos el gráfico de la regresión regularizada:

