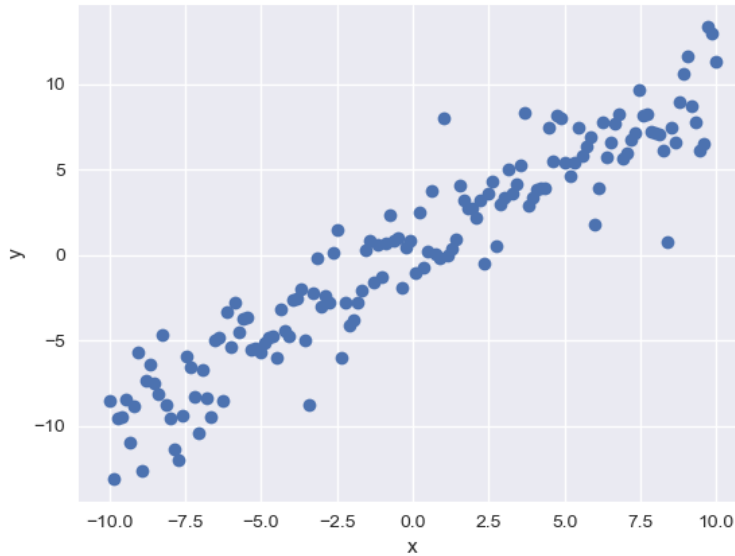


# Regresión Lineal

Carlos Malanche

13 de marzo de 2018

Ok, regresemos a lo que se tenía la clase pasada. El problema es encontrar una función estimadora  $\hat{f}$  que aproxime  $f$  dadas las parejas de datos  $\{x_i, y_i = f(x_i) + \epsilon\}_{i=1}^n$  donde, sin conocimiento previo,  $\epsilon$  es ruido de distribución Gaussiana con varianza  $\sigma^2$  y primer momento  $\mu = 0$ . Bien, pues supongamos que tomamos nuestra serie de pares de datos, y la graficamos, obteniendo lo siguiente:



Pues, a reserva de abstenciones, yo voto que es una línea con ruido, sí o no raza?

## 1. Caso de una variable

Por eso mismo sería una buena idea proponer  $\hat{f}(x) = mx + b$ . Ahora, dada la serie de datos, cómo estimamos los parámetros  $m$  y  $b$ ?

En la clase anterior vimos que la función de pérdidas cuadrática (para funciones escalares) está definida como

$$J(\hat{f}) = \frac{1}{2n} \sum_{i=1}^n (\hat{f}(x_i) - y_i)^2 \quad (1)$$

En nuestro caso ya no es necesario escribir  $J$  como un funcional, ya sabemos de qué variables va a depender el costo

$$J(m, b) = \frac{1}{2n} \sum_{i=1}^n (mx_i + b - y_i)^2 \quad (2)$$

Suena como una buena idea derivar la función de costo respecto a cada variable e igualar a cero para encontrar un *mínimo*. Primero respecto a  $m$

$$\frac{\partial J}{\partial m}(m, b) = \frac{1}{n} \sum_{i=1}^n (mx_i + b - y_i)x_i = 0 \quad (3)$$

Y ahora respecto a b

$$\frac{\partial J}{\partial b}(m, b) = \frac{1}{n} \sum_{i=1}^n (mx_i + b - y_i) = 0 \quad (4)$$

Esto es fácil de resolver, pues tenemos un sistema de dos ecuaciones con dos incógnitas, siendo la ecuación lineal en ambas variables. Abriremos la suma y factorizaremos para que sea más obvio (además de multiplicar ambas ecuaciones por n).

$$\left(\sum_{i=1}^n x_i^2\right)m + \left(\sum_{i=1}^n x_i\right)b = \sum_{i=1}^n x_i y_i \quad (5)$$

$$\left(\sum_{i=1}^n x_i\right)m + (n)b = \sum_{i=1}^n y_i \quad (6)$$

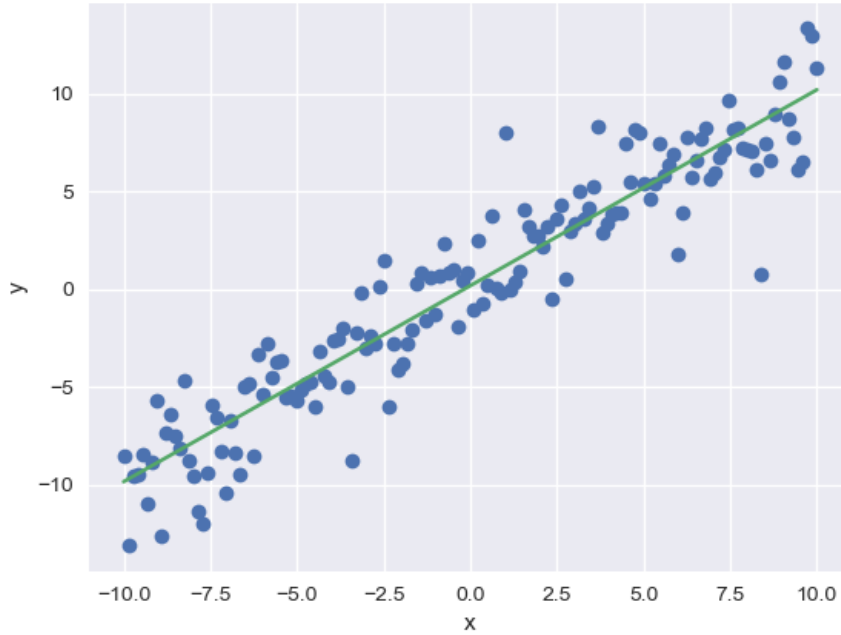
Pueden usar el método que más les guste para resolver este sistema. Por determinantes por ejemplo

$$m = \frac{n \sum_{i=1}^n x_i y_i - \left(\sum_{i=1}^n x_i\right)\left(\sum_{i=1}^n y_i\right)}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} = \frac{Cov(x, y)}{Var(x)} \quad (7)$$

Qué curioso resultado. Ahora veamos que pasa con b

$$b = \frac{\left(\sum_{i=1}^n x_i^2\right)\left(\sum_{i=1}^n y_i\right) - \left(\sum_{i=1}^n x_i\right)\left(\sum_{i=1}^n x_i y_i\right)}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} = \frac{\overline{(x \odot x)}\bar{y} - \bar{x}\overline{(x \odot y)}}{Var(x)} \quad (8)$$

Con eso, podemos construir la línea de predicciones del modelo, la cual se ve así



## 2. Caso multivariado

Antes de seguir, hay que definir las derivadas de escalares respecto a vectores

**Definición.** La derivada de una función vectorial  $\mathbf{y}(x)$  respecto a un escalar  $x$  será un vector renglón

$$\frac{\partial \mathbf{y}}{\partial x} = \left[ \frac{\partial y_1}{\partial x}, \frac{\partial y_2}{\partial x}, \dots, \frac{\partial y_m}{\partial x} \right]$$

**Definición.** La derivada de una función escalar  $y(\mathbf{x})$  respecto a un vector  $\mathbf{x}$  será un vector columna

$$\frac{\partial y}{\partial \mathbf{x}} = \left[ \frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \dots, \frac{\partial y}{\partial x_m} \right]^T$$

**Definición.** De las anteriores dos definiciones, se sigue que la derivada de una función vectorial  $\mathbf{y}(\mathbf{x})$  respecto a un vector  $\mathbf{x}$  es una matriz cuyas entradas están dadas por

$$\left( \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)_{ij} = \frac{\partial y_j}{\partial x_i}$$

Una implicación directa de esta convención (conocida como *denominator layout*, o formulación Hessiana) es que al derivar un producto de matrices, aparece una operación de transposición:

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{A}\mathbf{x}) = \frac{\partial \mathbf{x}}{\partial \mathbf{x}} \mathbf{A}^T = \mathbf{A}^T$$

Bien, lo que sigue es escribir la generalización del problema para  $\mathbf{x}$  un vector en  $\mathbb{R}^m$  en lugar de un escalar. Podemos escribir un plano en  $m$  dimensiones como el siguiente producto interior

$$\hat{f}(\mathbf{x}) = \mathbf{c} \cdot \mathbf{x} = \mathbf{c}^T \mathbf{x} \quad (9)$$

donde  $\mathbf{c}$  contiene los coeficientes a determinar. Falta decir algo sobre esta función, pero lo arreglaremos después de hacer un poco de trabajo.

Nuestra función de costo nuevamente ya tiene variables de las que depende, y podemos escribirla como

$$J(\mathbf{c}) = \frac{1}{2n} \sum_{i=1}^n (\mathbf{c}^T \mathbf{x}_i - y_i)^2 \quad (10)$$

Con un poco de imaginación, podemos ver que la suma de la función de costo es cuadrado de la norma de un vector, un vector cuyas entradas son

$$\left[ \mathbf{c}^T \mathbf{x}_1 - y_1, \mathbf{c}^T \mathbf{x}_2 - y_2, \dots, \mathbf{c}^T \mathbf{x}_n - y_n \right]$$

Pues, si inteligentemente definimos la matrix  $\mathbf{X}$  como

$$\mathbf{X} = \left[ \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \right]^T$$

entonces, definiendo  $\mathbf{y} := [y_1, y_2, \dots, y_n]^T$  tenemos que la función de costo se puede escribir como

$$J(\mathbf{c}) = \frac{1}{2n} \|\mathbf{X}\mathbf{c} - \mathbf{y}\|^2 = \frac{1}{2n} (\mathbf{X}\mathbf{c} - \mathbf{y})^T (\mathbf{X}\mathbf{c} - \mathbf{y}) \quad (11)$$

Vamos ahora a buscar el mínimo de esta función, encontrando el gradiente  $\partial J_{\mathbf{c}}$  e igualándolo a cero. Para esto necesitamos la regla de la cadena. Noten porfavor que por la formulación Hessiana, se tiene que para una función  $u : \mathbb{R}^n \rightarrow \mathbb{R}$ , una función  $\mathbf{g} : \mathbb{R}^m \rightarrow \mathbb{R}^n$  y una variable  $\mathbf{c} \in \mathbb{R}^m$ , la regla de la cadena es:

$$\frac{\partial}{\partial \mathbf{c}} u(\mathbf{g}(\mathbf{c})) = \frac{\partial \mathbf{g}}{\partial \mathbf{x}} \frac{\partial u}{\partial \mathbf{g}} \quad (12)$$

Ahora basta con aplicar la regla de la cadena para encontrar el gradiente de  $J$

$$\frac{\partial J}{\partial \mathbf{c}} = \frac{1}{n} \mathbf{X}^T (\mathbf{X}\mathbf{c} - \mathbf{y}) \quad (13)$$

Al igualar esta expresión a cero se puede resolver para  $\mathbf{c}$

$$\mathbf{c} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (14)$$

Al conjunto de ecuaciones descritas por la ecuación matricial (14) se les conoce como **ecuaciones normales de la regresión lineal**.

## 2.1. Interpretación geométrica de las ecuaciones normales

Desde haber escrito la función de costo en forma matricial, el problema se convirtió en minimizar el tamaño del vector  $\mathbf{X}\mathbf{c} - \mathbf{y}$ , es decir, encontrar un vector ortogonal al espacio generado por las columnas de  $\mathbf{X}$ , cuya longitud está dada por la longitud de  $\mathbf{y}$  y su ángulo a dicho espacio. Cada vector columna de la matriz  $\mathbf{X}$  es un vector correspondiente a una de las variables que estamos midiendo, y todas sus entradas son las múltiples mediciones que hicimos.

## 2.2. La matriz de Gram

Como una curiosa observación, pongámonle un poco más de atención a la matriz  $\mathbf{X}^T\mathbf{X}$ , la cual es conocida como *matriz de Gram*. Esta matriz puede dar problemas si resulta no ser invertible, por lo que diremos lo siguiente

**Teorema.** *La matriz de Gram es invertible si y sólo si  $\mathbf{X}$  es de rango completo.*

*Demostración.* ( $\Leftarrow$ ) Suponga que  $\mathbf{X}$  no es de rango completo, entonces existe un vector  $\mathbf{v} \neq \mathbf{0}$  tal que  $\mathbf{X}\mathbf{v} = \mathbf{0}$ , con lo que se tiene que  $\mathbf{X}^T\mathbf{X}\mathbf{v} = \mathbf{0}$ , lo que implica que  $\mathbf{X}^T\mathbf{X}$  tampoco es de rango completo y por lo tanto no es invertible.

( $\rightarrow$ ) Suponga que  $\mathbf{X}^T\mathbf{X}$  no es invertible, siendo ese el caso existe un vector  $\mathbf{v} \neq \mathbf{0}$  tal que  $\mathbf{X}^T\mathbf{X}\mathbf{v} = \mathbf{0}$ . Se sigue entonces que

$$\mathbf{v}^T\mathbf{X}^T\mathbf{X}\mathbf{v} = \|\mathbf{X}\mathbf{v}\|^2 = 0$$

Esto implica que hay dependencia lineal en las columnas de  $\mathbf{X}$  y por tanto es de rango deficiente □

El resultado de esto es que las columnas de la matriz  $\mathbf{X}$  deben ser linealmente independientes. Ojo que las columnas de la matriz no son los vectores de cada medición; hay un vector por cada variable en la que se hacen mediciones. Independencia lineal de estos vectores quiere decir que necesitamos que ni una de las variables que medimos sea la combinación lineal de otras. Aquí viene la importancia de la correlación! Y aunque la correlación entre columnas no sea 1, columnas con alta correlación volverán computacionalmente inestable a la inversión de la matriz de Gram. Cuidado!

Por último, esta matriz resulta ser positiva semidefinida. La prueba la pueden hacer ustedes.

## 2.3. Una manera de arreglar la inestabilidad

Cuando la inversión de la matriz de Gram es inestable, a veces se recurre a añadir términos de regularización a la función de costo. Es para que se queden *picados*, viene en las siguientes clases. Ah, también nunca verificamos que el problema es *convexo* y de minimización. Eso queda pendiente!