

MAT-468: Sesión 4, Cálculos en regresión II

Felipe Osorio

<http://fosorios.mat.utfsm.cl>

Departamento de Matemática, UTFSM



Método gradientes conjugados (GC) en regresión lineal

En el contexto de regresión lineal, considere:

$$\phi(\beta) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 = \frac{1}{2} (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta).$$

El objetivo del **procedimiento GC**¹ es producir la secuencia:

$$\beta^{(k+1)} = \beta^{(k)} + \lambda_k \mathbf{p}_k, \quad k = 0, 1, \dots \quad (1)$$

El algoritmo básico considera:

$$\lambda_k = \frac{\mathbf{p}_k^\top \mathbf{g}_k}{\mathbf{p}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{p}_k}, \quad \mathbf{g}_k = \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\beta^{(k)}).$$

(En efecto, $\partial\phi(\beta)/\partial\beta = -\mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\beta)$) y actualizamos la dirección de búsqueda como:

$$\mathbf{p}_{k+1} = \mathbf{g}_{k+1} + \delta_k \mathbf{p}_k, \quad \delta_k = -\frac{\mathbf{g}_{k+1}^\top \mathbf{p}_k}{\mathbf{p}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{p}_k}.$$

¹McIntosh (1982), Lecture Notes in Statistics 10.



Método gradientes conjugados (GC) en regresión lineal

Se ha sugerido usar:

$$\lambda_k = \frac{\mathbf{p}_k^\top \mathbf{X}^\top \mathbf{y}}{\mathbf{p}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{p}_k},$$

y actualizar

$$\mathbf{p}_{k+1} = \mathbf{g}_{k+1} + \delta_{k+1} \mathbf{p}_k, \quad \delta_{k+1} = -\frac{\mathbf{p}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{g}_k}{\mathbf{p}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{p}_k}.$$

Para hacer el proceso más simple es recomendable calcular

$$\mathbf{h}_k = \mathbf{X}^\top \mathbf{X} \mathbf{p}_k.$$

De este modo el **requerimiento de almacenamiento** del algoritmo es sólo $4p$.



Gradientes conjugados en regresión lineal

Algoritmo 1: Gradientes conjugados para regresión lineal.

Entrada : Datos X y y

Parámetros: Tolerancia τ .

```
1 begin
2   Hacer  $\beta = 0$ ,  $p = g = -X^T y$ ,  $\delta = 0$  y  $\gamma = \|g\|^2$ 
3   while  $\gamma > \tau$  do
4     Calcular  $h = X^T X p$  y  $u = p^T X^T X p = p^T h$ 
5     if  $k \neq 1$  then
6        $\delta = -h^T g / u$ 
7        $p = g + \delta p$ 
8        $\lambda = -p^T g / u$ 
9        $\beta = \beta + \lambda p$ 
10       $g = g + \lambda h$ 
11    end
12    return  $\hat{\beta} = \beta$ 
13 end
```



- ▶ Soluciones regularizadas: [Regresión ridge](#).
- ▶ Estimación vía [IRLS](#):
 - ▶ Modelos lineales generalizados.
 - ▶ Estimación L_1 .
 - ▶ Estimación M .
 - ▶ Regresión lineal considerando [distribuciones con colas pesadas](#).



Considere el modelo de regresión lineal

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

con $\mathbf{X} \in \mathbb{R}^{n \times p}$, $E(\boldsymbol{\epsilon}) = \mathbf{0}$ y $\text{Cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}_n$.

Es bien conocido que cuando \mathbf{X} es **mal-condicionada**, el sistema de ecuaciones

$$\mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X}^\top \mathbf{y},$$

puede ser **inestable** (ver, Stewart, 1987 y Belsley, 1991).



Considere la **descomposición valor singular (SVD)** de \mathbf{X} ,

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top,$$

donde $\mathbf{U} \in \mathbb{R}^{n \times r}$, tal que $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_r$, $\mathbf{D} = \text{diag}(d_1, \dots, d_r)$ con $d_1 \geq \dots \geq d_r > 0$, $\mathbf{V} \in \mathbb{R}^{r \times r}$ es matriz ortogonal y $r = \text{rg}(\mathbf{X})$.

La detección de colinealidad en el modelo lineal puede ser llevada a cabo por medio del **número condición**

$$\kappa(\mathbf{X}) = \|\mathbf{X}\| \|\mathbf{X}^+\| = \frac{d_1}{d_r},$$

y $\kappa(\mathbf{X})$ “grande” es indicador de colinealidad.



Número condición

Considere la matriz

$$\mathbf{A} = \begin{pmatrix} 1.000 & 0.500 \\ 0.667 & 0.333 \end{pmatrix}, \quad \mathbf{A}^{-1} = \begin{pmatrix} -666 & 1000 \\ 1344 & -2000 \end{pmatrix}.$$

El **número condición** se define como $\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$ para $\|\cdot\|$ alguna norma matricial.

Por ejemplo²,

$$\kappa_1(\mathbf{A}) = \|\mathbf{A}\|_1 \|\mathbf{A}^{-1}\|_1 = (1.667)(3000) = 5001$$

$$\kappa_\infty(\mathbf{A}) = \|\mathbf{A}\|_\infty \|\mathbf{A}^{-1}\|_\infty = (1.500)(3344) = 5016$$

$$\kappa_2(\mathbf{A}) = \frac{\max_{\mathbf{x} \neq 0} \|\mathbf{A}\mathbf{x}\| / \|\mathbf{x}\|}{\min_{\mathbf{x} \neq 0} \|\mathbf{A}\mathbf{x}\| / \|\mathbf{x}\|} = \left| \frac{\lambda_{\max}}{\lambda_{\min}} \right| = \frac{1.333375}{0.000375} = 3555.778$$

² $\|\mathbf{A}\|_p = \max_{\|\mathbf{x}\|_p=1} \|\mathbf{A}\mathbf{x}\|_p$, $\|\mathbf{A}\|_2 = \sqrt{\rho(\mathbf{A}^\top \mathbf{A})}$



Cemento Portland (Woods, Steinour y Starke, 1932)

Estudio experimental relacionando la emisión de calor durante la producción y endurecimiento de 13 muestras de cementos Portland. Woods et al. (1932) consideraron cuatro compuestos para los clinkers desde los que se produce el cemento.

La respuesta (Y) es la **emisión de calor** después de 180 días de curado, medido en calorías por gramo de cemento. Los regresores son los porcentajes de los cuatro compuestos: **aluminato tricálcico** (X_1), **silicato tricálcico** (X_2), **ferrito aluminato tetracálcico** (X_3) y **silicato dicálcico** (X_4).



Cemento Portland (Woods, Steinour y Starke, 1932)

Siguiendo a Woods et al. (1932) consideramos un modelo lineal sin intercepto (**modelo homogéneo**). El número condición escalado es $\kappa(\mathbf{X}) = 9.432$, esto es \mathbf{X} es **bien condicionada**. (variables centradas, $\kappa(\tilde{\mathbf{X}}) = 37.106$)

Por otro lado, Hald (1952), Gorman y Toman (1966) y Daniel y Wood (1980) adoptaron un modelo con intercepto (**modelo no homogéneo**). En cuyo caso, $\kappa(\mathbf{X}) = 249.578$, sugiriendo la presencia de **colinealidad**.

El aumento en el número condición se debe a que existe una **relación lineal aproximada**, pues

$$x_1 + x_2 + x_3 + x_4 \approx 100.$$

de modo que incluir el intercepto causa una **colinealidad severa**.



El **estimador ridge** (Hoerl y Kennard, 1970)³,

$$\hat{\beta}_\lambda = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}, \quad \lambda > 0.$$

puede ser visto como la solución del **problema regularizado**:

$$\min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 - \frac{\lambda}{2} \|\beta\|^2,$$

o bien como el problema **mínimos cuadrados con datos aumentados**:

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{X} \\ \sqrt{\lambda} \mathbf{I}_p \end{pmatrix} \beta + \begin{pmatrix} \epsilon \\ \epsilon_* \end{pmatrix}.$$

En este contexto λ es un **parámetro de regularización** (parámetro ridge).

³Technometrics 12, 55-67.



El mejor método para obtener $\hat{\beta}_\lambda$ es usar la descomposición SVD

$$\hat{\beta}_\lambda = \mathbf{V}\hat{\alpha}_\lambda,$$

con

$$\hat{\alpha}_\lambda = (\mathbf{D}^2 + \lambda \mathbf{I}_p)^{-1} \mathbf{D} \mathbf{z} = \begin{pmatrix} z_1 d_1 / (d_1^2 + \lambda) \\ \vdots \\ z_p d_p / (d_p^2 + \lambda) \end{pmatrix},$$

donde $\mathbf{z} = \mathbf{U}^\top \mathbf{y}$. Un procedimiento recomendado para seleccionar el **parámetro ridge** es **validación cruzada generalizada** (Golub, Heath y Wahba, 1979), definido como:

$$\text{GCV}(\lambda) = \frac{1}{n} \frac{\|\mathbf{y} - \mathbf{X}\hat{\beta}_\lambda\|^2}{\{\text{tr}(\mathbf{I}_n - \mathbf{H}(\lambda))\}^2},$$

con $\mathbf{H}(\lambda) = \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top$ y definimos $\text{edf} = \text{tr} \mathbf{H}(\lambda)$.



Es fácil mostrar que:

$$\text{tr } \mathbf{H}(\lambda) = \text{tr } \mathbf{D}^2(\mathbf{D}^2 + \lambda \mathbf{I}_p)^{-1} = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}.$$

Así,

$$\text{GCV}(\lambda) = \frac{\|\mathbf{z} - \mathbf{D}\hat{\boldsymbol{\alpha}}_\lambda\|^2/n}{(1 - \text{edf}/n)^2}.$$

Esto permite evaluar la función $\text{GCV}(\lambda)$ de forma simple. Para escoger un λ_{opt} se ha sugerido:

- ▶ Considerar una **grilla de valores** para λ .
- ▶ **Optimizar $\text{GCV}(\lambda)$** usando un procedimiento de minimización unidimensional.



Resultados de estimación:

Estimador	β_0	β_1	β_2	β_3	β_4	σ^2
LS (homogéneo)		2.193	1.153	0.759	0.486	4.047
LS (No homog.)	62.405	1.551	0.510	0.102	-0.144	3.682
Ridge, LW	17.189	2.016	0.976	0.578	0.313	3.874
Ridge, HKB	8.587	2.105	1.065	0.668	0.400	3.953
Ridge, GCV	0.085	2.165	1.159	0.738	0.489	4.055

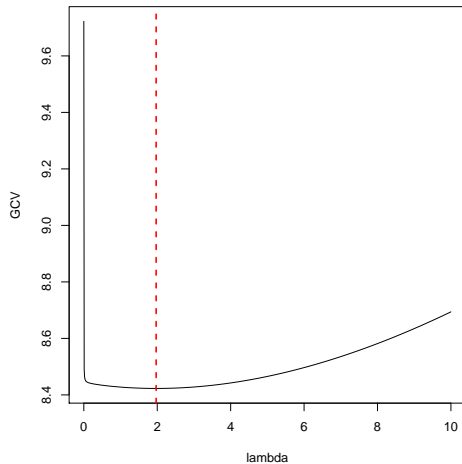
Se utilizó

$$\hat{\lambda}_{\text{HKB}} = ps^2 / \|\hat{\beta}_{\text{LS}}\|^2 = 0.00767 \quad (\text{Hoerl, Kennard y Baldwin, 1975}).$$

$$\hat{\lambda}_{\text{LW}} = ps^2 / \|\mathbf{X}\hat{\beta}_{\text{LS}}\|^2 = 0.00321 \quad (\text{Lawless y Wang, 1976}).$$

Además, se consideró una grilla de valores para $\lambda = 0.00, 0.01, 0.02, \dots, 10.00$, obteniendo $\hat{\lambda}_{\text{opt}} = 1.97$.





Selección del parámetro ridge usando GCV, $\lambda_{\text{opt}} = 1.97$.



Uno de los primeros **procedimientos robustos en regresión**⁴ corresponde al problema:

$$\min_{\beta} \sum_{i=1}^n |Y_i - \mathbf{x}_i^T \beta|.$$

Mínimo desvío absoluto (LAD) o **regresión L_1** puede ser planteado como un problema de **programación lineal**, considerando las partes positivas y negativas de los residuos, e^+ y e^- , respectivamente, y análogamente para β^+ , β^- .

Así, el problema puede ser expresado como (Charnes, Cooper y Ferguson, 1955):

$$\begin{aligned} \min_{\beta} \mathbf{1}^T (e^+ + e^-), \\ \text{sujeto a: } \mathbf{Y} = \mathbf{X}(\beta^+ - \beta^-) + (e^+ - e^-), \end{aligned}$$

con β^+ , β^- , e^+ , e^- deben ser todos ≥ 0 .

Observación: Barrodale y Roberts (1973, 1974) presentan un algoritmo de propósito especial para resolver este problema modificando el **método simplex** y la **estructura de datos** requerida.

⁴Este método es, de hecho, **anterior a LS!**



Uno de los primeros **procedimientos robustos en regresión**⁴ corresponde al problema:

$$\min_{\beta} \sum_{i=1}^n |Y_i - \mathbf{x}_i^T \beta|.$$

Mínimo desvío absoluto (LAD) o **regresión L_1** puede ser planteado como un problema de **programación lineal**, considerando las partes positivas y negativas de los residuos, e^+ y e^- , respectivamente, y análogamente para β^+ , β^- .

Así, el problema puede ser expresado como (Charnes, Cooper y Ferguson, 1955):

$$\begin{aligned} \min_{\beta} \mathbf{1}^T (e^+ + e^-), \\ \text{sujeto a: } \mathbf{Y} = \mathbf{X}(\beta^+ - \beta^-) + (e^+ - e^-), \end{aligned}$$

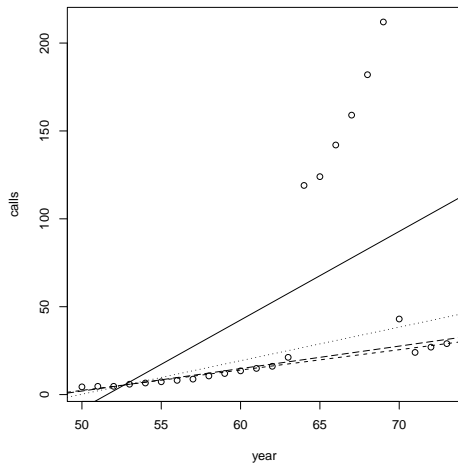
con β^+ , β^- , e^+ , e^- deben ser todos ≥ 0 .

Observación: Barrodale y Roberts (1973, 1974) presentan un algoritmo de propósito especial para resolver este problema modificando el **método simplex** y la **estructura de datos** requerida.

⁴Este método es, de hecho, **anterior a LS!**



Llamadas telefónicas en Bélgica 1950-73 (Rousseuw y Leroy, 1987)



Ajustes: LS, normal contaminada ($\epsilon = .15, \gamma = .10$), Student- t ($\nu = 2.5$), L_1 .



- ▶ Schlossmacher (1973) originalmente propuso calcular **estimadores LAD** usando IRLS.
- ▶ Lange y Sinsheimer (1993) y Phillips (2002) identificaron que este procedimiento IRLS corresponde a un **algoritmo EM**.
- ▶ Sin embargo, también ha sido reportado que este algoritmo puede ser **incapaz de detectar la observaciones básicas** de manera eficiente.



Regresión lineal: función LAD de L1pack

Considere el modelo

$$Y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + (\sqrt{2}\tau_i)^{-1}\epsilon_i, \quad i = 1, \dots, n$$

donde $\epsilon \stackrel{\text{ind}}{\sim} \text{N}(0, \phi)$ y τ_i tiene función de densidad

$$g(\tau_i) = \tau_i^{-3} \exp(-\frac{1}{2}\tau_i^{-2}).$$

El algoritmo EM procede a llevar a cabo la estimación de $\boldsymbol{\beta}$ y ϕ iterativamente mediante maximizar la función:

$$\begin{aligned} Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(k)}) &= -\frac{n}{2} \log \phi - \frac{1}{2\phi} \sum_{i=1}^n W_i^{(k)} (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \\ &= -\frac{n}{2} \log \phi - \frac{1}{2\phi} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top \mathbf{W}^{(k)} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \end{aligned}$$

donde $\mathbf{W}^{(k)} = \text{diag}(W_1^{(k)}, \dots, W_n^{(k)})$ y los pesos son dados por:

$$W_i^{(k)} = E(\tau_i^2 | Y_i, \boldsymbol{\theta}^{(k)}) = \sigma^{(k)} / \sqrt{2} |Y_i - \mathbf{x}_i^\top \boldsymbol{\beta}^{(k)}|,$$

para $|Y_i - \mathbf{x}_i^\top \boldsymbol{\beta}^{(k)}| > 0$.



Considere el modelo

$$Y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + (\sqrt{2}\tau_i)^{-1} \epsilon_i, \quad i = 1, \dots, n$$

donde $\epsilon \stackrel{\text{ind}}{\sim} N(0, \phi)$ y τ_i tiene función de densidad

$$g(\tau_i) = \tau_i^{-3} \exp(-\frac{1}{2}\tau_i^{-2}).$$

El **algoritmo EM** procede a llevar a cabo la estimación de $\boldsymbol{\beta}$ y ϕ iterativamente mediante maximizar la función:

$$\begin{aligned} Q(\boldsymbol{\theta} | \boldsymbol{\theta}^{(k)}) &= -\frac{n}{2} \log \phi - \frac{1}{2\phi} \sum_{i=1}^n W_i^{(k)} (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \\ &= -\frac{n}{2} \log \phi - \frac{1}{2\phi} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top \mathbf{W}^{(k)} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \end{aligned}$$

donde $\mathbf{W}^{(k)} = \text{diag}(W_1^{(k)}, \dots, W_n^{(k)})$ y los pesos son dados por:

$$W_i^{(k)} = E(\tau_i^2 | Y_i, \boldsymbol{\theta}^{(k)}) = \sigma^{(k)} / \sqrt{2} |Y_i - \mathbf{x}_i^\top \boldsymbol{\beta}^{(k)}|,$$

para $|Y_i - \mathbf{x}_i^\top \boldsymbol{\beta}^{(k)}| > 0$.



Paso de coeficientes:

- ▶ Calcular $\mathbf{r}^{(k)} = \mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^{(k)}$ y $\mathbf{W}^{(k)1/2} = \text{diag}(W_1^{(k)}, \dots, W_n^{(k)})$
- ▶ Obtener $\boldsymbol{\delta}^{(k)}$ como solución del problema WLS

$$\min_{\boldsymbol{\delta}} \|\mathbf{W}^{(k)1/2}(\mathbf{r}^{(k)} - \mathbf{X}\boldsymbol{\delta})\|^2$$

- ▶ Actualizar $\boldsymbol{\beta}^{(k+1)} = \boldsymbol{\beta}^{(k)} + \boldsymbol{\delta}^{(k)}$.

Paso de escala:

$$\phi^{(k+1)} = \frac{1}{n} \|\mathbf{W}^{(k)1/2} \mathbf{r}^{(k+1)}\|^2$$

Criterio de convergencia: basado en el criterio usado en la función `glm.fit()`.



Detalles de la implementación

Sea $X_* = \widehat{W}^{1/2} X$, $Y_* = \widehat{W}^{1/2} Y$ y calcular la descomposición QR de X_*

$$\text{(DGEQRF)} \quad X_* = Q \begin{pmatrix} R \\ \mathbf{0} \end{pmatrix}, \quad Q \in \mathcal{O}_n \text{ y } R \in \mathbb{R}^{p \times p} \text{ triangular superior,}$$

considere $c = Q^\top Y_*$, entonces

$$\text{(DORMQR)} \quad Q^\top e_* = Q^\top \widehat{W}^{1/2} (Y - X\beta) = \begin{pmatrix} c_1 - R\beta \\ c_2 \end{pmatrix} = \begin{pmatrix} r_1 \\ r_2 \end{pmatrix}$$

de este modo, δ es solución del sistema triangular

$$\text{(DTRTRS)} \quad R\delta = r_1 \quad \Rightarrow \quad R\beta^{(k+1)} = c_1,$$

actualizar $\beta^{(k+1)} = \beta^{(k)} + \delta$ (DAXPY) y $\phi^{(k+1)} = \|r_2\|^2/n$ (DNRM2). Finalmente, note que

$$\text{(DORMQR)} \quad \widehat{Y}_* = X_* \beta^{(k+1)} = Q \begin{pmatrix} R\beta^{(k+1)} \\ \mathbf{0} \end{pmatrix} = Q \begin{pmatrix} c_1 \\ \mathbf{0} \end{pmatrix}$$



Referencias bibliográficas



Golub, G.H., Heath, M., Wahba, G. (1979).

Generalized cross-validation as a method for choosing a good ridge parameter.
Technometrics 21, 215-223.



Hoerl, A.E., Kennard, R.W., Baldwin, K.F. (1975).

Ridge regression: some simulations.
Communications in Statistics 4, 105-123.



Lange, K., Sinsheimer, J.S. (1993).

Normal/independent distributions and their applications in robust regression.
Journal of Computational and Graphical Statistics 2, 175-198.



Lawless, J.F., Wang, P. (1976).

A simulation study of ridge and other regression estimators.
Communications in Statistics – Theory and Methods 14, 1589-1604.



Phillips, R.F. (2002).

Least absolute deviations estimation via the EM algorithm.
Statistics and Computing 12, 281-285.

