

**Análisis aplicado.**  
**Gradiente conjugado y proyecto.**  
10 de marzo, 2015. Guillermo Santiago Novoa Pérez

## Gradiente conjugado

El método de gradiente conjugado (GC) tiene propiedades matemáticas y computacionales que lo convierten en una herramienta muy versátil para resolver problemas en una amplia variedad de aplicaciones.

La principal motivación consiste en diseñar un método para resolver eficientemente el problema cuadrático

$$\underset{x}{\text{minimizar}} \quad q(x) = \frac{1}{2}x^T A x - b^T x, \quad (1)$$

en donde  $A$  es una matriz simétrica positiva definida de  $n \times n$ . El problema anterior es equivalente a resolver el sistema de ecuaciones lineales  $Ax = b$ . Una de las principales virtudes del método de GC es que alivia los conocidos inconvenientes del método de máximo descenso. En aritmética finita, GC parte de una aproximación inicial y termina con la solución exacta en no más de  $n$  iteraciones; en este sentido se le considera como un método directo. En aritmética finita, debido a la acumulación de errores por redondeo, el método podría no terminar en  $n$  pasos y por lo tanto se le considera un método iterativo.

Los métodos iterativos constituyen una alternativa muy atractiva con respecto a los métodos directos. Específicamente, cuando se trata de resolver una sucesión finita de sistemas cercanamente relacionados

$$A_j x = b_j, \quad j = 1, \quad A_{j+1} \approx A_j, \quad j = 1, \dots, N-1 \quad (2)$$

$x_j^*$ , la solución del sistema  $A_j x = b_j$ , se puede utilizar como punto inicial para resolver el sistema  $A_{j+1} x = b_{j+1}$ . En estos casos la solución de cada sistema se alcanza en unas cuantas iteraciones.

Si cada sistema de la sucesión (2) se resolviera con el método de Cholesky se requerirían  $\mathcal{O}(Nn^3)$  operaciones para resolver los  $N$  sistemas.

Antes de iniciar la construcción de GC necesitamos algunas definiciones.

1. Al vector  $r(x) = Ax - b$  lo denotamos como el *residuo* en  $x$ .<sup>1</sup>

2. Los elementos de un conjunto  $\mathcal{D} = \{d_0, d_1, \dots, d_k\}$  de vectores en  $\mathbb{R}^n$  son  $A$ -ortogonales si

$$d_i^T A d_j = 0, \quad i \neq j$$

Consideremos al vector  $x_0$  una aproximación inicial a la solución  $x^*$ . Entonces es inmediato observar que  $r(x_0)$  es una dirección de descenso para el problema de optimización (1); por lo tanto

---

<sup>1</sup>Notar que  $\nabla q(x) = Ax - b = r(x)$ . En particular, en  $x = x_k$ , tendremos  $r_k = Ax_k - b$ .

podemos obtener el iterando

$$x_1 = x_0 + \alpha_0 d_0,$$

en donde  $d_0 = -r_0$ ; el escalar  $\alpha_0$  se puede obtener resolviendo el problema unidimensional

$$\underset{\alpha}{\text{minimizar}} \quad q(x_0 + \alpha d_0),$$

es decir

$$\nabla q(x_0 + \alpha_0 d_0)^T d_0 = r_1^T d_0 = 0, \quad r_1^T r_0 = 0$$

y por lo tanto

$$\alpha_0 = \frac{-\nabla q(x_0)^T d_0}{d_0^T A d_0} = -\frac{r_0^T d_0}{d_0^T A d_0}$$

El iterando  $x_2$  se calcula con la dirección

$$d_1 = -r_1 + \beta_1 d_0, \quad \beta_1 = \frac{r_1^T A d_0}{d_0^T A d_0} \quad \text{II}$$

obtenida a partir de  $-r_1$  y  $d_0$  mediante  $A$ -ortogonalización de Gram-Schmidt.

$$x_2 = x_1 + \alpha_1 d_1,$$

en donde  $\alpha_1$  se obtiene resolviendo

$$\underset{\alpha}{\text{minimizar}} \quad q(x_1 + \alpha d_1),$$

es decir

$$\nabla q(x_1 + \alpha_1 d_1)^T d_1 = r_2^T d_1 = 0$$

y por lo tanto

$$\alpha_1 = \frac{-\nabla q(x_1)^T d_1}{d_1^T A d_1} = -\frac{r_1^T d_1}{d_1^T A d_1}$$

Vamos a probar que  $r_2$  también es ortogonal a  $d_0$ ; de la definición del residuo

$$r_2 = Ax_2 - b = A(x_1 + \alpha_1 d_1) - b = r_1 + \alpha_1 A d_1$$

---

<sup>II</sup>Observar que  $d_1$  está compuesta por la dirección de máximo descenso y una fracción de la dirección previa,  $d_0$ .

premultiplicando por  $d_0$  tenemos

$$d_0^T r_2 = d_0^T r_1 + \alpha_1 d_0^T A d_1 = r_0^T r_1 + \alpha_1 d_0^T A d_1 = 0.$$

Como las direcciones  $d_i$  son construidas mediante  $A$ -ortogonalización de Gram-Schmidt,  $d_0^T A d_1 = 0$ , además  $r_0^T r_1 = 0$ . Por lo tanto  $r_2$  es ortogonal a  $d_0$  también. Nos queda por investigar el producto  $r_2^T r_1$ :

$$r_2^T d_1 = 0 = r_2^T (-r_1 + \beta_1 d_0) = -r_2^T r_1 + r_2^T d_0$$

sabemos que  $r_2^T d_1 = 0$  y que  $r_2^T d_0 = 0$ , por lo tanto  $r_2^T r_1 = 0$ .

En conclusión tenemos que el conjunto  $\{r_0, r_1, r_2\}$  es ortogonal y por lo tanto podemos construir un vector  $A$ -ortogonal

$$d_2 = -r_2 + \frac{r_2^T A d_0}{d_0^T A d_0} d_0 + \beta_2 d_1$$

Utilizando la expresión

$$x_1 = x_0 + \alpha_0 d_0$$

es fácil probar que el producto  $r_2^T A d_0$  se anula. Premultiplicando la expresión anterior por  $A$  tenemos

$$\begin{aligned} x_1 &= x_0 + \alpha_0 d_0 \\ \implies A x_1 &= A x_0 + \alpha_0 A d_0 \\ \implies (A x_1 - b) - (A x_0 - b) &= \alpha_0 A d_0 \\ \implies r_1 - r_0 &= \alpha_0 A d_0 \quad * \\ \implies r_2^T (r_1 - r_0) &= \alpha_0 r_2^T A d_0 \\ \implies r_2^T A d_0 &= 0 \end{aligned}$$

\*III

Por lo tanto

$$d_2 = -r_2 + \beta_2 d_1, \quad \beta_2 = \frac{r_2^T A d_1}{d_1^T A d_1}.$$

En resumen, tenemos que los residuos  $r_0, r_1, r_2$  forman un conjunto ortogonal, y que las direcciones  $d_0, d_1, d_2$  forman un conjunto  $A$ -ortogonal. Un argumento inductivo evidente permite formular una versión preliminar del algoritmo de GC

---

<sup>III</sup>**Nota :** De este paso se obtiene la forma recursiva para  $r_i$

## Proyecto

Se empezará el proyecto terminando de probar el paso inductivo en la construcción del método de GC. Para ello, se tiene que demostrar que para toda dirección  $d_k^{IV}$  del método del GC, su obtención depende únicamente de la dirección anterior. Será importante notar que para probar que  $d_{k+1}$  solamente depende de la dirección anterior ( $d_k$ ), sólo hace falta probar que el  $k$ -ésimo residuo ( $r_k$ ) es  $A$ -ortogonal a toda dirección que no sea la inmediatamente anterior ( $d_k$ ).

En otras palabras :

$$r_i \perp_A d_j \quad \forall j \in \{0, 1, \dots, i-2\}$$

Nuestro **caso base** ya está dado anteriormente <sup>v</sup>, supongamos que para  $k$  se cumple lo siguiente:

Hipótesis de inducción

$$r_k \perp d_i \quad \forall i \in \{0, 1, \dots, k-1\} \quad (3)$$

$$r_k \perp r_i \quad \forall i \in \{0, 1, \dots, k-1\} \quad (4)$$

$$r_k \perp_A d_i \quad \forall i \in \{0, 1, \dots, k-2\} \quad (5)$$

Se separará a la prueba en 3 partes:

1. PD:  $r_{k+1}^T d_i = 0, \quad \forall i \in \{0, 1, \dots, k\}$
2. PD:  $r_{k+1}^T r_i = 0, \quad \forall i \in \{0, 1, \dots, k\}$
3. PD:  $r_{k+1}^T A d_i = 0, \quad \forall i \in \{0, 1, \dots, k-1\}$

De la definición de residuo se tiene que:

$$\begin{aligned} r_{k+1} &= A x_{k+1} - b \\ &= A(x_k + \alpha_k d_k) - b \\ &= r_k + \alpha_k A d_k \end{aligned}$$

Premultiplicando por  $d_i^T \quad i \in \{0, 1, \dots, k-1\}$

$$\implies d_i^T r_{k+1} = d_i^T r_k - \alpha_k d_i^T A d_k$$

como, por construcción,  $d_i \perp_A d_j \quad \forall i \neq j$

y, por (3)

$$\alpha_k d_i^T A d_k = 0 \quad d_i^T r_k = 0 \quad \forall i \in \{0, 1, \dots, k-1\}$$

$$\implies d_i^T r_{k+1} = 0;$$

---

<sup>IV</sup>la dirección es obtenida por medio de la  $A$ -ortogonalización de los negativos de los residuos ( $-r_i$ ) utilizando el método de Gram-Schmidt, es decir,  $\beta_i$  es la  $A$ -proyección de  $-r_{i+1}$  sobre  $d_i$

<sup>v</sup> $r_0, r_1, r_2, d_0, d_1, d_2$

Es decir,  $r_{k+1}$  es ortogonal a toda dirección anterior a la  $k$ -ésima dirección.

Para probar que  $r_{k+1}$  es ortogonal a la  $k$ -ésima dirección, sólo hace falta fijarse en la forma en la que se obtiene  $\alpha_k$ . Como  $x_{k+1} = x_k + \alpha_k d_k$  con  $\alpha_k$  obtenida mediante búsqueda lineal exacta, se tiene que:

$$\begin{aligned}\nabla f(x_k + \alpha_k d_k)^T d_k &= \nabla f(x_{k+1})^T d_k = r_{k+1}^T d_k = 0 = d_k^T r_{k+1} \\ \therefore d_i^T r_{k+1} &= 0 \quad \forall i \in \{0, 1, \dots, k\} \quad \square_1\end{aligned}$$

(Con esto se acaba de dar una prueba de que  $r_{k+1} \perp d_i \forall i \in \{0, 1, \dots, k\}$ )

Para la segunda parte de la prueba recordemos la forma que tienen las direcciones anteriores:

$$d_i = -r_i + \beta_i d_{i-1} \quad \forall i \in \{1, 2, \dots, k\}$$

como sabemos que  $d_i^T r_{k+1} = 0$ , sustituycamos  $d_i$  por la expresión anterior :

$$\begin{aligned}0 &= r_{k+1}^T d_i \stackrel{(5)}{=} r_{k+1}^T (-r_i + \beta_i d_{i-1}) &= r_{k+1}^T r_i + \beta_i r_{k+1}^T d_{i-1} \\ \text{por el resultado anterior } r_{k+1}^T d_i &= 0 &\implies \beta_i r_{k+1}^T d_{i-1} = 0 \\ &\therefore r_{k+1}^T r_i &= 0 \quad \forall i \in \{1, 2, \dots, k\} \\ \text{además, tenemos que } -r_0 &= d_0 &\& r_{k+1}^T d_0 = 0 \\ &\therefore r_{k+1}^T r_i &= 0 \forall i \in \{0, 1, \dots, k\} \quad \square_2\end{aligned}$$

(Con esto se acaba de dar una prueba de que  $r_{k+1} \perp r_i \forall i \in \{0, 1, \dots, k\}$ )

Para terminar la prueba será necesario usar una fórmula recursiva para el residuo. Esto es fácil de lograr si se recuerda el paso tomado anteriormente:

$$\begin{aligned}x_{i+1} &= x_i + \alpha_i d_i \\ \implies Ax_{i+1} &= Ax_i + \alpha_i Ad_i \\ \implies (Ax_{i+1} - b) &= (Ax_i - b) + \alpha_i Ad_i \\ \implies r_{i+1} &= r_i + \alpha_i Ad_i\end{aligned}$$

Con esa fórmula la prueba casi está completa.

Premultiplicamos a  $r_{i+1}$  por  $r_{k+1}$  con  $i \in \{0, 1, \dots, k-1\}$  y, por los resultados anteriores, termina la prueba.

$$\begin{aligned}r_{k+1}^T r_i &= r_{k+1}^T r_i + \alpha_i r_{k+1}^T Ad_i \quad \forall i \in \{0, 1, \dots, k-1\} \\ \implies r_{k+1}^T Ad_i &= 0 \quad \forall i \in \{0, 1, \dots, k-1\}\end{aligned}$$

$\square_3$

Con la prueba anterior finalizada se pasará a ver cualidades particulares del método del Gradiente Conjugado, así como su forma :

## 1. Método de GC

Sea  $x_0$  una aproximación inicial,  $d_0 \leftarrow r_0$ ,  $k \leftarrow 0$

Repetir mientras  $r_k \neq 0$

$$\begin{aligned}\alpha_k &\leftarrow -\frac{r_k^T d_k}{d_k^T A d_k} \\ x_{k+1} &\leftarrow x_k + \alpha_k d_k \\ r_{k+1} &\leftarrow r_k - \alpha_k A d_k \\ \beta_{k+1} &\leftarrow \frac{r_{k+1}^T A d_k}{d_k^T A d_k} \\ d_{k+1} &\leftarrow -r_{k+1} + \beta_{k+1} d_k \\ k &\leftarrow k + 1\end{aligned}$$

Sin embargo, existe otra forma en la que se puede mejorar al método en cuanto a su cálculo computacional simplemente notando ciertas relaciones algebraicas.

$$\blacksquare -r_k^T d_k = r_k^T r_k$$

$$\begin{aligned}d_k &= -r_k + \beta_k d_{k-1} \\ \implies r_k^T d_k &= -r_k^T r_k + \beta_k r_k^T d_{k-1} \\ \text{pero como } r_k &\perp d_{k-1} \\ \implies r_k^T d_k &= -r_k^T r_k; \\ \therefore \alpha_k &\longleftarrow \frac{r_k^T r_k}{d_k^T A d_k}\end{aligned}$$

$$\blacksquare \frac{r_{k+1}^T A p_k}{p_k^T A p_k} = \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k}$$

$$\begin{aligned}\beta_{k+1} &= \frac{r_{k+1}^T A d_k}{p_k^T A d_k} \\ (\text{Sustituyendo } A d_k \text{ de la fórmula recursiva de } r_{k+1}) & \\ &= \frac{r_{k+1}^T \left[ \frac{r_{k+1} - r_k}{\alpha_k} \right]}{p_k^T \left[ \frac{r_{k+1} - r_k}{\alpha_k} \right]} \\ &= \frac{r_{k+1}^T r_{k+1} - r_{k+1}^T r_k}{d_k^T r_{k+1} - d_k^T r_k} \\ &= \frac{r_{k+1}^T r_{k+1} - r_{k+1}^T r_k}{-d_k^T r_k} \\ &= \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k} + \frac{r_{k+1}^T r_k}{d_k^T r_k} \\ \text{pero acabamos de probar que } r_{k+1} &\perp r_k \quad \forall i \in \{0, 1, \dots, k\} \\ \therefore \beta_{k+1} &\longleftarrow \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k}\end{aligned}$$

Por lo tanto, el método de GC queda de la siguiente forma:

## 2. Algoritmo computacional de GC

Sea  $x_0$  una aproximación inicial,  $d_0 \leftarrow -r_0$ ,  $k \leftarrow 0$

Repetir mientras  $r_k \neq 0$

$$\alpha_k \leftarrow \frac{r_k^T r_k}{d_k^T A d_k}$$

$$x_{k+1} \leftarrow x_k + \alpha_k d_k$$

$$r_{k+1} \leftarrow r_k - \alpha_k A d_k$$

$$\beta_{k+1} \leftarrow \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k}$$

$$d_{k+1} \leftarrow -r_{k+1} + \beta_{k+1} d_k$$

$$k \leftarrow k + 1$$

## 3. A continuación estudiaremos algunas propiedades del método de GC que aparecen en el libro de Nocedal & Wright ( *Teoremas 5.1-5.5* )

**Teorema 5.1 :** Para cualquier aproximación inicial  $x_0 \in \mathbb{R}^n$  , la secuencia generada por el algoritmo del GC ( $\{x_k\}_{k=1}^{n-1}$ ) converge a la solución del sistema lineal  $Ax = b$  en a lo más n iteraciones ( $x_n = x^*$ ) .

Para poder probar el teorema 5.1 será útil probar antes el siguiente lema.

**Lema 5.1 :** Si  $A \in \mathbb{R}^{n \times n}$  es una matriz positiva definida y el conjunto de vectores  $\{v_0, v_1, \dots, v_{k-1}\}$  es  $A$ -ortogonal ( $\{v_i\} \neq 0$ ), entonces esos vectores son  $L.I.$ .

$$\begin{aligned} \text{Supongamos que } 0 &= \sum_{i=0}^{n-1} \alpha_i v_i \\ \text{Como } v_i &\neq 0 \quad \forall i \in \{0, 1, \dots, k-1\} \\ \text{y } A \text{ es s.p.d,} &\quad \text{premultiplicamos por } v_l A \\ \implies 0 &= v_l^T A \sum_{i=0}^{n-1} \alpha_i v_i \\ &= \sum_{i=0}^{n-1} \alpha_i v_l^T A v_i \\ &= \alpha_l v_l^T A v_l \end{aligned}$$

Pero como  $A$  es s.p.d. , la única manera en que la ecuación anterior esté balanceada es si  $\alpha_l = 0$ . Nótese que el procedimiento anterior es válido para cualquier  $l \in \{0, 1, \dots, k-1\}$ ,  $\therefore \{v_0, v_1, \dots, v_{k-1}\}$  son  $L.I.$ .  $\square_{\text{lema}_{5,1}}$  .

Volviendo al **Teorema 5.1**, tenemos que  $d_i \perp_A d_j \forall i \neq j$  con  $A$  una matriz s.p.d. Por el lema anterior, el conjunto de direcciones  $\{d_i\}_{i=0}^{n-1}$  es  $L.I.$  y por lo tanto, forma una base en  $\mathbb{R}^n$ . Esto

nos deja con que para todo vector  $z \in \mathbb{R}^n \exists \{\tau_j\}_{j=0}^{n-1}$  tal que  $z$  es combinación lineal del conjunto de direcciones con los coeficientes ( las proyecciones en cada dirección) definidos por  $\tau_i$ .

En particular, esto pasa para el vector  $(x^* - x_0)$ :

$$(x^* - x_0) = \tau_0 d_0 + \tau_1 d_1 + \dots + \tau_{n-1} d_{n-1}$$

Para lograr ver que  $x_n$  y  $x^*$  son los mismos, primero tendr mos que ver a  $x_n$  de la siguiente forma:

$$\begin{aligned} x_n &= x_{n-1} + \alpha_{n-1} d_{n-1} \\ &= (x_{n-2} + \alpha_{n-2} d_{n-2}) + \alpha_{n-1} d_{n-1} \\ &= \dots \\ &= x_0 + \alpha_0 d_0 + \alpha_1 d_1 + \dots + \alpha_{n-1} d_{n-1} \\ \therefore (x_n - x_0) &= \alpha_0 d_0 + \alpha_1 d_1 + \dots + \alpha_{n-1} d_{n-1} \end{aligned}$$

Esto quiere decir que el vector  $(x_n - x_0)$  tiene su propia caracterizaci n en el espacio de direcciones con los coeficientes  $\{\alpha_j\}_{j=0}^{n-1}$ . Para nuestro paso siguiente (que ser  comparar  $\alpha$ 's con  $\tau$ 's) tendremos que calcular la forma expl cita de las  $\tau$ 's. Premultiplicando  $(x^* - x_0)$  por  $d_j^T A$  con  $j \in \{0, 1, \dots, k-1\}$ :

$$\begin{aligned} d_j^T A(x^* - x_0) &= \tau_0 d_j^T A d_0 + \tau_1 d_j^T A d_1 + \dots + \tau_{n-1} d_j^T A d_{n-1} \\ &= \tau_j d_j^T A d_j \\ \implies \tau_j^* &= \frac{d_j^T A(x^* - x_0)}{d_j^T A d_j}; \quad *: d_j^T A d_j \geq 0 \end{aligned}$$

Sea  $k \in \{0, 1, \dots, n-1\}$  arbitraria, entonces, por la f rmula encontrada anteriormente,

$$(x_k - x_0) = \alpha_0 d_0 + \alpha_1 d_1 + \dots + \alpha_{k-1} d_{k-1}$$



Premultiplicando por la izquierda  $d_k^T A$ :

$$\begin{aligned}
d_k^T A(x_k - x_0) &= \alpha_0 d_k^T A d_0 + \alpha_1 d_k^T A d_1 + \dots + \alpha_{k-1} d_k^T A d_{k-1} = 0 \\
&\implies d_k^T A(x^* - x_0) - d_k^T A(x^* - x_0) + d_k^T A(x_k - x_0) = 0 \\
\implies d_k^T A(x_k - x_0) &= d_k^T A(x^* - x_0) - d_k^T A(x_k - x_0) \\
&= d_k^T A(x^* - x_0 - x_k + x_0) \\
&= d_k^T A(x^* - x_k) \\
\text{sustituyendo} & \quad Ax^* = b \\
&= d_k^T (b - Ax_k) \\
&= d_k^T (-r_k)
\end{aligned}$$

Juntando la forma de  $\tau_j$  con el resultado anterior, podemos ver que :

$$\tau_k = \frac{d_k^T A(x^* - x_0)}{d_k^T A d_k} = -\frac{r_k^T d_k}{d_k^T A d_k} \quad \square_{Teo5,1}$$

Con la prueba anterior se encontró que los vectores  $(x_n - x_0)$  y  $(x^* - x_0)$  tienen las mismas proyecciones sobre una base L.I. de  $\mathbb{R}^n$ , ésto, por el teorema de representación, nos permite decir que los dos vectores son el mismo en ese espacio, es decir,  $x_n = x^*$ .

**Teorema 5.2 :** Sea  $x_0$  cualquier aproximación inicial y  $\{x_1, x_2, \dots, x_n\}$  la secuencia generada por el algoritmo del GC, entonces se cumple lo siguiente:

- a)  $r_k^T d_i = 0 \quad \forall i \in \{0, 1, \dots, k-1\}$
- b)  $x_k$  es el minimizador de  $f(x) = \frac{1}{2}x^T A x - b^T x$  sobre el espacio

$$\{x | x = x_0 + \text{gen}\{d_0, d_1, \dots, d_{k-1}\}\}$$

La primera parte de la prueba ya ha sido elaborada anteriormente (en el paso inductivo), por lo que pasaremos a la segunda parte de la prueba, es decir, que  $x_k$  es el minimizador de  $f(x)$ . Para lograr esta prueba calcularemos la forma explícita que tiene el minimizador para después compararlo con  $x_k$ . Como A es s.p.d, sabemos que existe un mínimo en el subespacio generado por las direcciones del método con el origen en  $x_0$ . Cambiaremos de variable para facilitar la manipulación algebraica. Sea  $y = x - x_0$ , el problema que tenemos entonces es :

$$\begin{aligned}
&\underset{y}{\text{minimizar}} && f(y + x_0) \\
&\text{s.a. } y &= & \sum_{i=0}^{k-1} \gamma_i d_i; \\
&&& \text{con } \gamma_i \text{ siendo la proyección de } y \text{ sobre } d_i
\end{aligned}$$

Sea, también,  $g(\gamma) = f(y + x_0)$ , como  $g$  es una función continua y derivable ( en los coeficientes

$\gamma'_i s$ ), intentaremos derivar e igualar a cero para encontrar el mínimo de la función que sólo depende de  $\gamma$ .

$$\begin{aligned}\nabla g(\gamma) &= \sum d_i^T \nabla f(\sum \gamma_j d_j + x_0) \\ &= \sum d_i^T [A(\sum \gamma_j d_j + x_0) - b] \\ &= \sum d_i^T A \sum \gamma_j d_j + \sum d_i^T (Ax_0 - b)\end{aligned}$$

$$\begin{array}{ll}\text{sustituyendo} & r_0 = Ax_0 - b, \\ \text{y porque} & d_i \perp_A d_j; \quad \forall i \neq j, \\ \text{si} & \nabla g(\gamma) = 0 \\ \implies \gamma_j = & -\frac{d_j^T r_0}{d_j^T A d_j}\end{array}$$

Además, como las direcciones son  $A$ -ortogonales, podemos cambiar la relación pasada por una más conveniente :

$$\begin{aligned}d_i^T A d_j &= 0 \quad \forall i \neq j \\ \implies \frac{d_j^T r_0}{d_j^T A d_j} &= \frac{d_j^T (Ax_0 - b)}{d_j^T A d_j} \\ \text{pero } d_j^T (Ax_0 - b) &= d_j^T [A(x_0 + \alpha_0 d_0 + \dots + \alpha_{j-1} d_{j-1}) - b] \\ &= d_j^T (Ax_j - b) \\ \implies \frac{d_j^T r_0}{d_j^T A d_j} &= \frac{d_j^T r_j}{d_j^T A d_j} \\ \therefore y^* = (x_k^* - x_0) &= \sum - \left[ \frac{d_j^T r_j}{d_j^T A d_j} \right] d_j \\ &= \sum_{j=0}^{k-1} \alpha_j d_j \\ &\quad \square_{Teo5,2}.\end{aligned}$$

**Teorema 5.3:** Suponiendo que el  $k$ -ésimo iterando generado por el método de GC no es la solución  $x^*$ , se cumplen las siguientes propiedades:

- a)  $r_k^T r_i$  para  $i = 0, 1, \dots, k-1$ .
- b)  $\text{gen}\{r_0, r_1, \dots, r_k\} = \text{gen}\{r_0, Ar_0, \dots, A^k r_0\}$ .
- c)  $\text{gen}\{d_0, d_1, \dots, d_k\} = \text{gen}\{r_0, Ar_0, \dots, A^k r_0\}$ .
- d)  $d_k^T A d_i = 0$  para  $i \in \{0, 1, \dots, k-1\}$ .

y por lo tanto la secuencia  $\{x_k\}$  converge a  $x^*$  en a lo más  $n$  pasos.

La primera y la última parte de la prueba ya se han demostrado anteriormente (en el paso inductivo). Para la segunda y la tercera parte de la prueba se usará un argumento inductivo:

■ **Caso base ( $k = 0$ )**

Es trivial ver que tanto (b) como (c) se cumplen, es decir, el espacio generado por  $r_0$  pertenece a sí mismo y también al generado por su negativo ( $d_0 = -r_0$ )

■ **Hipótesis de Inducción :** Supongamos que tanto (b) como (c) se cumplen para  $k$ . P.D. se cumplen para  $k + 1$ . es decir:

- $\text{gen}\{r_0, r_1, \dots, r_k\} = \text{gen}\{r_0, Ar_0, \dots, A^k r_0\}$ .
- $\text{gen}\{d_0, d_1, \dots, d_k\} = \text{gen}\{r_0, Ar_0, \dots, A^k r_0\}$ .

Por nuestras hipótesis

$$r_k \in \text{gen}\{r_0, \dots, A^k r_0\} \text{ y } d_k \in \text{gen}\{r_0, \dots, A^k r_0\}$$

$$\begin{aligned} r_k &\in \text{gen}\{r_0, Ar_0, \dots, A^k r_0\}, \quad Ad_k \in \text{gen}\{Ar_0, A^2 r_0, \dots, A^{k+1} r_0\} \\ r_{k+1} = r_k + \alpha_k Ad_k &\implies r_{k+1} \in \text{gen}\{r_0, Ar_0, \dots, A^{k+1} r_0\} \\ \therefore \text{gen}\{r_0, r_1, \dots, r_{k+1}\} &\subseteq \text{gen}\{r_0, Ar_0, \dots, A^{k+1} r_0\} \end{aligned}$$

Por nuestras hipótesis también es cierto que:  $A^k r_0 \in \text{gen}\{d_0, d_1, \dots, d_k\}$  ,  
premultiplicando por A, tenemos que :  $A^{k+1} r_0 \in \text{gen}\{Ad_0, Ad_1, \dots, Ad_k\}$   
pero, por nuestra fórmula del residuo, tenemos que  $\alpha_j Ad_j = r_{j+1} - r_j$

$$\therefore \text{gen}\{r_0, r_1, \dots, r_k\} = \text{gen}\{r_0, Ar_0, \dots, A^k r_0\}.$$

Recordar que  $p_{k+1}$  es generado por  $r_{k+1}$  y  $p_k$ ,  
Además, utilizando las hipótesis de inducción, veamos que

$$\begin{aligned} \text{gen}\{d_0, d_1, \dots, d_k, d_{k+1}\} &= \text{gen}\{d_0, d_1, \dots, d_k, r_{k+1}\} \\ &= \text{gen}\{r_0, Ar_0, \dots, A^k r_0, r_{k+1}\} \\ (\text{por el apartado anterior}) &= \text{gen}\{r_0, Ar_0, \dots, A^{k+1} r_0\} \end{aligned}$$

□<sub>Teo5,3</sub>

**Lema de tasa de convergencia :** Para hacer la prueba del teorema 5.4 y el teorema 5.5 se necesitará de un lema acerca de la tasa de convergencia del método del GC. El resultado es un poco largo para probar pero, en sí, lo que hace es ver a la  $k$ -ésima  $x$  como el mejor polinomio de A que aproxima a  $x^*$  de grado  $k$ . Después se demuestra que con una representación adecuada por una base adecuada, ese polinomio que minimiza la A distancia entre  $x$  y  $x^*$  en realidad se puede ver como un polinomio sobre los valores propios de A. Eso es

importante porque quiere decir que cualquier otro polinomio de grado  $k$  (en los reales) es peor aproximación que el definido por  $x_k$ . Expresando, después, a  $x_k$  en término de  $x_0$  y usando propiedades de ortogonalidad de la base seleccionada se llega al siguiente resultado:

$$\|x_{k+1} - x^*\|_A^2 = \min_{P_k} \max_{1 \leq i \leq n} (1 + \lambda_i P_k(\lambda_i))^2 \|x_0 - x^*\|_A^2. \quad \text{VI} \quad \text{VII}$$

**Teorema 5.4:** Si  $A$  sólo tiene  $r$  valores propios distintos, entonces el algoritmo del GC terminará en a lo más  $r$  iteraciones.

Sea  $Q_r(\lambda)$ :

$$Q_r(\lambda) = \left[ \frac{(-1)^r}{\tau_1 \tau_2 \dots \tau_r} \right] (\lambda - \tau_1)(\lambda - \tau_2) \dots (\lambda - \tau_r)$$

Notar que  $Q_r(\lambda_i) = 0$  para  $i = 1, \dots, r$  y  $Q_r(0) = 1$

Por lo tanto,  $Q_r(\lambda) - 1$  es un polinomio de grado  $r$  con raíz en  $\lambda = 0$

Sea  $P_{r-1}$  de grado  $r-1$ :

$$P_{r-1}(\lambda) = \frac{Q_r(\lambda) - 1}{\lambda}.$$

$$(\iff Q_r(\lambda) = 1 + \lambda P_{r-1}(\lambda))$$

Por el lema sobre la tasa de convergencia anterior:

$$0 \leq \min_{P_k} \max_{1 \leq i \leq n} [1 + \lambda_i P_{r-1}(\lambda_i)]^2 \leq \max_{1 \leq i \leq n} [1 + \lambda_i P_{r-1}(\lambda_i)]^2 = \max_{1 \leq i \leq n} Q_r(\lambda_i) = 0$$

$$\text{y } \therefore \|x_r - x^*\|_A^2 = 0 \text{ y } x_r = x^*$$

y el algoritmo de GC termina en a lo más  $r$  iteraciones.  $\square_{Teo5,4}$

**Teorema 5.5:** Sean  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  los valores propios de  $A$ , entonces

$$\|x_{k+1} - x^*\|_A^2 \leq \left[ \frac{\lambda_{n-k} - \lambda_1}{\lambda_{n-k} + \lambda_1} \right]^2 \|x_0 - x^*\|_A^2$$

Algo importante del lema demostrado anteriormente es que nos deja jugar con la noción de polinomio aproximado en lugar de con  $x_k$ . Como la cota marca que la norma- $A$  al cuadrado del vector  $x_k - x^*$  siempre es menor o igual a la **mejor** aproximación en sólo **algunos** puntos delimitados (los valores propios de  $A$ ) al cuadrado por la misma norma al cuadrado

pero del vector  $x_0 - x^*$ .

Lo que esto nos da de libertad es que sólo necesitamos elegir un polinomio adecuado (ya que la cota es por el mejor de ellos, cualquier otro polinomio de grado  $k$  cumple la cota), que en los valores propios de  $A$  esté acotado adecuadamente y podemos llegar al resultado deseado.

Sea  $q(\lambda) = 1 + \lambda p(\lambda)$  un polinomio de grado  $m + 1$  que tiene una raíz en  $(\lambda_1 + \lambda_{n-k})/2$  y sus otras  $m$  raíces en los  $m$  valores propios más grandes de  $A$ . Notemos que, como  $q(\lambda)$  es un polinomio, entonces es diferenciable. Además,  $dq(\lambda)$  tiene  $m$  raíces (todas entre las raíces de  $q(\lambda)$ ) y  $d^2q(\lambda)$  tiene  $m - 1$  raíces entre las  $m$  raíces de  $dq(\lambda)$ . Una propiedad de los polinomios que nos van a servir es:

- Fuera de sus raíces, los polinomios crecen (o decrecen) a velocidad y forma  $x^k$  con  $k$  siendo el grado del polinomio.  
Para "acotar" nuestro polinomio, sólo es necesario analizarlo en una región (ventaja de la cota obtenida en el lema). Como, para toda  $\lambda \in \{\lambda_n, \lambda_{n-1}, \dots, \lambda_{n-k+1}\}$  el valor de  $q(\lambda)$  es cero, no será necesario analizar esos intervalos. Todos los valores que nos interesan (tanto para  $\lambda$  como para  $q(\lambda)$ ) se encuentran entre el cero (porque  $A$  es s.p.d) y  $\lambda_{n-k}$ . Existen dos opciones para nuestro polinomio, que  $k$  sea par (y por lo tanto nuestro polinomio sea de grado impar) o que  $k$  sea impar (y nuestro ...).
- Sea  $k$  par (impar), entonces se cumple que  $q(\lambda)$  es un polinomio de grado impar (par), por lo tanto,  $q(x) \leq 0$  ( $x \in [0, \frac{\lambda_1 + \lambda_{n-k}}{2}]$ ). Sin embargo, como  $d^2q(x) < 0$  ( $> 0$ ) en ese intervalo, la función es cóncava (convexa) en ese intervalo y se encuentra arriba (abajo) de la línea

$$-1 + \frac{2\lambda}{\lambda_1 + \lambda_{n-k}}$$

$$(1 - \frac{2\lambda}{\lambda_1 + \lambda_{n-k}})$$

- Para el intervalo que va de  $[\frac{\lambda_1 + \lambda_{n-k}}{2}, \lambda_{n-k}]$  pasa justo lo contrario, como las dos curvas se cruzan (el polinomio y la línea) en el mismo punto (la raíz), se necesitaría que hubieran dos cambios de signo en  $dq(\lambda)$  para que la línea recta terminara abajo del polinomio, pero como sólo hay una raíz de  $dq(\lambda)$  en ese intervalo, eso no sucede. Por lo tanto, la línea  $-1 + \frac{2\lambda}{\lambda_1 + \lambda_{n-k}}$  ( $1 - \frac{2\lambda}{\lambda_1 + \lambda_{n-k}}$ ) va por arriba (abajo) de  $q(\lambda)$  en ese intervalo.
- Juntando los dos intervalos, se llega a la conclusión de que:

$$|q(\lambda)| \leq \left| 1 - \frac{2\lambda}{\lambda_1 + \lambda_{n-m}} \right| \quad \text{en el intervalo } [\lambda_1, \lambda_{n-k}]$$

4. El método de GC tiene una propiedad muy útil en optimización. El problema (1) adopta la forma

$$\underset{p}{\text{minimizar}} \quad \frac{1}{2} p^T \nabla^2 f(x) p + \nabla f(x)^T p, \quad (6)$$

en donde se han omitido los índices por claridad. Cada paso del método de GC requiere de un producto matriz-vector  $\nabla^2 f(x)d_k$ . Investiga cómo aproximar el producto anterior mediante diferencias del gradiente. Al método resultante se le conoce como método de Newton libre de Hessiana.

##### 5. NOTA: GC para el método de Newton.

El sistema por resolver en cada iteración del método de Newton es

$$\nabla^2 f(x)p^N = -\nabla f(x),$$

el método de GC aproxima a  $p^N$  por medio de la sucesión

$$\{p_0, p_1, \dots, p_{n-1}\}.$$

Observar que los vectores  $p_k$  tienen el papel de los vectores  $x_k$  en los algoritmos anteriores, es decir las aproximaciones  $p_k$  se actualizan como sigue

$$p_{k+1} = p_k + \alpha_k d_k.$$

El método inicia con  $p_0 = 0$ , de donde es inmediato concluir que  $r_0 = -\nabla f(x)$ ; por lo tanto, la primera dirección para GC es la dirección opuesta al gradiente

$$d_0 = -\nabla f(x).$$

Si el producto  $d_0^T \nabla^2 f(x)d_0$  resultara negativo o cero, el método terminaría con

$$p = -\nabla f(x).$$

En iteraciones subsecuentes GC terminaría con la aproximación  $p_k$  si  $\|r_k\| < TOL$ , o bien

$$d_{k+1}^T \nabla^2 f(x)d_{k+1} \leq 0.$$

El método cumple con el teorema 5.5 para cálculos de  $p_k$  por lo que la mejoría en la dirección siempre tendrá una cota. A este método se le conoce como método del Gradiente Conjugado Parcial.

6. Estudiar la prueba de convergencia cuadrática del método de Newton (Nocedal & Wright, página 44). Justificar los pasos omitidos.
7. Versión computacional de GC que se detiene cuando ocurre alguna de las siguientes situaciones se encuentra anexado junto con este documento : a)  $\|r_k\| \leq TOL_1$ , b) el número de iteraciones excede un límite preestablecido; c)  $d_k^T \nabla^2 f(x)d_k \leq TOL_2$ .

Ahora analizaremos la construcción de **Powell-Symmetric-Broyden (PSB)**.

En este caso vamos a resolver el problema de optimización convexa:

$$\min \phi(m) = \frac{1}{2} \sum_{i=1}^n (m_{ii} - b_{ii})^2 + \sum_{i < j} (m_{ij} - b_{ij})^2$$

s. a.  $Am - y = 0$

Donde definimos la forma extendida de una matriz simétrica  $B$  pero considerando sólo la parte triangular superior de  $B$ , así pues obtenemos:

$$b = [b_{11} \ b_{12} \ \dots \ b_{1n} \ | \ b_{22} \ \dots \ b_{2n} \ | \dots | \ b_{nn}]$$

Y ahora tenemos una nueva versión de la matriz  $A$ :

$$A = \left[ \begin{array}{cccc|cccc} s_1 & s_2 & \dots & s_n & 0 & 0 & \dots & 0 \\ 0 & s_1 & \dots & 0 & s_2 & s_3 & \dots & s_n \\ \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & s_1 & 0 & 0 & \dots & s_2 \end{array} \right] \dots s_n$$

Las matrices de *Karush Kuhn – Tucker* son muy parecidas al caso anterior:

$$\begin{bmatrix} H & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} b_+ \\ -\lambda \end{bmatrix} = \begin{bmatrix} Hb \\ y \end{bmatrix}$$

Con la diferencia de que ahora está presente la matriz  $H$  que es una matriz diagonal y las entradas de la misma son 1 o 2. Las entradas con valor 1 corresponden a los elementos diagonales  $m_{ii}$  y las entradas con valor 2 corresponden a los elementos fuera de la diagonal  $m_{ij}$ .

Ahora resolvemos el primer bloque de las ecuaciones de *Karush Kuhn – Tucker*, es decir:  $b_+ = b + H^{-1}A^T\lambda$

Si llevamos a cabo la multiplicación de matrices  $H^{-1}A^T\lambda$  obtendremos que:

$$H^{-1}A^T\lambda = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & \frac{1}{2} & 0 & \dots & 0 \\ 0 & 0 & \frac{1}{2} & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \begin{bmatrix} s_1 & 0 & 0 & \dots & 0 \\ s_2 & s_1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & 0 \\ 0 & s_2 & 0 & \dots & 0 \\ 0 & s_3 & s_2 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 0 & 0 & 0 & \dots & s_n \end{bmatrix} \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_n \end{bmatrix} = \begin{bmatrix} \lambda_1 s_1 \\ \frac{\lambda_1 s_2 + \lambda_2 s_1}{2} \\ \vdots \\ \frac{\lambda_1 s_n + \lambda_n s_1}{2} \\ \vdots \\ \lambda_n s_n \end{bmatrix}$$

Por otra parte, tomemos la siguiente multiplicación:  $\frac{1}{2}(s\lambda^T + \lambda s^T)$  donde:

$$s\lambda^T = \begin{bmatrix} s_1\lambda_1 & s_1\lambda_2 & \dots & s_1\lambda_n \\ s_2\lambda_1 & s_2\lambda_2 & \dots & s_2\lambda_n \\ \vdots & \vdots & \dots & \vdots \\ s_n\lambda_1 & s_n\lambda_2 & \dots & s_n\lambda_n \end{bmatrix} \text{ y por otro lado } \lambda s^T = \begin{bmatrix} \lambda_1 s_1 & \lambda_1 s_2 & \dots & \lambda_1 s_n \\ \lambda_2 s_1 & \lambda_2 s_2 & \dots & \lambda_2 s_n \\ \vdots & \vdots & \dots & \vdots \\ \lambda_n s_1 & \lambda_n s_2 & \dots & \lambda_n s_n \end{bmatrix}$$

$$\text{Por lo que al tomar } \frac{1}{2}(s\lambda^T + \lambda s^T) \text{ tenemos: } \frac{1}{2}(s\lambda^T + \lambda s^T) = \begin{bmatrix} \lambda_1 s_1 & \frac{\lambda_1 s_2 + s_1 \lambda_2}{2} & \dots & \frac{\lambda_1 s_n + s_1 \lambda_n}{2} \\ \frac{\lambda_2 s_1 + s_2 \lambda_1}{2} & \lambda_2 s_2 & \dots & \frac{\lambda_2 s_n + s_2 \lambda_n}{2} \\ \vdots & \vdots & \dots & \vdots \\ \frac{\lambda_n s_1 + s_n \lambda_1}{2} & \frac{\lambda_n s_2 + s_n \lambda_2}{2} & \dots & \lambda_n s_n \end{bmatrix}$$

Es fácil observar que  $H^{-1}A^T\lambda$  es la expansión del triángulo superior de  $\frac{1}{2}(s\lambda^T + \lambda s^T)$

Organizando de forma matricial tenemos:

$$B_+ = B + \frac{1}{2}(s\lambda^T + \lambda s^T)$$

Ahora tomamos el producto  $H^{-1}A^T$  y lo premultiplicamos por  $A$  obteniendo así la matriz:

$$AH^{-1}A^T = \begin{bmatrix} s_1^2 + \sum_{i=1}^n s_i^2 & s_1 s_2 & s_1 s_3 & \dots & s_1 s_n \\ s_1 s_2 & s_2^2 + \sum_{i=1}^n s_i^2 & s_2 s_3 & \dots & s_2 s_n \\ \vdots & \vdots & \vdots & \dots & \vdots \\ s_1 s_n & s_2 s_n & s_3 s_n & \dots & s_n^2 + \sum_{i=1}^n s_i^2 \end{bmatrix}$$

Ahora usamos la fórmula de *Shermann - Morrison*:  $(A + uv^T)^{-1} = A^{-1} - \frac{A^{-1}uv^T A^{-1}}{1+v^T A^{-1}u}$  donde sustituimos de la fórmula:  $A = \frac{1}{2}s^T s \mathbb{I}$ ,  $u = s$  y  $v^T = s^T$

Obtenemos el siguiente resultado:

$$(AH^{-1}A^T)^{-1} = \frac{2}{s^T s} \left( \mathbb{I} - \frac{ss^T}{2s^T s} \right) \text{ de donde podemos despejar}$$

$$\lambda = \frac{2}{s^T s} \left( \mathbb{I} - \frac{ss^T}{2s^T s} \right) (y - Bs)$$

Finalmente, al sustituir  $\lambda$  obtenemos:

$$B_+ = B + \frac{(y-Bs)s^T + s(y-Bs)^T}{s^T s} - \frac{(y-Bs)^T s}{(s^T s)^2} ss^T \quad \square_{P.S.B.}$$