

MATEMÁTICA SUPERIOR APLICADA

SOLUCIÓN DE SISTEMAS DE ECUACIONES ALGEBRAICAS LINEALES

Planteo del Problema. Teoremas Básicos

Se trata de resolver n ecuaciones lineales simultáneas con n incógnitas:

$$\sum_{j=1}^n a_{ij} x_j = b_i \text{ con } i = 1, 2, \dots, n \quad (1)$$

o, en su forma matricial compacta:

$$\underline{\underline{A}} \cdot \underline{x} = \underline{b} \quad (2)$$

donde $\underline{\underline{A}} = [a_{ij}]$ es la matriz de coeficientes (cuadrada, $n \times n$) y $\underline{x}^t = (x_1, \dots, x_n)$ es el vector de incógnitas, siendo $\underline{b}^t = (b_1, \dots, b_n)$ el vector de términos independientes. Para el caso que nos interesa, tanto la matriz $\underline{\underline{A}}$ como el vector \underline{b} tienen componentes reales. En adelante se obviará la notación transpuesta de la matriz, suponiendo que en las operaciones entre matrices se disponen éstas de tal forma que sean compatibles para las mismas. Si definimos la denominada matriz aumentada como:

$$\underline{\underline{A}}_b = \left[\underline{\underline{A}} \quad \underline{b} \right] \text{ Matriz ampliada de } \underline{\underline{A}} \left[n \times (n+1) \right] \quad (3)$$

dado que \underline{b} es un vector columna; y recordando la definición de rango de una matriz: $r(\underline{\underline{A}})$, el teorema básico de existencia de solución establece:

- 1) El sistema de ecuaciones (Ecuación (2)) tiene solución sí y sólo sí: $r(\underline{\underline{A}}) = r(\underline{\underline{A}}_b)$.
- 2) Si $r(\underline{\underline{A}}) = r(\underline{\underline{A}}_b) = k < n$, luego las x_1, x_2, \dots, x_k son variables cuyas columnas son linealmente independientes en $\underline{\underline{A}}$, de modo que las restantes $(n-k)$ variables pueden asignarse arbitrariamente. O dicho de otra forma, hay una familia paramétrica de $(n-k)$ soluciones.
- 3) Si $r(\underline{\underline{A}}) = r(\underline{\underline{A}}_b) = n$, hay una única solución.

Corolario: Para el caso homogéneo ($\underline{b} = \underline{0}$), o sea $\underline{\underline{A}} \cdot \underline{x} = \underline{0}$ habrá solución no trivial, si y sólo si $r(\underline{\underline{A}}) < n$.

Para estos problemas, cosa que no ocurre en el caso no lineal, existe solución analítica (recordemos la denominada *Regla de Cramer*), pero la dificultad reside principalmente en computar esa solución. La evaluación de determinantes no hace práctico dicho procedimiento analítico; luego, el problema es desarrollar algoritmos computacionales más eficientes, es decir que sean más rápidos, sobre todo en el número de operaciones necesarias y que además

sean robustos de modo que la solución calculada sea lo más precisa posible.

Un punto vital es discutir cómo se espera que sea la matriz de coeficientes. En general, puede encontrarse entre alguna de estas dos categorías:

- i. Llena pero no muy grande. Es decir con muy pocos ceros, y en donde n no sea mayor que 100, por ejemplo.
- ii. Dispersa y relativamente muy grande, denominadas también ralas. En estos casos, son muy pocos (en relación al orden) los elementos distintos de cero y n puede ser mayor a 1000.

Naturalmente, los métodos desarrollados deben estar dirigidos a resolver alguna de estas dos categorías, y si es posible haciendo uso de sus características para incrementar su eficiencia.

Un problema, cual es la condición del sistema, se puede analizar considerando un vector residual \underline{r} , cuando se tiene una solución calculada $\underline{x}^{(c)}$. Esto es:

$$\underline{r} = \underline{b} - \underline{A} \cdot \underline{x}^{(c)} \quad (4)$$

se sabe que si \underline{x}^* es la solución exacta del SEAL:

$$\underline{A} \cdot \underline{x}^* = \underline{b} \text{ o bien } \underline{r} = \underline{b} - \underline{A} \cdot \underline{x}^* = \underline{0} \quad (5)$$

entonces será:

$$\underline{r} = \underline{A} \cdot (\underline{x}^* - \underline{x}^{(c)})$$

luego,

$$(\underline{x}^* - \underline{x}^{(c)}) = \underline{A}^{-1} \cdot \underline{r} \quad (6)$$

de donde se ve que aunque \underline{r} tenga elementos muy chicos, si \underline{A}^{-1} (matriz inversa) contiene coeficientes muy grandes, la diferencia entre \underline{x}^* y $\underline{x}^{(c)}$ puede ser aún muy grande. Esto permite anticipar la importancia de un escalado en los coeficientes de la matriz \underline{A} original, ya que aunque el vector residual \underline{r} impuesto sea pequeño, el error encontrado para la solución puede ser muy grande.

Métodos Directos

Un método directo para hallar la solución es uno en el cual, si todos los cálculos (computaciones) fueran llevados a cabo sin error de redondeo conduciría a la solución exacta del sistema dado. Prácticamente todos están basados en la *técnica de eliminación*. El error de truncamiento para estos métodos es intrascendente.

Eliminación Gaussiana

Desarrollando la Ecuación (1), se obtiene el sistema en la siguiente forma:

$$\begin{aligned}
 a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n &= b_1 = a_{1,n+1} \\
 a_{21}x_1 + a_{22}x_2 + a_{23}x_3 + \dots + a_{2n}x_n &= b_2 = a_{2,n+1} \\
 &\dots\dots\dots \\
 &\dots\dots\dots \\
 a_{n1}x_1 + a_{n2}x_2 + a_{n3}x_3 + \dots + a_{nn}x_n &= b_n = a_{n,n+1}
 \end{aligned}
 \tag{7}$$

Se supone que la matriz es no singular, y que $a_{11} \neq 0$ de manera de poder dividir la primer columna por a_{11} y así restar para las ecuaciones, donde $i=2, \dots, n$. Esto da como resultado:

$$\begin{aligned}
 a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n &= b_1 = a_{1,n+1} \\
 a_{22}^{(1)}x_2 + a_{22}^{(1)}x_3 + \dots + a_{2n}^{(1)}x_n &= b_2^{(1)} = a_{2,n+1}^{(1)} \\
 &\dots\dots\dots \\
 &\dots\dots\dots \\
 a_{n2}^{(1)}x_2 + \dots\dots\dots + a_{nn}^{(1)}x_n &= b_n^{(1)} = a_{n,n+1}^{(1)}
 \end{aligned}
 \tag{8}$$

siendo los $a^{(1)}_{ij}$, tal que:

$$a_{ij}^{(1)} = a_{ij} - \frac{a_{i1}}{a_{11}}a_{1j} ; i = 2, \dots, n ; j = 2, \dots, (n+1)
 \tag{9}$$

si fuese $a_{11} = 0$ se intercambian columnas, y se opera en consecuencia. Por comodidad se define: $m_{i1} = a_{i1}/a_{11}$ con $i=2, \dots, n$.

De igual forma se continúa con el procedimiento haciendo ahora (si $a^{(1)}_{22} \neq 0$) $m_{i2} = a^{(1)}_{i2}/a^{(1)}_{22}$ con $i=3, \dots, n$ resultando, al restar, el siguiente sistema:

$$\begin{aligned}
 a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n &= b_1 = a_{1,n+1} \\
 a_{22}^{(1)}x_2 + a_{22}^{(1)}x_3 + \dots + a_{2n}^{(1)}x_n &= b_2^{(1)} = a_{2,n+1}^{(1)} \\
 a_{32}^{(2)}x_3 + \dots + a_{3n}^{(2)}x_n &= b_3^{(2)} = a_{3,n+1}^{(2)} \\
 &\dots\dots\dots \\
 a_{n2}^{(2)}x_3 + \dots + a_{nn}^{(2)}x_n &= b_n^{(2)} = a_{n,n+1}^{(2)}
 \end{aligned}
 \tag{10}$$

donde:

$$a_{ij}^{(2)} = a_{ij}^{(1)} - m_{i2}a_{2j}^{(1)} ; i = 3, \dots, n ; j = 3, \dots, (n+1)$$

Y, continuando con el procedimiento hasta $(n-1)$ pasos llegamos al sistema final:

$$\begin{aligned}
a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n &= b_1 = a_{1,n+1} \\
a_{22}^{(1)}x_2 + a_{22}^{(1)}x_3 + \dots + a_{2n}^{(1)}x_n &= b_2^{(1)} = a_{2,n+1}^{(1)} \\
&+ a_{32}^{(2)}x_3 + \dots + a_{3n}^{(2)}x_n = b_3^{(2)} = a_{3,n+1}^{(2)} \\
&\dots \dots \dots \\
&+ a_{nn}^{(n-1)}x_n = b_n^{(n-1)} = a_{n,n+1}^{(n-1)}
\end{aligned} \tag{11}$$

con los elementos en la diagonal (distintos de cero), tal que:

$$\begin{aligned}
a_{ij}^{(k)} &= a_{ij}^{(k-1)} - m_{ik} a_{kj}^{(k-1)} ; m_{ik} = \frac{a_{ik}^{(k-1)}}{a_{kk}^{(k-1)}} \\
k &= 1, \dots, n-1 \\
j &= k+1, \dots, (n+1) \\
i &= k+1, \dots, n \\
a_{ij}^{(0)} &= a_{ij}
\end{aligned} \tag{12}$$

Luego, la solución es fácilmente calculada por sustitución hacia atrás, al terminar el procedimiento de eliminación. Esto es:

$$x_i = \frac{1}{a_{ii}^{(i-1)}} \left[a_{i,n+1}^{(i-1)} - \sum_{j=i+1}^n a_{ij}^{(i-1)} x_j \right] ; i = n, \dots, 1 \tag{13}$$

Método de Gauss-Jordan

Una variante muy importante es el *procedimiento de reducción o eliminación* de Gauss-Jordan. En este caso se procede a eliminar con los elementos diagonales en toda la columna, dejando sólo ese elemento; y se deriva en un sistema que, comparando con la Ecuación (10), en una primera pasada tiene la forma:

$$\begin{aligned}
a_{11}x_1 + 0 + a_{13}^{(1)}x_3 + \dots + a_{1n}^{(1)}x_n &= b_1 = a_{1,n+1}^{(1)} \\
a_{22}^{(1)}x_2 + a_{23}^{(1)}x_3 + \dots + a_{2n}^{(1)}x_n &= b_2^{(1)} = a_{2,n+1}^{(1)} \\
a_{33}^{(2)}x_3 + \dots + a_{3n}^{(2)}x_n &= b_3^{(2)} = a_{3,n+1}^{(2)} \\
&\dots \dots \dots \\
a_{n3}^{(2)}x_3 + \dots + a_{nn}^{(2)}x_n &= b_n^{(2)} = a_{n,n+1}^{(2)}
\end{aligned} \tag{14}$$

Siguiendo con el procedimiento se hacen cero todos los elementos, excepto los correspondientes a la diagonal, resultando:

$$\begin{aligned}
a_{11}x_1 + 0 + 0 + \dots + 0 &= b_1^{(n-1)} = a_{1,n+1}^{(n-1)} \\
a_{22}^{(1)}x_2 + 0 + \dots + 0 &= b_2^{(n-1)} = a_{2,n+1}^{(n-1)} \\
a_{33}^{(2)}x_3 + \dots + 0 &= b_3^{(n-1)} = a_{3,n+1}^{(n-1)} \\
&\dots \\
a_{nn}^{(n-1)}x_n &= b_n^{(n-1)} = a_{n,n+1}^{(n-1)}
\end{aligned} \tag{15}$$

de modo que la solución es simplemente:

$$x_i = \frac{a_{i,n+1}^{(n-1)}}{a_{ii}^{(i-1)}} ; i = 1, 2, \dots, n \tag{16}$$

A pesar de lo que resulta del procedimiento, la eliminación Gaussiana es la más eficiente de las dos, considerando sólo las multiplicaciones y divisiones (recordar, además, la acumulación del error de redondeo), llegando, para grandes sistemas ($n \gg 1$), el método Gauss-Jordan a requerir cerca de un 50% más de operaciones que el de Gauss.

Se ha trabajado mucho para lograr *formas compactas* del método de Gauss, no sólo para ahorrar espacio de almacenamiento (memoria) sino también para mejorar la precisión en los cálculos que más inciden en el resultado. Para ello se han definido matrices especiales en cuanto a la característica de su formulación, las cuales permiten ahorrar tiempo de cálculo una vez aplicadas. No trataremos este punto aquí, remitiendo al lector a la bibliografía recomendada al final del capítulo.

Análisis de Errores

Debido a que generalmente no es posible obtener la solución exacta, se considerarán los posibles errores y sus cotas. Las fuentes de error son variaciones en los elementos de \underline{A} y de \underline{b} , ya sean originales o debidas al redondeo. Estudiando primero el caso más simple, que considera cambios sólo en \underline{b} , será:

$$\underline{A} \cdot \underline{x} = \underline{b} + \underline{\Delta b} \tag{17}$$

si resulta a su vez:

$$\underline{x} = \underline{x}^* + \underline{\Delta x}$$

se ve que:

$$\underline{A} \cdot (\underline{x}^* + \underline{\Delta x}) = \underline{b} + \underline{\Delta b}$$

con lo que:

$$\underline{A} \cdot \underline{\Delta x} = \underline{\Delta b}$$

Si \underline{A} es no singular, entonces resulta:

$$\underline{\Delta x} = \underline{A}^{-1} \cdot \underline{\Delta b} \quad (18)$$

y/o, tomando normas:

$$\|\underline{\Delta x}\| \leq \|\underline{A}^{-1}\| \|\underline{\Delta b}\| \quad (19)$$

de igual modo:

$$\|\underline{b}\| \leq \|\underline{A}\| \|\underline{x}^*\| \quad (20)$$

Multiplicando miembro a miembro las Desigualdades (19) y (20) se obtiene una cota para el error relativo:

$$\frac{\|\underline{\Delta x}\|}{\|\underline{x}\|} \leq \|\underline{A}\| \|\underline{A}^{-1}\| \frac{\|\underline{\Delta b}\|}{\|\underline{b}\|} \quad (21)$$

La cantidad $\|\underline{A}\| \|\underline{A}^{-1}\|$ se denomina *número de condición de \underline{A}* y se indica $K(\underline{A})$; además se verifica que $K(\underline{A}) \geq 1$, de modo que:

$$\frac{\|\underline{\Delta x}\|}{\|\underline{x}\|} \leq K(\underline{A}) \frac{\|\underline{\Delta b}\|}{\|\underline{b}\|} \quad (22)$$

Así, si $K(\underline{A})$ es cercano a 1, se dice que \underline{A} está *bien condicionada*; y si es muy grande nos encontramos frente a un caso *mal condicionado*.

Si la fuente de error fueran los elementos de la matriz de coeficientes, esto es:

$$(\underline{A} + \underline{\Delta A}) \cdot (\underline{x}^* + \underline{\Delta x}) = \underline{b} \quad (23)$$

Por un procedimiento similar al anterior se llega a que:

$$\frac{\|\underline{\Delta x}\|}{\|\underline{x}^* + \underline{\Delta x}\|} \leq K(\underline{A}) \frac{\|\underline{\Delta A}\|}{\|\underline{A}\|} \quad (24)$$

con igual connotación que la Ecuación (21). Los efectos del valor $\|\underline{\Delta A}\|$ pueden subsanarse, en parte, trabajando con mayor precisión (mayor retención de cifras significativas durante el cálculo). En general se puede afirmar que $\frac{\|\underline{\Delta A}\|}{\|\underline{A}\|} \approx C(n)\varepsilon$ donde $C(n)$ es una constante del orden de la dimensión de la matriz y ε es el epsilon de la computadora.

Refinamiento Iterativo

Habíamos definido al vector residuo como $\underline{r} = \underline{b} - \underline{A} \cdot \underline{x}^{(c)}$ En situaciones en las que nos

interesa una solución muy cercana a la verdadera: $\underline{x}^* = \underline{A}^{-1} \cdot \underline{b}$ con un vector residuo \underline{r} que sea pequeño o, en caso en que $\underline{x}^{(c)}$ sea tal que:

$$\underline{e} = \underline{x}^* - \underline{x}^{(c)} \quad (25)$$

sea muy pequeño en términos relativos a \underline{x}^* ; si $\underline{b} = \underline{A} \cdot \underline{x}^*$, entonces resulta:

$$\underline{r} = \underline{A} \cdot \underline{x}^* - \underline{A} \cdot \underline{x}^{(c)} = \underline{A} \cdot (\underline{x}^* - \underline{x}^{(c)}) = \underline{A} \cdot \underline{e} \quad (26)$$

Tomando normas, y procediendo como ya se ha visto, se puede probar que:

$$\frac{1}{K(\underline{A})} \frac{\|\underline{r}\|}{\|\underline{b}\|} \leq \frac{\|\underline{e}\|}{\|\underline{x}^*\|} \leq K(\underline{A}) \frac{\|\underline{r}\|}{\|\underline{b}\|} \quad (27)$$

Por lo que, si el número de condición $K(\underline{A})$ es cercano a la unidad, pequeños errores relativos de \underline{b} y \underline{x}^* siguen la misma tendencia. Este hecho no ocurre para sistemas mal condicionados, como se había anticipado.

La respuesta al problema dependerá de la condición de \underline{A} y de la precisión de la aritmética (redondeo). Esto lleva a que el cómputo del residuo \underline{r} debe hacerse con doble precisión, por la razón de que \underline{r} suele ser del mismo orden de magnitud que el error de redondeo.

Formalizar un procedimiento que contemple el problema es lo que se denomina *refinamiento iterativo*. En realidad se arma un esquema iterativo con una solución calculada previamente, tal que el residuo para una etapa m de cálculo sea:

$$\underline{r}^{(m)} = \underline{b} - \underline{A} \cdot \underline{x}^{(m)} \quad ; \quad m=1,2,\dots \quad (28)$$

y, de acuerdo con la Ecuación (25),

$$\underline{x}^{(m+1)} = \underline{x}^{(m)} + \underline{e}^{(m)} \quad (29)$$

donde $\underline{e}^{(m)} = \underline{A}^{-1} \cdot \underline{r}^{(m)}$, según se desprende de la Ecuación (26). Si $\underline{x}^{(m+1)}$ no es satisfactoria se procede a la etapa $(m+1)$ de cálculo. En cambio, si $\left\| \frac{\underline{e}^{(m+1)}}{\underline{e}^{(m)}} \right\| < \varepsilon$, esto es, menor que una cierta tolerancia preestablecida aceptamos la solución $\underline{x}^{(m)}$. Así, se conforma un criterio de terminación que funciona muy bien en la práctica. Cada etapa de este proceso, además, es mucho más rápida que la solución del problema original. Se debe recordar que el cómputo de la Ecuación (28) se realiza en doble precisión, lo que implica un aumento en el requerimiento de memoria computacional.

Métodos Iterativos

Su basamento es idéntico al método de aproximaciones sucesivas visto en el caso de

sistemas de ecuaciones no lineales, con lo que, empezando con un vector inicial $\underline{x}^{(0)}$ se genera una sucesión de vectores, $\underline{x}^{(0)}, \underline{x}^{(1)}, \underline{x}^{(2)}, \dots, \underline{x}^{(n)}$, tal que:

$$\underline{x}^{(n+1)} = \underline{F}^{(n)}(\underline{x}^{(n)}, \underline{x}^{(n-1)}, \underline{x}^{(n-2)}, \dots, \underline{x}^{(n-k)}) \quad (30)$$

Si la función iteradora $\underline{F}^{(n)}$ no depende del nivel de iteración, la recurrencia se la llama estacionaria.

Para la mayoría de las matrices estos métodos requieren más cálculo, para un deseado grado de convergencia, que los métodos directos; pero para matrices ralas (de gran interés en aplicaciones) el esfuerzo computacional es comparable. Además, como hemos mencionado, para estas matrices es posible lograr una mejor utilización de la memoria computacional. Luego para matrices ralas grandes, los métodos iterativos son los más aconsejados por su eficiencia computacional.

Por cuestiones de eficiencia computacional, sólo nos interesarán los *procesos iterativos lineales*. Un procedimiento de iteración *matricial lineal en un punto* tiene la forma general:

$$\underline{x}^{(n+1)} = \underline{B}^{(n)} \cdot \underline{x}^{(n)} + \underline{c}^{(n)} \quad (31)$$

Si \underline{B} y \underline{c} son independientes del nivel n , la iteración matricial se denomina *estacionaria*.

El proceso iterativo continúa hasta que la diferencia (según una norma elegida) entre las aproximaciones obtenidas en dos iteraciones sucesivas sea menor que una cota δ prefijada. Esto es, se itera hasta que:

$$\|\underline{x}^{(n+1)} - \underline{x}^{(n)}\| \leq \delta \quad (32)$$

La notación $\|\bullet\|$ indica la norma elegida. Tanto δ como $\|\bullet\|$ son arbitrarias, adecuándose el criterio de error a las características físicas de la simulación que se realiza.

Normas Usualmente Empleadas

a) Norma Euclideana: $\|\bullet\|_2 \leq \sqrt{\sum_{i=1}^m (x_i^{(n+1)} - x_i^{(n)})^2}$

b) Norma de Máximo (Norma ∞): $\|\bullet\|_\infty \leq \max_{1 \leq i \leq m} |x_i^{(n+1)} - x_i^{(n)}|$

En ambos casos el subíndice indica la i -ésima componente del vector \underline{x} .

Método de Aproximaciones Sucesivas

Recordar que en los métodos iterativos basados en aproximaciones sucesivas (o Wegstein) debe explicitarse la variable independiente. A modo de introducción se puede ver dicha estructura sumando miembro a miembro en la Ecuación (2) el vector de incógnitas \underline{x} , así:

$$\underline{x} + \underline{A} \cdot \underline{x} = \underline{b} + \underline{x}$$

En forma equivalente,

$$(\underline{A} + \underline{I}) \cdot \underline{x} = \underline{b} + \underline{x}$$

de donde:

$$\underline{x} = (\underline{A} + \underline{I}) \cdot \underline{x} - \underline{b}$$

De manera que la ley de recurrencia adopta la siguiente forma:

$$\underline{x}^{(n+1)} = (\underline{A} + \underline{I}) \cdot \underline{x}^{(n)} - \underline{b} \quad (32)$$

que nos da una forma genérica simple de los métodos iterativos estacionarios que son los universalmente usados. De esta forma concluimos, que en este caso, la función iteradora es $\underline{F}(\underline{x}) = (\underline{A} + \underline{I}) \cdot \underline{x} - \underline{b}$, esto es, de acuerdo con la estructura de los esquemas iterativos lineales, resulta: $\underline{B} = (\underline{A} + \underline{I})$ y $\underline{c} = -\underline{b}$.

Convergencia del Método de Aproximaciones Sucesivas

Función Continua de Lipschitz en un Recinto R

Una función vectorial $\underline{F}(\underline{x})$ definida sobre un recinto R es una función continua de Lipschitz si verifica que:

$$\|\underline{F}(\underline{a}) - \underline{F}(\underline{b})\| \leq K \|\underline{a} - \underline{b}\| \quad ; \quad \text{con } K \in (0, 1) \text{ y } \forall \underline{a}, \underline{b} \in R \quad (33)$$

Si $\underline{F}(\underline{x})$ es una función continua de Lipschitz en el recinto R, entonces:

a) Llamando \underline{x}^* a la raíz de: $\underline{x} = \underline{F}(\underline{x})$, luego:

$$\|\underline{x}^{(n+1)} - \underline{x}^*\| = \|\underline{F}(\underline{x}^{(n)}) - \underline{F}(\underline{x}^*)\| \quad (34)$$

Si $\underline{F}(\underline{x})$ es una función continua de Lipschitz, la Ecuación (34) se convierte en:

$$\|\underline{x}^{(n+1)} - \underline{x}^*\| \leq K \|\underline{x}^{(n)} - \underline{x}^*\| \quad (35)$$

Esto es, el error cometido en una iteración es proporcional al error cometido en la iteración anterior. Se dice que un método presenta *convergencia lineal* si satisface la Ecuación (35).

b) Si se aplican en forma secuencial las Ecuaciones (34) y (35) se obtiene:

$$\|\underline{x}^{(n+1)} - \underline{x}^*\| \leq K^{(n+1)} \|\underline{x}^{(0)} - \underline{x}^*\| \quad (36)$$

Si $0 < K < 1$, el error decrecerá en las sucesivas iteraciones y el decrecimiento será mayor cuanto menor sea K , esto es, habrá mayor velocidad de convergencia ya que se requiere un menor número de iteraciones para lograr un error dado.

Convergencia de Funciones Continuas de Lifschitz

Sea R un recinto cerrado que contiene a \underline{x}^* , raíz en el recinto del sistema de ecuaciones $\underline{x} = \underline{F}(\underline{x})$. Cualquiera sea la aproximación inicial $\underline{x}^{(0)} \in R$, la sucesión de vectores generada por la ley de recurrencia:

$$\underline{x}^{(n)} = \underline{F}(\underline{x}^{(n-1)}) \text{ con } n = 1, 2, \dots \quad (37)$$

converge hacia \underline{x}^* siempre que $\underline{F}(\underline{x})$ satisfaga la condición de Lipschitz (Ec. (34)).

Sin embargo, la condición de Lipschitz es difícil de determinar para un problema dado. Una forma alternativa que genera una condición suficiente para la convergencia del método, es trabajar con la matriz Jacobiana, \underline{J} , cuyas componentes se definen como:

$$J_{ij}(\underline{x}) = \frac{\partial F_i(\underline{x})}{\partial x_j} \quad \forall i, j \quad (38)$$

Si las componentes de dicha matriz existen y

$$\|\underline{J}(\underline{x})\| \leq M < 1 \quad \forall \underline{x} \in R \quad (39)$$

queda establecida a través de la Ecuación (39) una condición suficiente para la convergencia del Método de Aproximaciones Sucesivas.

En un dominio unidimensional, la condición definida por la Ecuación (39) se reduce a:

$$|F'(x)| < 1 \quad (40)$$

Condición Necesaria de Convergencia

Una *condición necesaria* para que el método de aproximaciones sucesivas converja $\forall \underline{x}^{(0)} \in R$ es que el radio espectral de $\underline{J}(\underline{x}) \leq 1$. Se define el radio espectral de una matriz \underline{A} a:

$$r(\underline{A}) = \max_k |\lambda_k(\underline{A})| \quad (41)$$

Esto es, el radio espectral de una matriz (cuadrada) es el autovalor de la matriz de máximo valor absoluto.

Método de Jacobi

Escribimos la matriz \underline{A} de coeficientes de la siguiente manera:

$$\underline{A} = \underline{D} + \underline{L} + \underline{U} \quad (42)$$

siendo \underline{D} la matriz diagonal formada con los elementos diagonales de \underline{A} . \underline{L} y \underline{U} son matrices triangulares del tipo inferior y superior, respectivamente, formadas con el resto de los elementos de la matriz \underline{A} y ceros en su diagonal principal. Luego, reemplazamos la descomposición de la matriz \underline{A} (Ec. (42)) en el SEAL (Ecuación (2)), así:

$$\begin{aligned} (\underline{D} + \underline{L} + \underline{U}) \cdot \underline{x} &= \underline{b} \\ \underline{D} \cdot \underline{x} &= -(\underline{L} + \underline{U}) \cdot \underline{x} + \underline{b} \\ \underline{x} &= -\underline{D}^{-1} \cdot (\underline{L} + \underline{U}) \cdot \underline{x} + \underline{D}^{-1} \cdot \underline{b} \end{aligned}$$

Obtenemos la ley de recurrencia:

$$\underline{x}^{(n+1)} = -\underline{D}^{-1} \cdot (\underline{L} + \underline{U}) \cdot \underline{x}^{(n)} + \underline{D}^{-1} \cdot \underline{b} = \underline{D}^{-1} \cdot (\underline{D} - \underline{A}) \cdot \underline{x}^{(n)} + \underline{D}^{-1} \cdot \underline{b} \quad (43)$$

que supone que la diagonal principal no tiene ceros.

De no ser así, si \underline{A} es no singular, se debe permutar filas y columnas para obtener una forma que permita definir \underline{D} . Es deseable que sea de la forma diagonal dominante, esto es, que los elementos de la diagonal sean lo más grande posible respecto a los demás. Luego, el método de Jacobi puede entonces escribirse como:

$$\underline{x}^{(n+1)} = \underline{B} \cdot \underline{x}^{(n)} + \underline{c} \quad (44)$$

donde:

$$\underline{B} = -\underline{D}^{-1} \cdot (\underline{L} + \underline{U}) \quad (45)$$

$$\underline{c} = \underline{D}^{-1} \cdot \underline{b} \quad (46)$$

Por componentes:

$$x_i^{(n+1)} = \frac{b_i - \sum_{\substack{k=1 \\ k \neq i}}^n a_{ik} x_k^{(n)}}{a_{ii}} \quad (47)$$

Convergencia del Método de Jacobi

Siendo del Método de Jacobi un caso particular del Método de Aproximaciones Sucesivas, la condición necesaria y suficiente para su convergencia es que la función iteradora:

$$\underline{F}(\underline{x}) = -\underline{D}^{-1} \cdot (\underline{L} + \underline{U}) \cdot \underline{x} + \underline{D}^{-1} \cdot \underline{b} = \underline{B} \cdot \underline{x} + \underline{c} \quad (48)$$

sea una función continua de Lipschitz, lo que puede ser explicitado exigiendo que:

$$\|\underline{B}\|_{\infty} \leq \max_{1 \leq i \leq m} \sum_{j=1}^n |b_{ij}| \leq K < 1 \quad (49)$$

Lo cual significa sumar todos los elementos de una fila y elegir el máximo entre ellos.

Si se analiza la expresión de \underline{B} (Ecuación (48)), se observa que la expresión (49) es equivalente a:

$$\|\underline{B}\|_{\infty} \leq \max_i \sum_{\substack{j=1 \\ j \neq i}}^n \frac{|a_{ij}|}{|a_{ii}|} < 1 \quad (50)$$

Una expresión similar a la (50) puede generarse para el Método Gauss-Seidel.

Se concluye que el método es convergente solamente cuando se aplica a matrices cuya matriz de coeficientes sea o pueda ser transformada mediante intercambio de filas y/o columnas en diagonal dominante, lo que no es habitual en simulación de procesos químicos.

Método de Gauss-Seidel

Este es el método ampliamente usado que siempre converge si converge Jacobi, (y aún en casos en que éste no converge); y en general lo hace más rápido, siendo su implementación más eficiente.

La diferencia fundamental es que a medida que se calcula cada componente de $\underline{x}^{(n+1)}$ se utiliza inmediatamente en la misma iteración. Al proceder de ese modo, la ecuación obtenida para Jacobi, resulta:

$$\underline{x}^{(n+1)} = -(\underline{D}^{-1} \cdot \underline{L}) \cdot \underline{x}^{(n+1)} - (\underline{D}^{-1} \cdot \underline{U}) \cdot \underline{x}^{(n)} + \underline{D}^{-1} \cdot \underline{b} \quad (51)$$

reordenando y sacando factor común $\underline{x}^{(n+1)}$ se tiene:

$$(\underline{I} + \underline{D}^{-1} \cdot \underline{L}) \underline{x}^{(n+1)} = -(\underline{D}^{-1} \cdot \underline{U}) \cdot \underline{x}^{(n)} + \underline{D}^{-1} \cdot \underline{b}$$

por lo que, premultiplicando por \underline{D} , obtenemos:

$$(\underline{D} + \underline{L}) \cdot \underline{x}^{(n+1)} = -\underline{U} \cdot \underline{x}^{(n)} + \underline{b}$$

Luego resulta:

$$\underline{x}^{(n+1)} = -(\underline{D} + \underline{L})^{-1} \cdot \underline{U} \cdot \underline{x}^{(n)} + (\underline{D} + \underline{L})^{-1} \cdot \underline{b} \quad (52)$$

entonces la expresión generalizada será:

$$\underline{x}^{(n+1)} = \underline{B} \cdot \underline{x}^{(n)} + \underline{c} \quad (53)$$

donde:

$$\underline{B} = -(\underline{D} + \underline{L})^{-1} \cdot \underline{U} \quad (54)$$

y

$$\underline{c} = (\underline{D} + \underline{L})^{-1} \cdot \underline{b} \quad (55)$$

Por componentes, resulta:

$$x_i^{(n+1)} = \frac{b_i - \sum_{k=1}^{i-1} a_{ik} x_k^{(n+1)} - \sum_{k=i+1}^n a_{ik} x_k^{(n)}}{a_{ii}} \quad (56)$$

Existe un teorema que prueba que: *Si una matriz es definida positiva, el procedimiento de Gauss-Seidel converge independientemente del vector inicial propuesto.*

Otro hecho importante es que los métodos matriciales iterativos convergen linealmente. Es natural pensar que el error de redondeo es mayor en un método iterativo que en uno directo; sin embargo como siempre se usa la matriz de coeficientes original, el error por redondeo en que se incurre en un método iterativo es sólo aquel producido en la última iteración. En consecuencia, ambos tipos de métodos tienen un error equivalente, tan serio para unos como para otros.

El error de truncamiento, en el caso de los métodos iterativos, suele ser de un orden mayor de magnitud que el mencionado de redondeo. No obstante, dado que en los métodos iterativos se tiene un criterio de control de error, este es acotado convenientemente.

Método de Sobrerrelajación

Una condición necesaria y suficiente para que el método de Jacobi (y eventualmente el de Gauss-Seidel), planteado a través del algoritmo $\underline{x}^{(n+1)} = \underline{B} \cdot \underline{x}^{(n)} + \underline{c}$ converja, es que el máximo autovalor de \underline{B} sea en valor absoluto menor que 1 (radio espectral menor que 1).

Cuanto menor sea dicho radio espectral, más rápida será la convergencia ya que se puede demostrar que:

$$\frac{\|\underline{x}^{(n)} - \underline{x}^*\|}{\|\underline{x}^{(0)} - \underline{x}^*\|} \leq |\lambda|_{m\acute{a}x}^n \quad (57)$$

El método de sobrerrelación pretende modificar la matriz B y consecuentemente su

autovalor de manera de acelerar la convergencia. Para ello modifica la Ec. (44) conforme a:

$$\underline{x}^{(n+1)} = \underline{B} \cdot \underline{x}^{(n)} + \underline{c} \quad (58)$$

Siendo los valores estimados para una nueva iteración:

$$\underline{x}^{(n+1)} = \alpha \underline{x}^{(n+1)} + (1-\alpha) \underline{x}^{(n)} \quad (59)$$

Si se sustituye la Ec. (58) en la Ec. (59) queda:

$$\underline{x}^{(n+1)} = [\alpha \underline{B} + (1-\alpha) \underline{I}] \underline{x}^{(n)} + \alpha \underline{c} \quad (60)$$

Siendo ahora el máximo autovalor en valor absoluto de la matriz:

$$\underline{B}' = \alpha \underline{B} + (1-\alpha) \underline{I} \quad (61)$$

el que determinará la velocidad de convergencia. Pero, la ecuación característica de \underline{B}' vendrá expresada por:

$$\det(\underline{B}' - \lambda' \underline{I}) = \det \left\{ \alpha \left[\underline{B} + \frac{(1-\alpha-\lambda')}{\alpha} \underline{I} \right] \right\} = \alpha^n \det \left[\underline{B} + \frac{(1-\alpha-\lambda')}{\alpha} \underline{I} \right] = 0 \quad (62)$$

Si definimos:

$$-\lambda = \frac{1-\alpha-\lambda'}{\alpha} \quad (63)$$

la Ecuación (62) también dará los autovalores de \underline{B} , entonces:

$$\lambda'_i = 1 + \alpha(\lambda_i - 1) \quad (64)$$

Adoptando $\alpha > 1$ para $\lambda_i > 0$ es posible reducir $|\lambda'|_{m\acute{a}x}$ y consecuentemente acelerar la convergencia del método.

El método de relajación puede también aplicarse al método de Gauss-Seidel y una discusión del mismo tenor que la precedente puede llevarse a cabo en dicha circunstancia.