

# CAPÍTULO III

## REVISIÓN DE MÉTODOS NUMÉRICOS APLICABLES EN SIMULACIÓN DE PROCESOS EN ESTADO ESTACIONARIO

**Por**  
**H. J. Espinosa, P. Aguirre y G. A. Pérez**

### III.1 CONCEPTOS BÁSICOS

La solución de una ecuación no lineal o de una función trascendente, como método, supone la búsqueda de un valor tal que satisfaga la ecuación o función en cuestión.

Si  $f(x)$  es una función no lineal genérica y  $x^*$  es la solución o raíz, entonces se cumple que

$$f(x^*) = 0 \quad (1)$$

Otra forma alternativa de formular el problema es:

$$x^* = F(x^*) \quad (2)$$

donde  $F(x)$  es una función diferente de la original, y supone que puede explicitarse de la primitiva, cumpliéndose para ello el teorema de funciones implícitas.

Una *solución iterativa* significa comenzar con un valor inicial,  $x_0$ , y generar una sucesión (secuencia)  $x_0, x_1, x_2, \dots, x_n$ , tal que:

$$\{x_n\} / \lim_{n \rightarrow \infty} x_n = x^* \quad (3)$$

donde  $n$  representa el número de iteración (término de la correspondiente sucesión).

El *error exacto en la iteración  $n$*  es:

$$E_n = |x_n - x^*| \quad (4)$$

Las hipótesis usuales son que  $f(x)$ ,  $F(x)$  y  $x^*$  satisfacen las siguientes restricciones:

- i.-  $x^*$  se sitúa dentro de:  $I = [a, b] / f(x) \wedge F(x) \in C_I$ .

**Modelado, Simulación y Optimización de Procesos Químicos**

Autor: Nicolás J. Scenna y col.

ISBN: 950-42-0022-2 - ©1999

ii.- Hay una sola raíz (única) en  $I$  y es real ( $x^* \in \mathbb{R}$ ).

donde el símbolo  $\wedge$  representa el operador lógico y  $C_j$ : funciones continuas.

A pesar de que no se diga, todos los métodos usuales se basan en dichas hipótesis. De no ser así se requerirán métodos especiales.

### III.2 MÉTODOS BÁSICOS. DISCUSIÓN DE LA CONVERGENCIA

En general, los métodos de mayor *orden de convergencia* llegan a la solución en un número menor de iteraciones. Las excepciones se producen al inicializarlos (o proponer las primeras estimaciones) con valores numéricos muy inapropiados. De todos modos, el menor número de iteraciones no necesariamente significa el menor tiempo computacional, dado que éste depende del esfuerzo de cálculo involucrado en cada iteración. El cálculo de  $f(x)$  es el mayor consumidor de tiempo junto con el cálculo de su derivada, *que para el caso de una ecuación o función única pueden suponerse equivalente*. El número de estas evaluaciones suele ser una medida más adecuada de la eficiencia del método de solución.

Los métodos iterativos hallan la solución exacta (si eso ocurre) *sólo en un número infinito de iteraciones* (Ecuación (3)). En la práctica, las iteraciones se detienen cuando el error es menor que una adecuada tolerancia, impuesta por el usuario. El valor exacto del error (Ecuación (4)) no puede usarse como criterio de terminación, porque normalmente no se conoce.

Por lo tanto, como criterio de finalización se usa una estimación del denominado error exacto, y es así como interviene la tolerancia de error adoptada,  $E_d$ .

Los criterios generalmente usados son:

$$|f(x_n)| < E_d \quad (5.1)$$

$$|x_n - x_{n-1}| < E_d \quad (5.2)$$

$$|x_n - x_{n-1}| < E_d |x_n| \quad (5.3)$$

Los dos primeros comparan errores absolutos, mientras que el tercero analiza errores relativos. Estos criterios pueden tener limitaciones, como veremos, y existen propuestas alternativas que resuelven el problema.

Una comparación de los métodos más usados puede verse en la tabla siguiente:

**Tabla III.1: Métodos más usuales.**

MÉTODO	ORDEN DE CONVERGENCIA	INFORMACIÓN PARA CALCULAR $x_{n+1}$
BISECCIÓN	LINEAL (1)	$f(x_n), f(x_{n+1})$
SECANTE	SUPERLINEAL (1.618)	$f(x_n), f(x_{n+1})$
NEWTON-RAPHSON	CUADRÁTICO (2)	$f(x_n), f'(x_n)$
SUSTITUCIÓN DIRECTA	LINEAL (1)	$F(x_n)$

Evidentemente, en nuestro campo de aplicaciones importan las raíces reales de  $f(x) = 0$ . El hecho de analizar una función trascendente primero (o una función o ecuación no lineal) es debido a la relación con los sistemas de ecuaciones, tanto lineales como no lineales, y también porque introduce y educa en la generación de *algoritmos*, y en los problemas asociados con la convergencia: factibilidad y velocidad.

Un concepto básico asociado con los métodos iterativos es el llamado *orden de convergencia* del método. Si definimos el error en la iteración  $n$ , como  $E_n$ , entonces si existe un número real  $p \geq 1$  tal que:

$$\lim_{n \rightarrow \infty} \frac{|x_{n+1} - x^*|}{|x_n - x^*|^p} \equiv \lim_{n \rightarrow \infty} \frac{|E_{n+1}|}{|E_n|^p} = K \neq 0 \quad (6)$$

se dice que *el método es de orden  $p$* , en  $x^*$ . La constante  $K$  se llama *constante de error asintótica*, y depende de  $f(x)$ . A mayor orden de convergencia, el método convergirá a mayor velocidad. Pero no implica, obviamente, garantía de convergencia.

Otra medida usual es aquella que se hace para saber cuánto se debe computar, para alcanzar una cierta precisión en la raíz buscada. En este caso, además del orden interviene el *costo computacional por iteración* que es importantísimo para definir la *eficiencia* de un método.

### III.3 PRINCIPALES MÉTODOS

#### III.3.1 El Método de Bisección

Es interesante considerarlo porque naturalmente es lo primero que uno haría, si fuera capaz de graficar (o conocer la gráfica) de la función  $f(x)$  en cuestión. Además, su basamento se puede probar, y dentro de ciertas pautas da garantía de éxito.

Si las raíces de interés práctico son aisladas, esto es los valores de  $x$  han sido divididos en distintos intervalos, en cada uno de los cuales sospechamos que se encuentra alguna raíz de importancia para nuestros fines, posibilita garantía de convergencia a la misma. Existe un teorema, que dice: *Si  $f(x)$  es continua desde  $x=a$  hasta  $x=b$ , y si  $f(a)$  y  $f(b)$  tienen signos opuestos, luego hay como mínimo una raíz*

real de  $f(x) = 0$ , entre  $a$  y  $b$ .

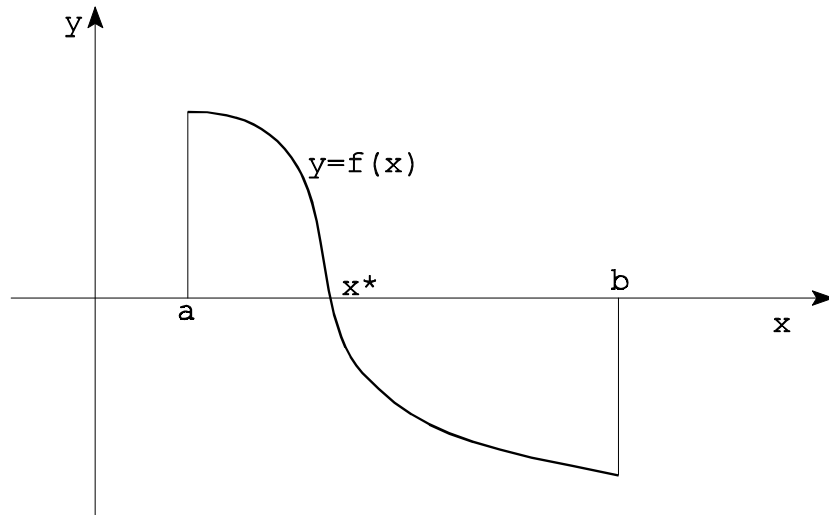


Figura III.1: Raíz aislada de una ecuación.

Basado en el teorema anterior, supongamos (tal como se aprecia en la Fig. (III.1)) que  $f(x)$  es continua y positiva en  $x = a$  y negativa en  $x = b$ . Luego, habrá una raíz entre  $a$  y  $b$ . Si, entonces, dividimos el intervalo  $I = [a, b]$  por la mitad, podemos calcular  $f[(a+b)/2]$ , teniendo tres posibilidades, a saber:

- i.- Que sea cero (dentro de cierto error), en cuyo caso es la raíz.
- ii.- Que sea negativo, en cuyo caso la raíz está entre  $x = (a+b)/2$  y  $x = a$ .
- iii.- Que sea positivo, entonces se encuentra entre  $x = (a+b)/2$  y  $x = b$ .

De modo que el procedimiento continúa bisectando el intervalo cada vez, hasta lograr la raíz dentro de la precisión deseada.

Cada repetición reduce el error máximo por un factor de 2 (dos), de manera que tres iteraciones producen una mejora aproximada de un orden de magnitud. Su gran virtud es que *asegura la convergencia*, que constituye una propiedad que no encontramos en los otros métodos. A cambio de ello, puede ser tremendamente lento lo que lo puede transformar en ineficiente, debido a que es lineal (orden 1).

Otro método que converge para toda función continua es el denominado *Método de la Falsa Posición o Regula Falsi* (ver Figura (III.2)).

*Mecánica* (algoritmo)

$$x_1 \text{ y } x_2 / f(x_1) \cdot f(x_2) < 0$$

luego  $\rightarrow x_3$

si:  $f(x_3)f(x_1) < 0 ; i = 1, 2$

(en este caso  $i = 1$ )

se repite  $x_4$  y continúa.

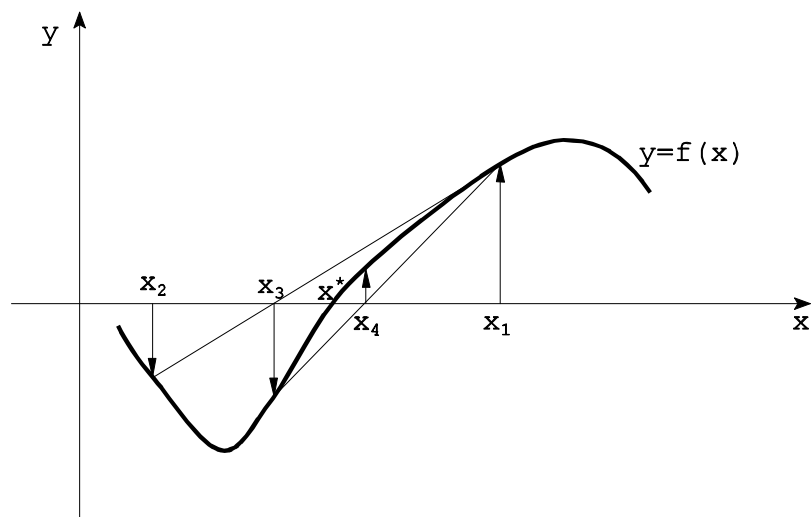


Figura III.2: Método Regula Falsi.

Se ve que:

$$x_3 = \frac{y_2}{y_2 - y_1} x_1 + \frac{y_1}{y_1 - y_2} x_2 \quad (7.1)$$

$$x_4 = \frac{y_3}{y_3 - y_1} x_1 + \frac{y_1}{y_1 - y_3} x_3 \quad (7.2)$$

Si  $f(x)$  es cóncava entre  $x_1$  y  $x_2$  se dice que el método es *estacionario*; esto es, el punto  $x_1$  es siempre uno de los dos puntos usados para la siguiente iteración (ver Figura (III.3)). Lo mismo ocurriría si fuese convexa en las inmediaciones de la raíz, provocando una convergencia lineal en estos casos.

Una mejora que lo hace más eficiente consiste en aplicar la fórmula hallada (Ecuación (7.2)) a los puntos  $x_{i-1}$  y  $x_{i+1}$ , pero reemplazando

$$y_{i-1} \text{ por } \alpha y_{i-1} = \bar{y}_{i-1}, \text{ tal que } 0 < \alpha < 1 \quad (8)$$

Esto es:

$$x_{i+2} = \frac{\alpha y_{i-1}}{\alpha y_{i-1} - y_{i+1}} x_{i+1} + \frac{y_{i+1}}{y_{i+1} - \alpha y_{i-1}} x_{i-1} \quad (9)$$

Si cambia el signo debe retornarse el método original, de lo contrario se continúa con la mejora presentada. Las elecciones más simples para  $\alpha$  son:

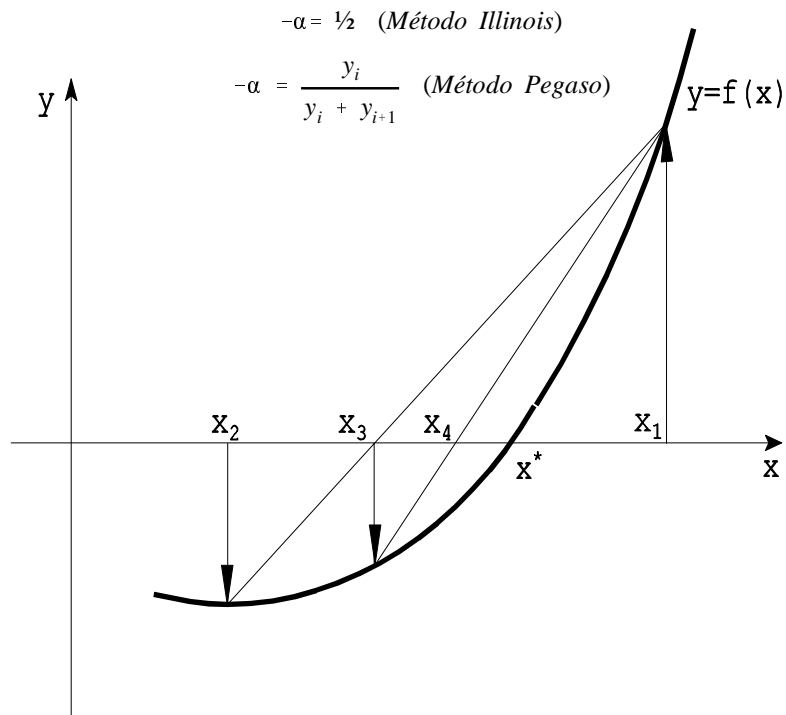


Figura III.3: Método Regula Falsi para funciones cóncavas.

Estas mejoras incrementan el orden, y los hace superlineales (Illinois: 1.44 y Pegaso: 1.64). Son algoritmos muy buenos cuando tenemos dos puntos en los cuales  $f(x)$  tiene signo opuesto. En caso de no querer hacer un esfuerzo computacional en hallarlos, se pueden usar *interpolación lineal inversa* con los dos últimos puntos computados para generar el siguiente. En este caso resulta:

$$x_{i+1} = \frac{y_i}{y_i - y_{i-1}} x_{i-1} + \frac{y_{i-1}}{y_{i-1} - y_i} x_i \quad (10)$$

A este método estacionario se lo denomina *método de la secante*. Requiere de dos aproximaciones iniciales (recordar Tabla (III.1)), y sólo si son cercanas a la raíz el método es convergente. Pero tiene la ventaja de ser super lineal (orden de convergencia = 1.618). Su nombre se debe, según la Ecuación (10), a que se puede escribir:

$$x_{i+1} = x_i - \frac{x_i - x_{i-1}}{y_i - y_{i-1}} y_i \quad (11)$$

expresión que será de utilidad cuando discutamos métodos para sistemas de ecuaciones.

En la Ecuación (11) se pone de manifiesto un problema típico de errores en cálculo numérico. Se ve que cerca de la solución  $y_i$  e  $y_{i-1}$  son cantidades parecidas con lo que aparece el *problema de cancelación sustractiva* en ese segundo término; que es un término de corrección que estará aportando muy pocos dígitos significativos, lo que obliga a usar doble precisión (término computacional que implica la cantidad de dígitos retenidos) si queremos acercarnos a  $x^*$  con mucha precisión.

### III.3.2 Métodos de Newton-Raphson. Usos de la Derivada de la Función

Si basados en la Ecuación (11), expresión de la secante generalizada, modificamos el término corrector de modo que aparezca la pendiente en el punto, obtendremos el conocido método de Newton-Raphson. Es decir, si:

$$x_{i+1} = x_i - \frac{y_i}{\frac{y_i - y_{i-1}}{x_i - x_{i-1}}} = x_i - \frac{f(x_i)}{\frac{f(x_i) - f(x_{i-1})}{x_i - x_{i-1}}} \quad (12)$$

se puede hacer que:

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} \quad (13)$$

expresión del algoritmo que nos interesa.

Otra forma de entender la propuesta es suponer una expansión en serie de Taylor, truncada en el primer término (linealización de la función). Esto es:

$$f(x_{i+1}) - f(x_i) \approx f'(x_i)(x_{i+1} - x_i) \dots \quad (14)$$

$y = f(x_{i+1})$  es tal que es un valor muy cercano a  $x^*$ , tal que si:  $y_{i+1}/f(x_{i+1}) = 0$ , tenemos:

$$f(x_i) - f'(x_i)(x_{i+1} - x_i) \quad (15)$$

que conduce, reordenando, a la Ecuación (13) básica del método.

Su ventaja es el orden de convergencia ( $p = 2$ , cuadrática) y su desventaja principal es la evaluación de la derivada de la función, ya sea analítica o numérica, como se verá en los casos de interés práctico.

**Ejemplo:**

Sea la siguiente función, que por lo sencilla pueden conocerse sus raíces, de tal forma de ilustrar la mecánica del método:  $f(x) = x^2 - 3x = 0$ , con raíces 0 y 3. Encontrar una de sus raíces por N-R. Para ello como vimos necesitamos un valor inicial supuesto, o semilla. Sea éste  $x_0 = 1$ , y el criterio de tolerancia:  $\mathcal{J}(x) \neq 10^{-3}$ .

$f(x_0) = -2$ , distinto de cero, luego:

$$x_1 = x_0 + \frac{f(x_0)}{f'(x_0)}$$

$f'(x) = 2x - 3$ , luego:

$$x_1 = 1 - (-2) / (-1) = -1.$$

$f(x_1) = 1 + 3 = 4 \dots 0$ , luego:

$$x_2 = x_1 + \frac{f(x_1)}{f'(x_1)} = (-1) + 4 / (-1) = -5$$

$f(-5) = 0.64, \dots 0$ , luego:

$$x_3 = x_2 + \frac{f(x_2)}{f'(x_2)} = (-5) + (0.64) / (-3.4) = -5.188$$

$f(x_3) = -0.0354$

$$x_4 = x_3 + \frac{f(x_3)}{f'(x_3)} = (-5.188) + (0.0354) / (-3.0235) = -5.277 \cdot 10^{-5}$$



ello implica que  $x_i$  es raíz de la función dada, dentro del margen de tolerancia especificado.

### III.3.3 Sustitución Directa o Aproximaciones Sucesivas

Dada una función o una ecuación, si puede resolverse en forma explícita para una variable, es decir si:

$$\text{Dado: } f(x) = 0$$

$$\text{Se propone: } x = F(x)$$

esto es, explicitando la variable independiente, se puede, entonces, establecer la fórmula para un *algoritmo de un solo punto y estacionario*, es decir:

$$x_{i+1} = F(x_i) \tag{16}$$

Su éxito dependerá, evidentemente, del arreglo logrado para la ecuación (esto es:

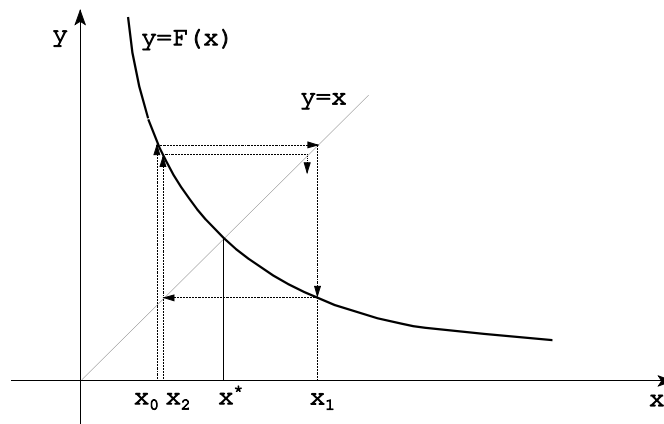


Figura III.4: Método de aproximaciones sucesivas.

$F(x)$ , cuyo análisis permite ver la factibilidad de solución. Esto es, de acuerdo a la Figura (III.4), cuando es convergente (estabilidad del método). Se puede probar -y ver gráficamente- que si:

$$|F'(x)| < 1 \quad ; \quad x \in I \setminus x^{\text{fl}}$$

es convergente.

**Ejemplo:**

Sea la misma función del ejemplo anterior. Encontrar una raíz por el método de sustitución directa. Para explicitar  $x$  operamos algebraicamente, y obtenemos:

$$x = \sqrt{3x} = F(x), \quad \text{con } x \geq 0$$

Hacemos  $x_0 = 1$ , criterio de error ( $5 \cdot 10^{-3}$ )

$x_{i+1} = F(x_i)$ , entonces:

$$x_1 = F(x_0) = 3^{1/2} = 1.732051$$

$$x_2 = F(x_1) = 2.279510$$

$$x_3 = F(x_2) = 2.615060$$

$$x_4 = F(x_3) = 2.800923$$

$$x_5 = F(x_4) = 2.898753$$

$$x_6 = F(x_5) = 2.948942$$

$$x_7 = F(x_6) = 2.974361$$

$$x_8 = F(x_7) = 2.987153$$

$$x_9 = F(x_8) = 2.993570$$

$$x_{10} = F(x_9) = 2.996783$$

$|x_{10} - x_9| = 3.2 \cdot 10^{-3}$ , luego,  $x_{10}$  es raíz de  $f(x)$  dentro del margen de error especificado.

$$f(x_{10}) = 0.0096 \neq 0$$

**Aceleradores de Convergencia: Casos  $p=1$**

Son técnicas basadas en extrapolaciones. Si la iteración es convergente, será:

$$x^{\text{fl}} - x_{i+1} = K_i (x^{\text{fl}} - x_i); \quad |K_i| < 1 \quad (17)$$

donde:  $|K_i| \approx K$ : constante de error asintótico.

Cerca de la convergencia  $K_i$  es prácticamente constante (y parecido al valor final de  $K$ ). De modo que:

$$(x^{\text{fl}} - x_{i+1}) = K (x^{\text{fl}} - x_i) \quad (18)$$

y eliminando  $K$  entre dos iteraciones

$$\frac{(x^* - x_{i+2})}{(x^* - x_{i+1})} \cong \frac{(x^* - x_{i+1})}{(x^* - x_i)} \quad (19)$$

resolviendo para  $x^*$  se obtiene:

$$x^* \cong \frac{x_i x_{i+2} - x_{i+1}^2}{x_{i+2} - 2 x_{i+1} + x_i} \quad (20)$$

Este resultado es el *procedimiento de Aitken*, de aceleración de convergencia.

De manera que:  $x_{i+3} \cong x^*$  es una mejor aproximación (valor extrapolado); y así se continúa con el proceso iterativo, que logra mejorar la performance.

### III.3.4 Procedimiento de Wegstein

Es el algoritmo más utilizado para acelerar el método de aproximaciones sucesivas. Incluso, es de gran importancia su implementación en problemas de sistemas de ecuaciones no lineales, como veremos en el próximo capítulo.

La base del mismo es proponer a la clásica iteración de aproximaciones sucesivas un valor *mejorado*, según la siguiente ecuación:

$$\bar{x}_{i+1} = q \bar{x}_i + (1 - q) x_{i+1} \quad (21)$$

de modo que:

$$x_{i+2} = F(\bar{x}_{i+1})$$

se corrige  $x_{i+2}$  y continúa. De la Ecuación (21) se ve que es necesario generar dos valores según el esquema tradicional, y conociendo  $q$  comenzar con esta propuesta.

Para analizar el cálculo de  $q$  puede ser útil la Figura (III.5), donde se aprecia la aproximación propuesta.

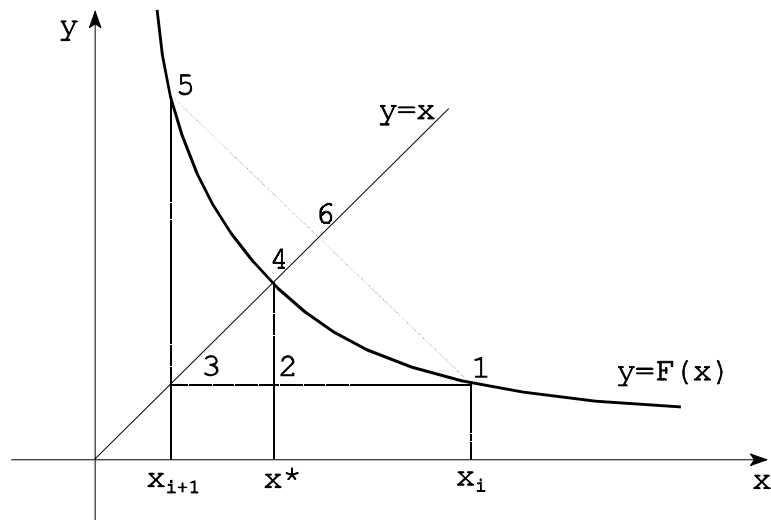


Figura III.5: Método de Wegstein.

Teniendo como idea que:

$$\bar{x}_{i+1} \approx x^*$$

es la definición de  $q$ , según:

$$x^* = q x_i + (1 - q) x_{i+1}$$

o sea

$$x^* - x_{i+1} = q (x_i - x_{i+1}) \quad (22)$$

y, de igual modo,

$$x^* - x_i = (1 - q) (x_{i+1} - x_i) \quad (23)$$

resultando

$$\frac{x^* - x_{i+1}}{x^* - x_i} \equiv - \frac{q}{(1 - q)} \quad (24)$$

Como  $x^*$  no se conoce, se debe aproximar. Esto es:

$$\frac{q}{(1 - q)} = \frac{\overline{23}}{12} = \frac{\overline{24}}{12} \cong \frac{\overline{35}}{13}$$

donde los segmentos indicados en los cocientes están determinados en Figura (III.5).  
Luego:

$$\frac{q}{(1 - q)} \cong \frac{x_{i+2} - x_{i+1}}{x_i - x_{i+1}} \quad (25)$$

de donde obtenemos la expresión para el cálculo de q:

$$q \cong \frac{x_{i+2} - x_{i+1}}{x_i - 2 x_{i+1} + x_{i+2}} \quad (26)$$

En el capítulo siguiente veremos ejemplos de aplicación tanto para una sola ecuación como para sistemas de ecuaciones no lineales, resueltas por este método.

### III.3.5 Uso de Fracciones Continuas

Para problemas de una variable la representación a través de fracciones continuadas ha demostrado tener ciertas ventajas sobre otros métodos tradicionales. La principal es el reducido número de operaciones que requiere su ejecución, con la correspondiente disminución del tiempo de computación.

Este tipo de representación ha sido usado para problemas de interpolación de datos; y se puede extender al caso de búsqueda de raíces que nos interesa.

Si se conocen valores tabulados de una función, es decir:  $(x_0, f_0), (x_1, f_1), \dots, (x_n, f_n)$ , para aproximar  $f(x)$  se puede proponer la siguiente expresión:

$$f(x) \cong \varphi_0(x) = a_0 + \frac{x - x_0}{a_1 + \frac{x - x_1}{a_2 + \frac{x - x_2}{a_3 + \frac{x - x_3}{\dots \dots \dots a_{n-1} + \frac{x - x_{n-1}}{a_n}}}} \quad (27)$$

Que, en forma compacta, puede ser escrita como:

$$f(x) \cong \varphi_0(x) = a_0 + \frac{x - x_0}{\varphi_1(x)} \quad (28)$$

donde:

$$\begin{aligned} \varphi_1(x) &= a_1 + \frac{x - x_1}{\varphi_2(x)} \\ \varphi_2(x) &= a_2 + \frac{x - x_2}{\varphi_3(x)} \\ &\dots\dots\dots \\ \varphi_{i-1}(x) &= a_{i-1} + \frac{x - x_{i-1}}{\varphi_i(x)} \\ &\dots\dots\dots \\ \varphi_n(x) &= a_n \end{aligned} \quad (29)$$

Los coeficientes  $a_i$  pueden evaluarse con la Ecuación (29). Esto es, si  $\varphi_0(x_i)$  se iguala a los valores  $f_i$ , luego se resuelven las ecuaciones que resultan. Es decir: en  $x = x_0$

$$\varphi_0(x_0) = f(x_0) \equiv f_0 = a_0 \quad (30)$$

en  $x = x_1$

$$\varphi_0(x_1) = f_1 = f_0 + \frac{x_1 - x_0}{a_1}$$

luego,

$$a_1 = \frac{x_1 - x_0}{f_1 - f_0} \quad (31)$$

De igual forma se pueden obtener todos los coeficientes.

En el caso de búsqueda de raíces, es decir resolver para los valores de  $x$  tal que:  $f(x) = 0$ , se plantea

$$x = \psi(f) \quad (32)$$

ecuación de la función inversa. Si ésta se aproxima por fracciones continuadas,

$$x = \psi_0(f) = a_0 + \frac{f - f_0}{a_1 + \frac{f - f_1}{a_2 + \frac{f - f_2}{a_3 + \frac{\dots}{a_{n-1} + \frac{f - f_{n-1}}{a_n}}}}} \quad (33)$$

Luego, sustituimos  $f = 0$  en Ecuación (33), de modo que el lado izquierdo sea una aproximación a la raíz buscada,

$$\bar{x} = a_0 - \frac{f_0}{a_1 - \frac{f_1}{a_2 - \frac{f_2}{a_3 - \frac{\dots}{a_{n-1} - \frac{f_{n-1}}{a_n}}}}} \quad (34)$$

La determinación de los coeficientes  $a_i$  se lleva a cabo en una forma similar a la ya vista. Esto es, si:

$$x = \Omega_0(f) = a_0 + \frac{f - f_0}{\Omega_1(f)} \quad (35)$$

donde:

$$\begin{aligned} \Omega_1(f) &= a_1 + \frac{f - f_1}{\Omega_2(f)} \\ &\dots\dots\dots \\ \Omega_{i-1}(f) &= a_{i-1} + \frac{f - f_{i-1}}{\Omega_i(f)} \\ &\dots\dots\dots \\ \Omega_n(f) &= a_n \end{aligned} \quad (36)$$

la evaluación se hace con la Ecuación (36), poniendo  $\Omega_0(f_i)$  igual a  $x_i$  y resolviendo

sucesivamente. En  $f = f_0$ ,  $\psi(f_0) = a_0$  (Ecuación 33).

Un algoritmo muy interesante, basado en este método, es el de Shacham (1989). Se denomina de memoria mejorada y se aplica a cualquier ecuación no lineal.

El reemplazo de los polinomios de interpolación de Lagrange por fracciones continuadas mejora los métodos de interpolación inversa, sobre todo en el uso de memoria que requieren. El mencionado algoritmo, para el caso de una ecuación del tipo  $f(x) = 0$ , se puede describir de la siguiente manera:

- 1.- Elija dos valores iniciales,  $x_0$  y  $x_1$ . Calcule  $y_0 = f(x_0)$ ,  $y_1 = f(x_1)$ .
- 2.- Compute  $x_2$ , con las siguientes ecuaciones:  
 $a_0 = x_0$   
 $a_1 = (y_1 - y_0)/(x_1 - x_0)$   
 $x_2 = a_0 - y_0/a_1$
- 3.- Haga  $n = 2$  y calcule  $y_2 = f(x_2)$ .
- 4.- Con las ecuaciones recursivas se calcula  $x_{n+1}$ , de modo que:  
 $b_0 = x_n$   
 $b_i = (y_n - y_{i-1})/(b_{i-1} - a_{i-1}) ; i = 1, 2, 3, \dots, n-1$   
 $\psi_n = a_n = (y_n - y_{n-1})/(b_{n-1} - a_{n-1})$   
 $\psi_{i-1} = a_{i-1} - y_{i-1}/\psi_i ; i = n, n-1, \dots, 1$   
 $x_{n+1} = \psi_0$
- 5.- Calcule  $y_{n+1} = f(x_{n+1})$ .
- 6.- Verifique la convergencia. Si no converge, haga  $n = n+1$ , y vuelva a 4.

### III.4 SOLUCIÓN DE SISTEMAS DE ECUACIONES LINEALES SIMULTÁNEAS

#### III.4.1 Planteo del Problema. Teoremas Básicos

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

$$\sum_{j=1}^n a_{ij} x_j = b_i ; i = 1, 2, \dots, n \quad (37)$$

o, en su forma compacta matricial

$$\underline{A} \cdot \underline{x} = \underline{b} \quad (38)$$

donde  $\underline{A} = [a_{ij}]$  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  $\underline{A}$  como  $\underline{b}$  son 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:

$$\underline{\underline{A}}_b = [\underline{\underline{A}} \quad \underline{b}] ; n \cdot (n + 1) \quad (39)$$

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

- 1.- El sistema de ecuaciones (Ecuación(38)) tiene solución *sí y sólo sí*:  $v(\underline{\underline{A}}) = v(\underline{\underline{A}}_b)$ .
- 2.- Si  $v(\underline{\underline{A}}) = v(\underline{\underline{A}}_b) = k < n$ , luego las  $x_{i1}, x_{i2}, \dots, x_{ik}$  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  $v(\underline{\underline{A}}) = v(\underline{\underline{A}}_b) = n$ , 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  $v(\underline{\underline{A}}) < n$ .

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. Estas matrices son típicas al resolver problemas con ecuaciones diferenciales parciales y también aparecen al plantear modelos de simulación de plantas completas, o bien de procesos con múltiple etapas en serie, como veremos en los próximos capítulos.

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  $\mathbf{r}$ , cuando se tiene una solución calculada  $\mathbf{x}_c$ . Es decir:

$$\mathbf{r} = \underline{\mathbf{b}} - \underline{\mathbf{A}} \cdot \mathbf{x}_c \quad (40)$$

se sabe que:

$$\underline{\mathbf{A}} \cdot \mathbf{x}^* = \underline{\mathbf{b}}, \text{ o } \underline{\mathbf{b}} - \underline{\mathbf{A}} \cdot \mathbf{x}^* = 0 \quad (41)$$

entonces será:

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

luego,

$$(\mathbf{x}^* - \mathbf{x}_c) = \underline{\mathbf{A}}^{-1} \cdot \mathbf{r} \quad (42)$$

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

### III.4.2 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 (42), se obtiene el sistema en la siguiente forma:

$$\begin{aligned} a_{11} x_1 + a_{12} x_2 + \dots + a_{1n} x_n &= b_1 \equiv a_{1, n+1} \\ a_{21} x_1 + a_{22} x_2 + \dots + a_{2n} x_n &= b_2 \equiv a_{2, n+1} \\ &\dots\dots\dots \\ &\dots\dots\dots \\ a_{n1} x_1 + a_{n2} x_2 + \dots + a_{nn} x_n &= b_n \equiv a_{n, n+1} \end{aligned} \quad (43)$$

Se supone que la matriz es no singular, y que  $a_{i1} \neq 0$  de manera de poder dividir la primera columna por  $a_{i1}$  y así restar para las ecuaciones, donde  $i=2, \dots, n$ . Esto da como resultado:

$$\begin{aligned}
 a_{11} x_1 + a_{12} x_2 + \dots + a_{1n} x_n &= a_{1, n+1} \\
 a_{22}^{(1)} x_2 + \dots + a_{2n}^{(1)} x_n &= a_{2, n+1}^{(1)} \\
 &\dots \\
 &\dots \\
 a_{n2}^{(1)} x_2 + \dots + a_{nn}^{(1)} x_n &= a_{n, n+1}^{(1)}
 \end{aligned}
 \tag{44}$$

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

$$a_{ij}^{(1)} = a_{ij} - \frac{a_{1j}}{a_{11}} a_{i1} ; i=2, \dots, n ; j=2, \dots, n+1
 \tag{45}$$

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  $a_{22}^{(1)} \neq 0$ )  $m_{i2} = a_{i2}^{(1)}/a_{22}^{(1)}$  con  $i=3, \dots, n$  resultando, al restar, el siguiente sistema:

$$\begin{aligned}
 a_{11} x_1 + a_{12} x_2 + \dots + a_{1n} x_n &= a_{1, n+1} \\
 a_{22}^{(1)} x_2 + \dots + a_{2n}^{(1)} x_n &= a_{2, n+1}^{(1)} \\
 a_{33}^{(2)} x_3 + \dots + a_{3n}^{(2)} x_n &= a_{3, n+1}^{(2)} \\
 &\dots \\
 a_{n3}^{(2)} x_3 + \dots + a_{nn}^{(2)} x_n &= a_{n, n+1}^{(2)}
 \end{aligned}
 \tag{46}$$

donde:

$$a_{ij}^{(2)} = a_{ij}^{(1)} - m_{i2} a_{2j}^{(1)} ; i=3, \dots, n \quad j=3, \dots, n+1
 \tag{47}$$

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

$$\begin{aligned}
 a_{11} x_1 + a_{12} x_2 + \dots + a_{1n} x_n &= a_{1, n+1} \\
 a_{22}^{(1)} x_2 + \dots + a_{2n}^{(1)} x_n &= a_{2, n+1}^{(1)} \\
 &\dots \\
 a_{nn}^{(n-1)} x_n &= a_{n, n+1}^{(n-1)}
 \end{aligned}
 \tag{48}$$

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

$$\begin{aligned}
 a_{ij}^k &= a_{ij}^{k-1} - m_{ik} a_{kj}^{k-1}; & k &= 1, \dots, n-1 \\
 & & j &= k + 1, \dots, n+1 \\
 & & i &= k + 1, \dots, n \\
 a_{ij}^0 &= a_{ij}
 \end{aligned}
 \tag{49}$$

con  $m_{ik} = a_{ik}^{(k-1)} / a_{kk}^{(k-1)}$

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

$$x_i = \frac{1}{a_{ii}^{(i-1)}} \left[ a_{i,n+1}^{(i-1)} - \sum_{j=i+1}^n a_{ij}^{(i-1)} x_j \right], \quad i = n, \dots, 1
 \tag{50}$$

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 Ecuación (44), en una primera pasada tiene la forma:

$$\begin{aligned}
 a_{11}^{(2)} x_1 + a_{13}^{(2)} x_3 + \dots + a_{1n}^{(2)} x_n &= a_{1,n+1}^{(2)} \\
 a_{22}^{(1)} x_2 + a_{23}^{(1)} x_3 + \dots + a_{2n}^{(1)} x_n &= a_{2,n+1}^{(1)} \\
 a_{33}^{(2)} x_3 + \dots + a_{3n}^{(2)} x_n &= a_{3,n+1}^{(2)} \\
 &\dots \\
 a_{n3}^{(2)} x_3 + \dots + a_{nn}^{(2)} x_n &= a_{n,n+1}^{(2)}
 \end{aligned}
 \tag{51}$$

Notar que:

$$a_{2,n+1}^{(2)} = a_{2,n+1}^{(1)}$$

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

$$\begin{aligned}
 a_{11} x_1 &= a_{1, n+1}^{(n-1)} \\
 a_{22}^{(1)} x_2 &= a_{2, n+1}^{(n-1)} \\
 &\dots\dots\dots \\
 a_{nn}^{(n-1)} x_n &= a_{n, n+1}^{(n-1)}
 \end{aligned}
 \tag{52}$$

de modo que la solución es simplemente:

$$x_i = \frac{a_{i, n+1}^{(n-1)}}{a_{ii}^{(i-1)}} ; \quad i=1, \dots, n
 \tag{53}$$

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 acumulación de error, además), 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 **A** y de **b**, ya sea originales o debidas al redondeo. Estudiando primero el caso más simple, que considera cambios sólo en **b**, será:

$$\underline{A} \cdot \underline{x} = \underline{b} + \underline{\delta b}
 \tag{54}$$

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:

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

$$\underline{\delta x} = \underline{A}^{-1} \cdot \underline{\delta b} \quad (55)$$

y/o, tomando normas:

$$\|\underline{\delta x}\| \leq \|\underline{A}^{-1}\| \|\underline{\delta b}\| \quad (56)$$

de igual modo:

$$\|\underline{b}\| \leq \|\underline{A}\| \cdot \|\underline{x}^*\|$$

tal que:

$$\|\underline{x}^*\| \geq \frac{\|\underline{b}\|}{\|\underline{A}\|}$$

resultando el error relativo:

$$\frac{\|\underline{\delta x}\|}{\|\underline{x}^*\|} \leq \|\underline{A}\| \|\underline{A}^{-1}\| \frac{\|\underline{\delta b}\|}{\|\underline{b}\|} \quad (57)$$

La cantidad  $\|\underline{A}\| \|\underline{A}^{-1}\|$  se denomina el *número de condición de A*,  $K(\underline{A}) \geq 1$  siempre. De modo que:

$$\frac{\|\underline{\delta x}\|}{\|\underline{x}^*\|} \leq K(\underline{A}) \frac{\|\underline{\delta b}\|}{\|\underline{b}\|} \quad (58)$$

Así, si  $K(\underline{A})$  es cercano a 1, se dice que  $\underline{A}$  es *bien condicionado*; 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:

$$(\underline{A} + \underline{\delta A}) \cdot (\underline{x}^* + \underline{\delta x}) = \underline{b} \quad (59)$$

por un procedimiento similar al anterior se llega a que:

$$\frac{\|\delta \underline{x}\|}{\|\underline{x}^* + \delta \underline{x}\|} \leq K(\underline{A}) \frac{\|\delta \underline{A}\|}{\|\underline{A}\|} \quad (60)$$

con igual connotación que en Ecuación (57). Los efectos del valor  $\|\delta \underline{A}\|$  pueden subsanarse, en parte, trabajando con mayor precisión (mayor retención de cifras significativas durante el cálculo).

### Refinamiento iterativo

Hemos visto que:

$$\underline{r} = \underline{b} - \underline{A} \cdot \underline{x}_c$$

En las situaciones que nos interesa una solución muy cercana a la verdadera  $\underline{x}^* = \underline{A}^{-1} \cdot \underline{b}$ , con un vector tal que  $\underline{r}$  sea pequeño o, en caso en que  $\underline{x}_c$  sea tal que:

$$\underline{e} = \underline{x}^* - \underline{x}_c \quad (61)$$

y sea muy pequeño en términos relativos a  $\underline{x}^*$ . Si:

$$\underline{b} = \underline{A} \cdot \underline{x}^*$$

Resulta:

$$\underline{r} = \underline{A} \cdot \underline{x}^* - \underline{A} \cdot \underline{x}_c = \underline{A} \cdot (\underline{x}^* - \underline{x}_c) = \underline{A} \cdot \underline{e} \quad (62)$$

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

$$\frac{1}{K(\underline{A})} \frac{\|\underline{r}\|}{\|\underline{b}\|} \leq \frac{\|\underline{e}\|}{\|\underline{x}^*\|} \leq K(\underline{A}) \frac{\|\underline{r}\|}{\|\underline{b}\|} \quad (63)$$

Por lo que, si el número de condición  $K(\underline{A})$  es cercano a la unidad, pequeños errores relativos de  $\underline{r}$  y  $\underline{e}$  siguen la misma tendencia (en coincidencia). 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 computada previamente, tal que el residuo para una etapa  $m$  de cálculo sea:

$$\underline{r}^{(m)} = \underline{b} - \underline{A} \cdot \underline{x}^{(m)} ; m = 1, 2, \dots, \quad (64)$$

y de acuerdo a la Ecuación (61)

$$\underline{x}^{(m+1)} = \underline{x}^{(m)} + \underline{e}^{(m)} \quad (65)$$

con  $\underline{e}^{(m)}$  calculada según la ecuación anterior.

Si  $\underline{x}^{(m+1)}$  no es satisfactorio se procede a la etapa  $(m+1)$  de cálculo. En cambio, si  $\|\underline{e}^{(m+1)}\|/\|\underline{e}^{(m)}\|$  es menor que una tolerancia establecida 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 la computación de Ecuación (64) se realiza en doble precisión, lo que implica aumento en el requerimiento de memoria.

### III.4.3 Métodos Iterativos

Su basamento es idéntico al método de aproximaciones sucesivas visto en el caso de funciones no lineales, con lo que empezando con un vector inicial se genera una sucesión de vectores, tal que:

$$\underline{x}_{i+1} = \underline{E}_i [\underline{x}_i, \underline{x}_{i-1}, \dots, \underline{x}_{i-k}] \quad (66)$$

Si la  $F_i$  (función iteradora) no es dependiente sobre  $i$  (nivel de iteración) se llama, a la recurrencia, estacionaria.

Para la mayoría de las matrices estos métodos requieren más computación, para un deseado grado de convergencia, que los métodos directos; pero para la 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. Luego para matrices ralas grandes, los métodos iterativos son los mas aconsejados por su eficiencia computacional.

Sólo nos interesarán *procesos iterativos lineales* por consideraciones de eficacia, una vez más. Un procedimiento de iteración *matricial lineal en un punto* tiene la forma general.

$$\underline{x}_{i+1} = \underline{B}_i \cdot \underline{x}_i + \underline{c}_i \quad (67)$$

dónde  $\underline{B}_i$  y  $\underline{c}_i$  son independientes de  $i$  (iteración estacionaria).

Recordar que en los métodos iterativos basados en substitución directa (o Wegstein) debe explicitarse la variable independiente. A modo de introducción se puede ver dicha forma haciendo:



$$\begin{aligned} \underline{A} \cdot \underline{x} &= \underline{b} \\ \underline{A} \cdot \underline{x} + \underline{x} &= \underline{b} + \underline{x} \\ (\underline{A} + \underline{I}) \cdot \underline{x} &= \underline{x} + \underline{b} \\ \underline{x} &= (\underline{A} + \underline{I}) \cdot \underline{x} - \underline{b} \end{aligned}$$

o, en la ley de recurrencia:

$$\underline{x}_{i+1} = (\underline{I} + \underline{A}) \cdot \underline{x}_i - \underline{b} \quad (68)$$

que nos da una forma genérica simple de los métodos iterativos estacionarios que son los universalmente usados.

### Método de Jacobi

Se escribe la matriz  $A$  como:

$$\underline{A} = \underline{D} + \underline{L} + \underline{U} \quad (69)$$

siendo  $D$  la matriz diagonal (formada precisamente con esos elementos).  $L$  y  $U$  son triangulares inferior y superior, respectivamente, con el resto de los elementos; y ceros en su diagonal principal. Luego, si:

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

obtenemos:

$$\underline{x}_{i+1} = -\underline{D}^{-1} \cdot (\underline{L} + \underline{U}) \cdot \underline{x}_i + \underline{D}^{-1} \cdot \underline{b} \quad (70)$$

por supuesto que supone que la diagonal principal no tiene todos ceros (en  $A$ ).

De no ser así, si  $A$  es no singular, se debe permutar filas y columnas para obtener una forma que permita definir  $D$ . Es deseable que sea de la forma diagonal dominante, esto que los elementos de la diagonal sean lo más grandes posibles respecto a los demás. Para este método de Jacobi la matriz  $B$  puede entonces escribirse como:

$$\underline{B} = - \underline{D}^{-1} \cdot (\underline{L} + \underline{U}) \quad (71)$$

haciendo:

$$\underline{x}_{i+1} = \underline{B} \cdot \underline{x}_i + \underline{D}^{-1} \cdot \underline{b} \quad (72)$$

### 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 es calculada cada componente de  $\underline{x}_{i+1}$  es usada inmediatamente en la misma iteración. Al proceder de ese modo, la ecuación obtenida para Jacobi, resulta:

$$\underline{x}_{i+1} = - \underline{D}^{-1} \cdot (\underline{L} \cdot \underline{x}_{i+1}) + \underline{D}^{-1} \cdot \underline{U} \cdot \underline{x}_i + \underline{D}^{-1} \cdot \underline{b} \quad (73)$$

reordenando y sacando factor común  $\underline{x}_{i+1}$  se tiene:

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

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

$$(\underline{D} + \underline{L}) \cdot \underline{x}_{i+1} = - \underline{U} \cdot \underline{x}_i + \underline{b}$$

luego:

$$\underline{x}_{i+1} = - (\underline{D} + \underline{L})^{-1} \cdot \underline{U} \cdot \underline{x}_i + (\underline{D} + \underline{L})^{-1} \cdot \underline{b} \quad (74)$$

entonces será, tomando  $\underline{B}$  para la expresión generalizada:

$$\underline{x}_{i+1} = + \underline{B} \cdot \underline{x}_i + (\underline{D} + \underline{L})^{-1} \cdot \underline{b} \quad (75)$$

$$\underline{B} = - (\underline{D} + \underline{L})^{-1} \cdot \underline{U} \quad (76)$$

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.

### PROBLEMAS PROPUESTOS

P1) Sea la siguiente ecuación definida para  $x > 0$ :

$$x^2 - \ln(3x) + 1/9 = 0$$

Encontrar las raíces de la misma por los siguientes métodos:

- a) Wegstein
- b) Newton-Raphson

P2) Para el flujo isoentrópico de un gas perfecto que fluye desde un reservorio a través de una boquilla convergente - divergente operando con velocidad sónica en la constricción, se puede demostrar que:

$$\frac{A_c^2}{A^2} = \left( \frac{\gamma + 1}{2} \right)^{(\gamma + 1) / (\gamma - 1)} \left( \frac{\gamma - 1}{2} \right) \left[ \left( \frac{P}{P_r} \right)^{2 / \gamma} - \left( \frac{P}{P_r} \right)^{(\gamma + 1) / \gamma} \right]$$

donde  $P$  es la presión sobre el área transversal  $A$  de la boquilla,  $P_r$  es la presión en el reservorio,  $A_c$  es la sección transversal en la constricción y  $\gamma$  es la relación entre el calor específico a presión constante y el calor específico a volumen constante.

Si  $A_c$ ,  $\gamma$ ,  $P_r$  y  $A$  ( $> A_c$ ) se conocen, idear un esquema para calcular las *dos posibles presiones*  $P$  que satisfacen la ecuación de arriba. Implemente su método en la computadora.

*Datos sugeridos para prueba:*

$$A_c = 0.1 \text{ pie}^2, \gamma = 1.41, P_r = 100 \text{ psia y } A = 0.12 \text{ pie}^2.$$

P3) Una bolsa esférica de gas a alta presión, inicialmente de radio  $r_0$  y presión  $p_0$ , se expande radialmente en una explosión submarina adiabática. Para el caso especial de un gas con  $\gamma = 4/3$  (relación de calor específico a presión constante a calor específico a volumen constante), el radio  $r$  en cualquier instante de tiempo  $t$  posterior a la explosión se expresa como:

$$\frac{t}{r_0} \sqrt{\frac{p_0}{\rho}} = \left( 1 + \frac{2}{3} \alpha + \frac{1}{5} \alpha^2 \right) (2 \alpha)^{1/2}$$

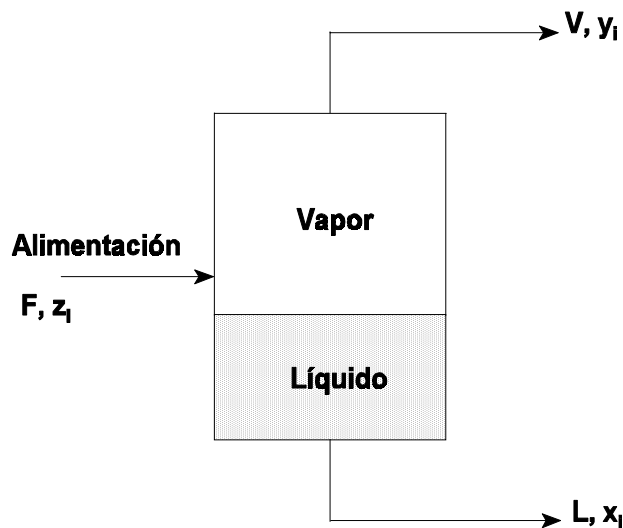
donde  $\alpha = (r/r_0) - 1$  y  $\rho$  es la densidad del agua. Durante la expansión adiabática, la presión del gas se expresa como  $(p/p_0) = (r_0/r)^{\gamma}$ .

Desarrollar un procedimiento para calcular la presión y el radio de la bolsa para cualquier instante de tiempo.

*Datos sugeridos para prueba:*

$p_0 = 10^4 \text{ lb}_f/\text{pulgada}^2$ ,  $\rho = 64 \text{ lb}_m/\text{pie}^3$ ,  $r_0 = 1 \text{ pie}$ ,  $t = 0.5, 1, 2, 3, 5$  y  $10$  milisegundos.

P4)  $F$  moles/hr. de una corriente de gas natural licuado de  $n$ -componentes se introduce como corriente de alimentación a un tanque de vaporización flash como se indica en la figura siguiente:



Las corrientes resultantes de vapor y de líquido se extraen a las velocidades de  $V$  y  $L$  moles/hr., respectivamente. Las fracciones molares de los componentes en la alimentación, en el vapor y en el líquido, se designan,  $z_i$ ,  $y_i$ , e  $x_i$ , respectivamente ( $i = 1, 2, \dots, n$ ). Suponiendo equilibrio líquido/vapor y operación en estado estacionario (en el capítulo IX desarrollaremos con mayor extensión el análisis de este equipo),

tenemos:

Balance de materia:  $F = L + V$

Balance individual:  $z_i F = x_i L + y_i V$

Relación de equilibrio:  $K_i = y_i/x_i$

Aquí,  $K_i$ , es la constante de equilibrio para la  $i$ -ésima componente a la temperatura y presión prevalecientes en el tanque. A partir de estas ecuaciones y del hecho que:

$$\sum_{i=1}^n x_i = \sum_{i=1}^n y_i = 1$$

probar que:

$$\sum_{i=1}^n \frac{z_i (K_i - 1)}{\Theta (K_i - 1) + 1} = 0$$

donde  $\Theta = V/F$  representa la fracción vaporizada.

Escriba un programa que lea los valores de  $F$ ,  $z_i$  y los  $K_i$  como datos y entonces utilice el método de Newton para resolver esta última ecuación en  $\Theta$ . El programa debería también calcular la fracción de líquido  $L/F$ , los  $x_i$  e  $y_i$  utilizando las tres primeras ecuaciones dadas más arriba. Los datos de prueba, que se muestran en la tabla siguiente, están relacionados con el flasheo de una corriente de gas natural a 1600 psia y 120 °F.

Componente	i	$z_i$	$K_i$
<b>Dióxido de carbono</b>	1	0.0046	1.650
<b>Metano</b>	2	0.8345	3.090
<b>Etano</b>	3	0.0381	0.720
<b>Propano</b>	4	0.0163	0.390
<b>Isobutano</b>	5	0.0050	0.210
<b>n-Butano</b>	6	0.0074	0.175
<b>Pentanos</b>	7	0.0287	0.093
<b>Hexanos</b>	8	0.0220	0.065
<b>Heptanos</b>	9	0.0434	0.036
<b>Total</b>		1.0000	

Suponer que  $F = 1000$  moles/hr. Además deberían ser leídos como datos la tolerancia

---

**Modelado, Simulación y Optimización de Procesos Químicos**

Autor: Nicolás J. Scenna y col.

ISBN: 950-42-0022-2 - ©1999

$\varepsilon$  y un número máximo de iteraciones. ¿Cuál sería un buen valor de arranque para  $\theta_i$ ?

P5) Para el flujo turbulento de un fluido a través de un tubo liso, es posible establecer la siguiente relación entre el factor de fricción  $c_f$  y el número de Reynolds  $Re$ :

$$\sqrt{\frac{1}{c_f}} = -0.4 + 1.74 \ln \left( Re \sqrt{c_f} \right)$$

Calcular  $c_f$  para  $Re = 10^4, 10^5$  y  $10^6$ .

P6) Para el flujo estacionario de un fluido incompresible a través de un tubo rugoso de longitud  $L$  y diámetro interior  $D$ , la caída de presión viene expresada por la siguiente relación:

$$\Delta p = \frac{f_M \rho u_M^2 L}{2 D}$$

donde  $\rho$  es la densidad del fluido,  $u_m$  es la velocidad media del fluido y  $f_M$  es el factor de fricción de Moody (adimensional). El factor de fricción de Moody es una función de la rugosidad  $\varepsilon$  y del número de Reynolds,

$$Re = \frac{D \rho u_M}{\mu}$$

donde  $\mu$  es la viscosidad del fluido. Para  $Re \leq 2000$ ,

$$f_M = \frac{64}{Re}$$

mientras que para  $Re > 2000$ ,  $f_M$  viene expresada por la Ecuación de Colebrook,

$$\sqrt{\frac{1}{f_M}} = -2 * \log_{10} \left( \frac{\varepsilon}{3.7 D} + \frac{2.51}{Re \sqrt{f_M}} \right)$$

Un buen punto de arranque para la solución iterativa de esta ecuación puede encontrarse a partir de la Ecuación de Blasius,

$$f_M = 0.316 Re^{-0.25}$$

apropiada para flujo turbulento en tubos lisos.

Escribir una función (function), *PDELTA*, que podría ser llamada a través de

la sentencia:  $DELTA = PDELTA(Q, D, L, RHO, MU, E)$  donde el valor de  $PDELTA$  representa la caída de presión para el flujo volumétrico (caudal)  $Q$  de un fluido con densidad  $RHO$  y viscosidad  $MU$  a través de un tubo de longitud  $L$ , diámetro interior  $D$  y rugosidad  $E$ . Observe que  $MU$  y  $L$  deben definirse como variables reales.

Escriba un programa que lea los valores de  $Q, D, L, RHO, MU$  y  $E$ , llame a  $PDELTA$  para calcular la caída de presión, imprima los datos y resultados y vuelva a leer otro conjunto de datos.

*Datos sugeridos para prueba:*

Q	gal/min	170	4
D	pulgada	3.068	0.622
L	pie	10000	100
RHO	lb <sub>m</sub> /pie <sup>3</sup>	62.4	80.2
MU	lb <sub>m</sub> /pie seg.	0.0007	0.05
E	pulgada	0.002	0.0005

P7) Sea el siguiente sistema de ecuaciones lineales:

$$\begin{aligned} 9x_1 + 3x_2 + 5x_3 &= -3 \\ 2x_2 - x_3 &= 0 \\ 4x_1 - 3x_2 + 7x_3 &= -1 \end{aligned}$$

Encontrar la solución por los siguientes métodos:

- Gauss
- Jacobi
- Gauss-Seidel
- Sustitución directa
- Wegstein

P8) Sea el siguiente sistema:

$$\begin{aligned} 9x_1 + 2x_2 &= 3 \\ 2x_1 + 3x_2 + x_3 &= 4 \\ 5x_2 + 4x_3 + 8x_4 &= -8 \\ 3x_3 - x_4 &= 0 \end{aligned}$$

Encontrar la solución mediante el método de Gauss. ¿Puede aprovechar la estructura particular del sistema para facilitar el cálculo?.

P9) Sea el siguiente sistema:

$$\begin{aligned} 2x + 3y - z &= 0 \\ 4x - 3y - 9z &= 2 \\ 3x - 3y - z &= 1/3 \end{aligned}$$

Aplicar el método de Gauss-Seidel para encontrar la solución.

P10) Sea el siguiente sistema expresado en forma matricial:

$$\begin{bmatrix} 1 & 3 & 5 \\ 2 & 4 & 3 \\ 1 & 1 & 2 \end{bmatrix} \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \Delta x_3 \end{bmatrix} = \begin{bmatrix} 2 \\ 0 \\ 3 \end{bmatrix}$$

Encontrar los valores para  $\Delta x_1$ ,  $\Delta x_2$  y  $\Delta x_3$  por los siguientes métodos:

- Gauss
- Gauss-Seidel
- Jacobi

P11) Confeccionar un diagrama de flujos para un programa que ingrese como datos la matriz de coeficientes y el término independiente de un sistema de ecuaciones lineales de dimensión  $N$ , y calcule la solución. Suponer que dispone de subrutinas de inversión de matrices.

P12) Utilizando software comercial (rutinas IMSL, Numerical Recipes, etc.) confeccionar el programa correspondiente al diagrama de flujos anterior.

### BIBLIOGRAFÍA RECOMENDADA

- ▶ Carnahan, B., H. A. Luther y J. O. Wilkes, *Applied Numerical Methods*, John Wiley and Sons, Inc., New York (1969).
- ▶ Cohen, A. M., J. F. Cutts, R. Fielder, D. E. Jones, J. Ribbans y E. Stuart, *Análisis Numérico*, Editorial Reverté S. A., Barcelona (1977).
- ▶ Hamming, R. W., *Numerical Methods for Scientists and Engineers*, McGraw-Hill, New York (1962).
- ▶ Hildebrand, F. B., *Introduction to Numerical Analysis*, McGraw-Hill, New York (1956).
- ▶ Hornbeck, R. W., *Numerical Methods*, Quantum Publishers, Inc., New York.
- ▶ Lapidus, L., *Digital Computation for Chemical Engineers*, McGraw-Hill, New York (1962).

---

### Modelado, Simulación y Optimización de Procesos Químicos

Autor: Nicolás J. Scenna y col.

ISBN: 950-42-0022-2 - ©1999



- ▶ Luthe, R., A. Olivera y F. Schutz, *Métodos Numéricos*, Editorial Limusa, México (1978).
- ▶ Milne, W. E., *Numerical Solution of Differential Equations*, John Wiley and Sons, New York (1953).
- ▶ O'Neill, P. V., *Advanced Engineering Mathematics*, 3ra. Ed., Wadsworth, Inc., USA (1991).
- ▶ Perry, J. H., *Chemical Engineers Handbook*, McGraw-Hill, New York.
- ▶ Ralston, A. y P. Rabinowitz, *First Course in Numerical Analysis*, McGraw-Hill Inc., USA (1978).