



---

**Universidad de Valladolid**

Facultad de Ciencias

## **TRABAJO FIN DE GRADO**

Grado en Matemáticas

**Una introducción a los métodos multigrad.**

*Autor: Marcos Santana Pastor*

*Tutora: Begoña Cano Urdiales*



# Índice general

<b>Introducción</b>	<b>5</b>
<b>1. Problema modelo y métodos iterativos clásicos</b>	<b>7</b>
1.1. Problema modelo: Problema Poisson en 2D . . . . .	7
1.2. Convergencia de los métodos iterativos clásicos . . . . .	9
1.2.1. Convergencia del método de Jacobi . . . . .	9
1.2.2. Convergencia del método de Gauss-Seidel . . . . .	14
1.2.3. Convergencia del método S.O.R. . . . .	15
<b>2. Método multigrad</b>	<b>21</b>
2.1. Ideas principales . . . . .	21
2.2. Procedimiento de suavizado del error . . . . .	24
2.2.1. Iteración tipo Jacobi . . . . .	24
2.2.2. Iteración tipo Gauss-Seidel . . . . .	27
2.3. Análisis local de Fourier de las propiedades de suavizado . . . . .	28
2.3.1. Primeros conceptos para el análisis del suavizado . . . . .	28
2.3.2. Análisis del suavizado . . . . .	30
2.4. Introducción al ciclo de 2 mallas . . . . .	36
2.4.1. Aproximación de la solución de la ecuación del defecto . . . . .	36
2.4.2. Corrección en malla grosera . . . . .	38
2.4.3. Estructura del operador de dos mallas . . . . .	40
2.5. Componentes del método multigrad . . . . .	40
2.5.1. Elección del operador en la red grosera $L_H$ . . . . .	41
2.5.2. Elección del operador restricción . . . . .	41
2.5.3. Elección del operador interpolador . . . . .	42
2.6. Análisis local de Fourier del ciclo de 2 mallas . . . . .	42
2.6.1. Factor de convergencia asintótico . . . . .	46
2.7. El ciclo multigrad . . . . .	47
2.7.1. Definición recursiva del método multigrad . . . . .	47
2.7.2. Coste computacional . . . . .	50
2.7.3. Convergencia y eficiencia del método multigrad . . . . .	52
2.8. $h$ - Independencia de la convergencia . . . . .	55
<b>3. Conclusiones Finales</b>	<b>59</b>

**A. Programas en MATLAB**

**61**

# Introducción

En el análisis numérico, el método multigrad es un algoritmo muy eficiente y ampliamente utilizado para aproximar la solución de problemas de ecuaciones diferenciales en el espacio utilizando distintas mallas de puntos con distintos diámetros donde el problema ha sido discretizado previamente en la malla más fina. La idea principal de los métodos multigrad es aumentar la velocidad de convergencia de los métodos iterativos clásicos (Jacobi o Gauss-Seidel). Los métodos clásicos tienen la propiedad de reducir las componentes de alta frecuencia de los errores de la aproximación. La mejora que aportan los métodos multigrad viene de la corrección de la aproximación en una malla más gruesa, es decir, resolver un problema similar pero en una malla con diámetro mayor. El problema en la malla más gruesa, aunque es más barato de resolver, es similar al problema en la malla más fina en el sentido de que también tendrá componentes del error de alta y baja frecuencia. Para resolver esto, el ciclo multigrad aplica recursivamente este razonamiento sobre el problema hasta alcanzar la red más gruesa, donde el coste computacional de resolver es despreciable frente al de aplicar el método clásico en la malla más fina.

Los primeros estudios realizados sobre el método multigrad, en un sentido estricto, son de Fedorenko [6], [7] entre los años 1962 y 1964 y después vinieron los estudios de Bakhvalov [1] en 1966. Mientras Fedorenko se restringió en el estudio de la convergencia para problemas discretos de segundo orden con valores frontera y con coeficientes variables en el cuadrado unidad, Bakhvalov considero la posibilidad de combinar los métodos multigrad con métodos que ayudaran a elegir el mejor iterante inicial. La actual eficiencia del método fue reconocida y mostrada por Brandt entre 1970 y 1977 en [3], [4] y [5]. Algunas de las principales contribuciones de Brandt fueron la introducción del método multigrad no lineal, discusiones del método en dominios generales o el uso del análisis local de Fourier como herramienta para el estudio del diseño del ciclo multigrad. Desde 1980, muchos investigadores han contribuido al desarrollo del campo de los métodos multigrad. El ciclo multigrad ha sido utilizado para la ecuación 3D Helmholtz para la predicción meteorológica o también para las ecuaciones de elasticidad de Lamé o las ecuaciones de Navier-Stokes.

En este trabajo se pretende introducir los métodos multigrad para aproximar la solución del problema de Poisson con condiciones frontera en un cuadrado. Partiendo del problema modelo continuo, discretizaremos mediante la fórmula de los cinco puntos y tendremos un sistema lineal que resolver. Veremos cómo

los métodos multigrad son los más adecuados para aproximar la solución del problema discreto cuando son comparados con los métodos clásicos de Jacobi, Gauss-Seidel o S.O.R. Muchos de los razonamientos y resultados que se encuentran en este texto son generalizables a otras discretizaciones de problemas lineales elípticos. Algunos libros que tratan con mayor profundidad los métodos multigrad son [2], [8] y [9].

La bibliografía más utilizada para la realización del trabajo es [9], [10] y [11]. Más concretamente, [10] y [11] para el estudio de los métodos clásicos en el capítulo 1 y [9] para el estudio del ciclo multigrad en el capítulo 2. Además, es constante el uso de los contenidos de las asignaturas de Ampliación de Análisis Numérico y de Álgebra Lineal del Grado.

La estructura del trabajo es la siguiente:

En el capítulo 1 introduciremos el problema modelo, que será el ejemplo sobre el que aplicaremos y analizaremos el método multigrad a lo largo de todo el trabajo. También haremos un estudio exhaustivo de la convergencia de los métodos iterativos clásicos sobre dicho problema modelo, centrándonos en su dependencia del diámetro de la malla.

En el capítulo 2 desarrollaremos detenidamente las ideas principales en las que se basa el método multigrad. Analizaremos las propiedades de suavizado de los métodos clásicos, daremos la estructura del ciclo en dos mallas y veremos por separado las distintas componentes del método multigrad. Tras este estudio en el ciclo de dos mallas, generalizaremos los resultados obtenidos a la multimalla con la intención de obtener así el método multigrad. Estudiaremos el coste computacional que requiere el método y compararemos su eficiencia con los métodos clásicos iterativos. Por último, veremos la independencia que existe entre la convergencia del método multigrad y el diámetro de la malla para un tipo concreto de métodos descritos (W-ciclos).

En el capítulo 3, resumiremos las conclusiones finales obtenidas durante en el trabajo.

Finalmente, en el apéndice mostramos la implementación del método multigrad en MATLAB, así como los métodos iterativos clásicos aplicados a nuestro problema modelo.

# Capítulo 1

## Problema modelo y métodos iterativos clásicos

En este capítulo introducimos el problema de Poisson en 2 dimensiones con condiciones frontera tipo Dirichlet y su discretización por la fórmula de los cinco puntos. El dominio en el que nos centraremos será el cuadrado unidad  $[0, 1] \times [0, 1]$ , donde tendremos un malla de puntos a la cual dotaremos de un cierto orden.

Analizaremos el comportamiento de los métodos iterativos clásicos para resolver el sistema lineal obtenido al discretizar el problema. Estos métodos iterativos serán el método de Jacobi, Gauss-Seidel y el método S.O.R. para el orden de la red de puntos lexicográfico.

La conclusión a la que llegaremos es que la convergencia de los métodos clásicos puede ser rápida cuando la distancia entre los puntos de la malla es relativamente grande. Sin embargo, si hacemos dicha distancia cada vez más pequeña tendremos que el radio espectral de las matrices de iteración de los métodos será cada vez más cercano a 1 y, por lo tanto, la convergencia será cada vez más lenta.

### 1.1. Problema modelo: Problema Poisson en 2D

Uno de los problemas elípticos más sencillos es el problema de Poisson

$$\Delta u(x, y) = u_{xx}(x, y) + u_{yy}(x, y) = f^\Omega(x, y), \quad (x, y) \in \Omega,$$

junto a la condición frontera tipo Dirichlet

$$u(x, y) = f^\Gamma(x, y), \quad (x, y) \in \Gamma$$

donde  $\Omega := (0, 1) \times (0, 1)$  y  $\Gamma$  su frontera.

La discretización por la fórmula de los 5 puntos la representaremos de la forma:

8CAPÍTULO 1. PROBLEMA MODELO Y MÉTODOS ITERATIVOS CLÁSICOS

$$\begin{cases} \Delta_h u_h(x, y) = f_h^\Omega(x, y), & (x, y) \in \Omega_h, \\ u_h(x, y) = f_h^\Gamma(x, y), & (x, y) \in \Gamma_h, \end{cases} \quad (1.1)$$

donde si  $G_h$  es la red infinita de puntos en  $\mathbb{R}^2$  definida de la forma

$$G_h := \{(x, y) : x = x_i = ih, \quad y = y_j = jh : (i, j) \in \mathbb{Z}^2\},$$

entonces  $\Omega_h := \Omega \cap G_h$ ,  $\Gamma_h := \Gamma \cap G_h$ , con  $h = 1/N$  para algún  $N \in \mathbb{N}$ . Además,  $u_h$  es una función tomando valores en la red  $\Omega_h$  y la discretización del Laplaciano  $\Delta_h$  viene dado por

$$\Delta_h u(x, y) = \frac{u(x-h, y) + u(x+h, y) + u(x, y-h) + u(x, y+h) - 4u(x, y)}{h^2},$$

que es un método en diferencias finitas consistente, ya que

$$\|\Delta u - \Delta_h u\|_{\infty, \Omega_h} = O(h^2), \quad h \rightarrow 0,$$

siempre que  $u$  sea suficientemente derivable para hacer su desarrollo de Taylor (al menos son necesarias las derivadas hasta orden 4 de  $u$  en  $\bar{\Omega}$ ).

Los elementos  $f_h^\Omega, f_h^\Gamma$  de la ecuación, son las funciones red que representan a  $f^\Omega, f^\Gamma$  en la discretización, es decir, son las funciones  $f^\Omega, f^\Gamma$  evaluadas en la malla de puntos  $\Omega_h$  y en su frontera  $\Gamma_h$ , respectivamente.

**Observación 1.1.1** *Es importante tener en cuenta el orden que tienen los puntos dentro de la red a la hora de resolver el sistema asociado a (1.1), especialmente para los métodos de S.O.R y Gauss-Seidel. Nosotros, a lo largo de este trabajo, consideraremos el orden lexicográfico dentro de la red de puntos. Una breve explicación de lo que esto significa en nuestro problema modelo es que el primer elemento es el que se encuentra más abajo y más a la izquierda en  $\Omega_h$ , el siguiente será justo el que esté a su derecha y así sucesivamente hasta completar la fila y comenzar con la siguiente.*

Para finalizar con los detalles sobre el problema modelo, daremos la matriz del operador discreto con el orden lexicográfico. A esta matriz la denotamos por  $A_h$  y es de la forma

$$A_h := \begin{pmatrix} T_h & -I_h & & & \\ -I_h & T_h & -I_h & & \\ & -I_h & T_h & -I_h & \\ & & \ddots & \ddots & \ddots \\ & & & -I_h & T_h \end{pmatrix},$$

donde la matriz  $I_h$  es la matriz identidad de orden  $N-1$  y  $T_h$  es la matriz

$$T_h := \begin{pmatrix} 4 & -1 & & & \\ -1 & 4 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & & -1 & 4 \end{pmatrix}_{(N-1) \times (N-1)}.$$

Nótese que la dimensión de la matriz  $A_h$  es justamente el número de puntos de la malla. Cuanto más fina sea la malla, más grande será la matriz.

## 1.2. Convergencia de los métodos iterativos clásicos

A continuación, vamos a desarrollar los cálculos necesarios para obtener los radios espectrales de las matrices de iteración de los métodos iterativos de Jacobi, Gauss-Seidel y el método S.O.R. a la hora de aproximar una solución para nuestro problema modelo. En particular, hallaremos el  $\omega$  óptimo para que el método S.O.R. converja lo más rápido posible a la solución del problema.

Lo común para los tres métodos iterativos es que estamos aproximando la solución del sistema lineal:

$$\begin{aligned} h^2 f^\Omega(x, y) &= u_h(x+h, y) + u_h(x, y+h) + u_h(x-h, y) + u_h(x, y-h) \\ &\quad - 4u_h(x, y), \quad (x, y) \in \Omega_h \\ f_h^\Gamma(x, y) &= u_h(x, y), \quad (x, y) \in \Gamma_h \end{aligned} \quad (1.2)$$

Probaremos que el método de Gauss-Seidel converge más rápido que el método de Jacobi, pues el radio espectral de la matriz de iteración de Gauss-Seidel es menor que el de la matriz de iteración de Jacobi. De la misma forma, obtendremos que el método S.O.R. es el mejor de entre los tres métodos iterativos clásicos.

### 1.2.1. Convergencia del método de Jacobi

Partiendo del iterante  $u_h^n$  en toda la región  $\Omega_h \cup \Gamma_h$  donde  $u_h^n(x, y) = f_h^\Gamma(x, y)$  para todo  $(x, y) \in \Gamma_h$ , el siguiente iterante, con el método de Jacobi, viene determinado por

$$\begin{aligned} u_h^{n+1}(x, y) &= \frac{u_h^n(x+h, y) + u_h^n(x, y+h) + u_h^n(x-h, y) + u_h^n(x, y-h)}{4} \\ &\quad - \frac{h^2}{4} f_h^\Omega(x, y), \quad (x, y) \in \Omega_h \\ u_h^{n+1}(x, y) &= f_h^\Gamma(x, y) \quad (x, y) \in \Gamma_h \end{aligned}$$

Si  $u_h$  es la solución exacta del sistema lineal (1.2), entonces tenemos la ecuación del error

$$e_h^n(x, y) = u_h^n(x, y) - u_h(x, y), \quad (1.3)$$

que obviamente satisface

$$e_h^{n+1}(x, y) = \frac{e_h^n(x+h, y) + e_h^n(x, y+h) + e_h^n(x-h, y) + e_h^n(x, y-h)}{4}$$

para los puntos  $(x, y) \in \Omega_h$  y,  $e_h^{n+1}(x, y) = 0$  para  $(x, y) \in \Gamma_h$ .

Dados los vectores  $U_h^n$  y  $U_h$  conteniendo los valores nodales de  $u_h^n$  y  $u_h$  en  $\Omega_h$  en el orden lexicográfico, definimos el vector  $E_h^n := U_h^n - U_h$  del espacio vectorial  $\mathbb{R}^K$ , con  $K = (N-1)^2$ , como el vector formado por la evaluación de los errores en los puntos de la malla. Y, por lo tanto,

$$E_h^{n+1} = B_h E_h^n$$

siendo  $B_h$  la matriz de iteración del método de Jacobi.

Ahora nos centramos en buscar los autovalores de la matriz  $B_h$ . Para ello, buscamos funciones  $v_h$  en  $\Omega_h \cup \Gamma_h$  no nulas verificando para algún  $\mu_h$  que

$$\begin{aligned} \mu_h v_h(x, y) &= \frac{v_h(x+h, y) + v_h(x, y+h) + v_h(x-h, y)}{4} \\ &\quad + \frac{v_h(x, y-h)}{4}, \quad (x, y) \in \Omega_h, \\ v_h(x, y) &= 0, \quad (x, y) \in \Gamma_h. \end{aligned} \quad (1.4)$$

Consideremos las funciones

$$v_{p,q}(x, y) = \sin(p\pi x) \sin(q\pi y), \quad (1.5)$$

con  $1 \leq p \leq (N-1)$ ,  $1 \leq q \leq (N-1)$  y veamos que cumplen lo anterior para cierto valor  $\mu$ . Más concretamente, utilizando la fórmula del seno de la suma de ángulos y sacando factores comunes tenemos

$$\begin{aligned} &\frac{1}{4} [\sin(p\pi(x+h)) \sin(q\pi y) + \sin(p\pi x) \sin(q\pi(y+h)) + \sin(p\pi(x-h)) \sin(q\pi y) \\ &\quad + \sin(p\pi x) \sin(q\pi(y-h))] \\ &= \frac{1}{4} [\sin(p\pi x) \sin(q\pi y) (2 \cos(p\pi h) + 2 \cos(q\pi h))] \\ &= \frac{1}{2} [(\cos(p\pi h) + \cos(q\pi h)) v_{p,q}(x, y)] \end{aligned}$$

Por lo tanto, cada función  $v_{p,q}(x, y)$  tiene asociado el autovalor

$$\mu = \mu_{h,p,q} = \frac{\cos(p\pi h) + \cos(q\pi h)}{2}. \quad (1.6)$$

El correspondiente autovector  $V_{h,p,q} \in \mathbb{R}^K$  estará formado por las siguientes componentes.

$$V_{h,p,q} := \begin{pmatrix} v_{p,q}(h, h) \\ v_{p,q}(2h, h) \\ \vdots \\ v_{p,q}((N-1)h, h) \\ \vdots \\ v_{p,q}((N-1)h, (N-1)h) \end{pmatrix}.$$

Estos vectores son linealmente independientes. Esto lo probaremos demostrando que son ortogonales. Observemos que

$$\begin{aligned} \langle V_{h,p,q}, V_{h,p',q'} \rangle &= \sum_{(x,y) \in \Omega_h} v_{p,q}(x,y) v_{p',q'}(x,y) \\ &= \sum_{(x,y) \in \Omega_h} \sin(p\pi x) \sin(q\pi y) \sin(p'\pi x) \sin(q'\pi y) \\ &= \left( \sum_{i=1}^{N-1} \sin\left(\frac{p\pi i}{N}\right) \sin\left(\frac{p'\pi i}{N}\right) \right) \left( \sum_{j=1}^{N-1} \sin\left(\frac{q\pi j}{N}\right) \sin\left(\frac{q'\pi j}{N}\right) \right). \end{aligned}$$

Notemos ahora que, si  $(p, q) \neq (p', q')$ , aplicando el siguiente lema, al menos uno de los factores anteriores es nulo, con lo que queda probada la ortogonalidad de los vectores.

**Lema 1.2.1** *Si  $p \neq p'$ , entonces  $\sum_{k=1}^{N-1} \sin\left(\frac{p\pi k}{N}\right) \sin\left(\frac{p'\pi k}{N}\right) = 0$ .*

**Demostración.**

Debido a la siguiente igualdad trigonométrica

$$\sin(\alpha) \sin(\beta) = \frac{1}{2} (\cos(\alpha - \beta) - \cos(\alpha + \beta)),$$

tenemos que el interior del sumatorio lo podemos expresar de la forma que sigue

$$\sin\left(\frac{p\pi j}{N}\right) \sin\left(\frac{p'\pi j}{N}\right) = \frac{1}{2} \left( \cos\left(\frac{(p-p')\pi j}{N}\right) - \cos\left(\frac{(p+p')\pi j}{N}\right) \right).$$

Distinguimos dos casos:

- Si  $p - p'$  es impar. Entonces  $p + p'$  será también impar, puesto que  $p + p' =$

$(p - p') + 2p'$ . En cualquier caso, si  $k$  es un número impar tenemos:

$$\begin{aligned} \sum_{j=1}^{N-1} \cos\left(\frac{k\pi j}{N}\right) &= \sum_{j=1}^{\lfloor \frac{N-1}{2} \rfloor} \cos\left(\frac{k\pi j}{N}\right) + \sum_{j=\lfloor \frac{N-1}{2} \rfloor+1}^{N-1} \cos\left(\frac{k\pi j}{N}\right) \\ &= \sum_{j=1}^{\lfloor \frac{N-1}{2} \rfloor} \cos\left(\frac{k\pi j}{N}\right) + \sum_{l=1}^{N-\lfloor \frac{N-1}{2} \rfloor-1} \cos\left(k\pi - \frac{k\pi l}{N}\right) \\ &= \sum_{j=1}^{\lfloor \frac{N-1}{2} \rfloor} \cos\left(\frac{k\pi j}{N}\right) - \sum_{l=1}^{N-\lfloor \frac{N-1}{2} \rfloor-1} \cos\left(\frac{k\pi l}{N}\right), \end{aligned}$$

donde, para la segunda igualdad, hemos utilizado en el segundo sumatorio el cambio de índice  $j = N - l$  y, para la tercera igualdad, hemos tenido en cuenta que  $k$  es impar.

Si  $N - 1$  es par, entonces

$$\lfloor \frac{N-1}{2} \rfloor = \frac{N-1}{2},$$

y también

$$N - \lfloor \frac{N-1}{2} \rfloor - 1 = \frac{N-1}{2}.$$

Por lo tanto, los sumatorios anteriores quedan:

$$\sum_{j=1}^{\frac{N-1}{2}} \cos\left(\frac{k\pi j}{N}\right) - \sum_{l=1}^{\frac{N-1}{2}} \cos\left(\frac{k\pi l}{N}\right) = 0.$$

Si ahora  $N - 1$  es impar, entonces

$$\lfloor \frac{N-1}{2} \rfloor = \frac{N}{2} - 1$$

y también

$$N - \lfloor \frac{N-1}{2} \rfloor - 1 = \frac{N}{2}.$$

Por lo tanto, los sumatorios quedan de la forma:

$$\sum_{j=1}^{\frac{N}{2}-1} \cos\left(\frac{k\pi j}{N}\right) - \sum_{l=1}^{\frac{N}{2}} \cos\left(\frac{k\pi l}{N}\right) = \sum_{j=1}^{\frac{N}{2}-1} \cos\left(\frac{k\pi j}{N}\right) - \sum_{l=1}^{\frac{N}{2}-1} \cos\left(\frac{k\pi l}{N}\right) = 0$$

puesto que el último sumando del segundo sumatorio es 0.

- Si  $p - p'$  es par. Con un razonamiento similar al anterior se prueba que  $p + p'$  es par.

Antes que nada, observemos que  $e^{\frac{2\pi ik}{N}}$ , para  $k \neq mN$  con  $m \in \mathbb{Z}$ , es una raíz  $N$ -ésima de la unidad distinta de 1. Por lo tanto, verifica la ecuación

$$x^N - 1 = (x - 1)(x^{N-1} + x^{N-2} + \dots + x + 1) = 0.$$

Entonces, como  $e^{\frac{2\pi ik}{N}}$  es distinto de 1, tenemos

$$\sum_{j=0}^{N-1} e^{\frac{2\pi ijk}{N}} = 0.$$

Si igualamos la parte real a 0, llegamos a la siguiente expresión

$$\sum_{j=1}^{N-1} \cos\left(\frac{2\pi jk}{N}\right) = -1.$$

Por lo tanto, puesto que  $p-p' = 2k_1$ ,  $p+p' = 2k_2$  con  $k_1 \in \{\lfloor \frac{-N+2}{2} \rfloor, \dots, \lfloor \frac{N-2}{2} \rfloor\}$  y  $k_2 \in \{1, 2, \dots, N-1\}$ , entonces

$$\sum_{j=1}^{N-1} \cos\left(\frac{(p-p')\pi j}{N}\right) = \sum_{j=1}^{N-1} \cos\left(\frac{(p+p')\pi j}{N}\right) = -1.$$

Concluimos que en cualquiera de los dos casos

$$\sum_{k=1}^{N-1} \sin\left(\frac{p\pi k}{I}\right) \sin\left(\frac{p'\pi k}{I}\right) = \frac{1}{2} \sum_{k=1}^{N-1} \left( \cos\left(\frac{(p-p')\pi k}{N}\right) - \cos\left(\frac{(p+p')\pi k}{N}\right) \right) = 0.$$

□

En resumen, lo que tenemos es que los vectores  $V_{h,p,q}$  forman una base de vectores propios de la matriz de iteración del método de Jacobi, cuyos autovalores son  $\mu_{h,p,q}$ , con  $p, q \in \{1, \dots, N-1\}$ .

Para obtener el radio espectral de la matriz  $B_h$ , buscamos los  $p, q$  que hagan lo más grande posible la siguiente cantidad

$$|\mu_{h,p,q}| = \frac{1}{2} \left| \cos\left(\frac{p\pi}{N}\right) + \cos\left(\frac{q\pi}{N}\right) \right|,$$

y esto lo conseguimos cuando  $p = 1$ ,  $q = 1$ . Concluimos que el radio espectral de la matriz de iteración del método de Jacobi es

$$\rho(B_h) = \mu_{h,1,1} = \frac{1}{2} \left( \cos\left(\frac{\pi}{N}\right) + \cos\left(\frac{\pi}{N}\right) \right) = \cos\left(\frac{\pi}{N}\right) = 1 - \frac{h^2}{2} \pi^2 + O(h^4).$$

Por lo tanto, si  $h \rightarrow 0$  el radio espectral de la matriz  $B$  se acerca asintóticamente a 1, lo que implica una convergencia lenta del método cuando  $h$  es pequeño (esto es, cuando hacemos la malla más fina).

**Observación 1.2.1** *La razón de convergencia de un método iterativo es  $R = -\log_{10}(\rho(M))$  donde  $M$  es la matriz de iteración del método. Lo que la razón de convergencia nos da es que, cuando el número de iteraciones es lo suficientemente grande, una iteración del método divide la norma del error por un factor de  $10^R$ . Por lo tanto, para dividir el error por  $10^m$  son necesarias  $m/R$  iteraciones.*

*En nuestro caso particular, haciendo desarrollo de Taylor del logaritmo en torno al 0, tenemos que*

$$R(B_h) = -\log_{10}(\rho(B_h)) = \frac{h^2}{2}\pi^2 + O(h^4)$$

*Por lo tanto, para dividir el error entre  $10^m$  en general harán falta aproximadamente  $2m/(h^2\pi^2)$  iteraciones cuando  $h$  tiende a 0.*

### 1.2.2. Convergencia del método de Gauss-Seidel

En las mismas condiciones del método de Jacobi, tendremos que los iterantes de Gauss-Seidel vienen dados por la siguiente expresión:

$$\begin{aligned} u_h^{n+1}(x, y) &= \frac{u_h^n(x+h, y) + u_h^n(x, y+h) + u_h^{n+1}(x-h, y) + u_h^{n+1}(x, y-h)}{4} \\ &\quad - \frac{h^2}{4} f_h^\Omega(x, y), \quad (x, y) \in \Omega_h, \\ u_h^{n+1}(x, y) &= f_h^\Gamma(x, y), \quad (x, y) \in \Gamma_h. \end{aligned} \tag{1.7}$$

Por tanto, los errores del método de Gauss-Seidel cumplen

$$\begin{aligned} e_h^{n+1}(x, y) &= \frac{e_h^n(x+h, y) + e_h^n(x, y+h) + e_h^{n+1}(x-h, y) + e_h^{n+1}(x, y-h)}{4}, \\ &\quad (x, y) \in \Omega_h, \\ e_h^{n+1}(x, y) &= 0, \quad (x, y) \in \Gamma_h. \end{aligned} \tag{1.8}$$

Formando, de nuevo, los vectores de  $\mathbb{R}^K$  cuyas entradas son las evaluaciones de los errores en los puntos de la malla, tendremos la siguiente expresión matricial

$$E_h^{n+1} = L_h^B E_h^{n+1} + U_h^B E_h^n$$

donde la matriz  $L_h^B$  es estrictamente triangular inferior y  $U_h^B$  es estrictamente triangular superior y de forma que verifican  $L_h^B + U_h^B = B_h$ . Por lo tanto, la matriz de iteración del método viene dada por  $\mathcal{L}_h = (I - L_h^B)^{-1} U_h^B$  y buscamos los autovalores  $\lambda_h$  y autovectores  $w_h$  de forma que  $\mathcal{L}_h w_h = \lambda_h w_h$ , es decir,  $\lambda_h w_h = \lambda_h L_h^B w_h + U_h^B w_h$ . A esta última expresión le corresponde la siguiente

ecuación

$$\lambda_h w_h(x, y) = \frac{w_h(x+h, y) + w_h(x, y+h) + \lambda_h w_h(x-h, y)}{4} + \frac{\lambda_h w_h(x, y-h)}{4}, \quad (x, y) \in \Omega_h, \quad (1.9)$$

y para los puntos en  $\Gamma_h$  tenemos  $w_h(x, y) = 0$ .

Simplemente sustituyendo, tenemos que si  $v_h(x, y)$  y  $\mu_h$  verifican la ecuación que teníamos para el método de Jacobi (1.4), entonces

$$w_h(x, y) = \lambda_h^{\frac{x+y}{2h}} v_h(x, y)$$

satisface la ecuación que tenemos para Gauss-Seidel (1.9) con  $\lambda_h = \mu_h^2$ .

Concluimos que los autovalores de la matriz  $\mathcal{L}_h$  son

$$\lambda_h = \lambda_{h,p,q} = \mu_{h,p,q}^2 = \frac{1}{4} \left( \cos\left(\frac{p\pi}{N}\right) + \cos\left(\frac{q\pi}{N}\right) \right)^2.$$

Quedándonos con el mayor en valor absoluto ( $p = 1, q = 1$ ), tenemos que el radio espectral de la matriz de iteración del método de Gauss-Seidel es

$$\rho(\mathcal{L}_h) = \cos^2(\pi h).$$

Terminamos esta sección advirtiendo del hecho de que el coseno de un ángulo al cuadrado será menor o igual que el coseno del mismo ángulo. Por lo tanto, el radio espectral de la matriz de iteración del método de Jacobi es mayor que el de la matriz de iteración de Gauss-Seidel, con lo que la convergencia en Gauss-Seidel es más rápida.

Notemos además que

$$\rho(\mathcal{L}_h) = 1 - (h\pi)^2 + O(h^4),$$

con lo que  $R(\mathcal{L}_h) = (\pi h)^2 + O(h^4)$  y siguiendo la observación 1.2.1, hacen falta  $m/(h\pi)^2$  iteraciones aproximadamente para dividir el error entre  $10^m$ .

### 1.2.3. Convergencia del método S.O.R.

Para un parámetro  $\omega \neq 0$  dado, el método S.O.R. para nuestro problema modelo viene dado por

$$u_h^{n+1}(x, y) = \frac{\omega}{4} [u_h^n(x+h, y) + u_h^n(x, y+h) + u_h^{n+1}(x-h, y) + u_h^{n+1}(x, y-h) - f_h^\Omega(x, y)] + (1-\omega)u_h^n(x, y), \quad (x, y) \in \Omega_h$$

$$u_h^{n+1}(x, y) = f_h^\Gamma(x, y), \quad (x, y) \in \Gamma_h.$$

En esta sección, buscaremos el  $\omega$  óptimo que nos dé la mejor convergencia para el método S.O.R. y además estudiaremos, para ese  $\omega$  óptimo, el comportamiento del radio espectral cuando  $h$  se hace pequeño.

Realizando los mismos pasos que en las secciones anteriores, tendremos que la ecuación del error para el método S.O.R. es

$$e_h^{n+1}(x, y) = \frac{\omega}{4}[e_h^n(x+h, y) + e_h^n(x, y+h) + e_h^{n+1}(x-h, y) + e_h^{n+1}(x, y-h)] \\ + (1-\omega)e_h^n(x, y),$$

que nos da la expresión matricial

$$E_h^{n+1} = \omega L_h^B E_h^{n+1} + ((1-\omega)I + \omega U_h^B) E_h^n$$

donde las matrices  $L_h^B$  y  $U_h^B$  son las que aparecen anteriormente en la descomposición de la matriz de iteración del método Jacobi  $B_h$ . Denotamos por  $\mathcal{L}_{\omega, h}$  a la matriz de iteración del método S.O.R., donde

$$\mathcal{L}_{\omega, h} = (I - \omega L_h^B)^{-1}((1-\omega)I + \omega U_h^B).$$

Nos centramos en buscar entonces los autovectores  $V_h$  y autovalores  $\lambda_h$  de la matriz  $\mathcal{L}_{\omega, h}$ , es decir,  $\lambda_h$  y  $V_h \neq 0$  que cumplen

$$\mathcal{L}_{\omega, h} V_h = \lambda_h V_h,$$

o lo que es lo mismo,

$$\lambda_h V_h = \omega \lambda_h L_h^B V_h + ((1-\omega)I + \omega U_h^B) V_h.$$

Esto puede escribirse en términos de la función nodal  $v_h$  como sigue

$$\lambda_h v_h(x, y) = \frac{\omega}{4}[v_h(x+h, y) + v_h(x, y+h) + \lambda_h v_h(x-h, y) + \lambda_h v_h(x, y-h)] \\ + (1-\omega)v_h(x, y), \quad (x, y) \in \Omega_h, \\ v_h(x, y) = 0, \quad (x, y) \in \Gamma_h,$$

o equivalentemente, puesto que  $\omega \neq 0$ ,

$$\frac{(\lambda_h + \omega - 1)}{\omega} v_h(x, y) = \frac{1}{4}[v_h(x+h, y) + v_h(x, y+h) + \lambda_h v_h(x-h, y) + \lambda_h v_h(x, y-h)].$$

Haciendo entonces el cambio de variable  $v_h(x, y) = \lambda_h^{\frac{x+y}{2h}} w_h(x, y)$ , esto es equivalente a encontrar  $\omega$  tal que

$$\frac{(\lambda_h + \omega - 1)}{\omega \lambda_h^{1/2}} w_h(x, y) = \frac{1}{4}[w_h(x+h, y) + w_h(x, y+h) + w_h(x-h, y) + w_h(x, y-h)].$$

Esta última ecuación nos da de nuevo una relación entre los autovalores de la matriz de iteración del método de Jacobi y los autovalores de la matriz de iteración del método S.O.R. y es que, por cada autovalor  $\mu_h$  de la matriz de iteración del método de Jacobi, tenemos dos del método S.O.R. teniendo en cuenta que

$$\frac{(\lambda_h + \omega - 1)}{\omega \lambda_h^{1/2}} = \mu_h \Leftrightarrow \lambda_h - \lambda_h^{1/2} \omega \mu_h + (\omega - 1) = 0, \quad (1.10)$$

que es una ecuación cuadrática en  $\lambda_h^{1/2}$ .

**Observación 1.2.2** *Antes de seguir con el estudio de la ecuación cuadrática veamos que, para que la matriz de iteración  $\mathcal{L}_{\omega,h}$  sea no singular debe ocurrir que  $\omega \neq 1$ . Esto se deduce del hecho de que*

$$\det(\mathcal{L}_{\omega,h}) = \det(I_h - \omega L_h^B)^{-1} \det((1 - \omega)I_h + \omega U_h^B) = (1 - \omega)^K,$$

donde  $K = (N - 1)^2$  es la dimensión de la matriz. Además, una condición necesaria para la convergencia es, obviamente, que  $0 < \omega < 2$ .

Resolviendo la ecuación (1.10),

$$\lambda_h = \left( \frac{\omega\mu_h \pm \sqrt{(\omega\mu_h)^2 - 4(\omega - 1)}}{2} \right)^2. \quad (1.11)$$

Como antes advertimos, (1.11) nos da, por cada autovalor de la matriz de iteración del método Jacobi, dos de la matriz de iteración del método S.O.R. Ahora bien, los autovalores de la matriz de iteración del método de Jacobi verifican que  $\mu_h^{pq} = -\mu_h^{(N-p)(N-q)}$ . En consecuencia, existe una correspondencia entre los pares de soluciones de la ecuación (1.10) y los autovalores  $\mu_h^{pq}, \mu_h^{(N-p)(N-q)}$ .

Ahora vamos a buscar el  $\omega$  óptimo que minimice  $|\lambda_h^{1/2}|$  cuando  $\lambda_h^{1/2}$  es real. Para ello nos quedamos con la raíz cuadrada positiva en (1.11) y supondremos, sin pérdida de generalidad, que  $\mu_h > 0$ .

Observamos que  $\lambda_h^{1/2}$  es real, cuando

$$(\omega\mu_h)^2 - 4(\omega - 1) = \left( \omega\mu_h - \frac{2}{\mu_h} \right)^2 - 4 \left( \frac{1}{\mu_h^2} - 1 \right) \geq 0.$$

Para determinar cómo  $\lambda_h^{1/2}$  varía, en función de  $\omega$ , derivamos y obtenemos

$$\begin{aligned} \frac{\partial}{\partial \omega} \lambda_h^{1/2} &= \frac{1}{2}\mu_h + \frac{1}{2}(\omega\mu_h^2 - 2)(\omega^2\mu_h^2 - 4\omega + 4)^{-1/2} \\ &= \frac{\mu_h}{2} \left( 1 - \frac{\frac{2}{\mu_h} - \mu_h\omega}{\sqrt{(2/\mu_h - \omega\mu_h)^2 - 4(\mu_h^{-2} - 1)}} \right) < 0. \end{aligned} \quad (1.12)$$

La última desigualdad se debe al hecho de que  $2/\mu_h - \mu_h\omega > 0$ , que es equivalente, por ser  $\mu_h > 0$ , a  $\omega < 2/\mu_h^2$ , lo cual se cumple para  $0 < \omega < 2$ . Además,  $\sqrt{(2/\mu_h - \omega\mu_h)^2 - 4(\mu_h^{-2} - 1)} < 2/\mu_h - \mu_h\omega$  (para convencerse de esto basta con observar que lo que hay dentro de la raíz cuadrada es el cuadrado de  $2/\mu_h - \omega\mu_h$  menos una cantidad positiva, pues  $0 < \mu_h < 1$ . La monotonía de la función raíz cuadrada de  $x$  nos ayuda a concluir).

Como la derivada es negativa tenemos que la función decrece en  $\omega$ . Entonces, para minimizar  $\lambda_h^{1/2}$  debemos avanzar en  $\omega$  hacia valores más grandes. Buscaremos el mayor valor de  $\omega$  para el que  $\lambda_h^{1/2}$  sea real y cumpliendo  $0 < \omega < 2$ .

18CAPÍTULO 1. PROBLEMA MODELO Y MÉTODOS ITERATIVOS CLÁSICOS

Puesto que  $\omega^2\mu_h^2 - 4\omega + 4$  es una parábola en  $\omega$  con un mínimo y las dos raíces de dicha parábola son

$$\omega = \frac{2}{\mu_h^2} \left( 1 \pm \sqrt{1 - \mu_h^2} \right),$$

se tiene que, como  $\mu_h < 1$ , la asociada al signo positivo es mayor que 2. Por tanto, nos interesa la otra raíz, que también puede reescribirse como

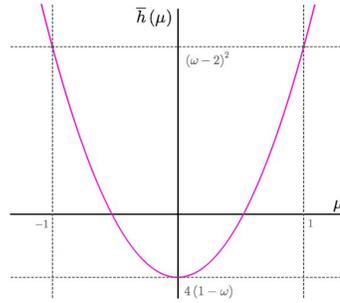
$$\omega = \frac{2}{1 + \sqrt{1 - \mu_h^2}}. \quad (1.13)$$

Tras este estudio cuando  $\lambda_h^{1/2}$  es real, vamos a buscar el  $\omega$  óptimo cuando  $\lambda_h^{1/2}$  es complejo. Puesto que (1.10) tiene coeficientes reales, las expresiones entre paréntesis en (1.11) correspondientes a  $-\mu_h$  y  $\mu_h$  son complejos conjugados. Además

$$|\lambda_h| = |\lambda_h^{1/2}|^2 = \frac{1}{4} [\omega^2\mu_h^2 + 4(\omega - 1) - \omega^2\mu_h^2] = \omega - 1.$$

Por lo tanto, en el caso complejo, para reducir el valor de  $|\lambda_h^{1/2}|$  debemos reducir el valor de  $\omega$ . Con lo cual, buscaremos el mínimo valor de  $\omega$  para el cual  $\lambda_h^{1/2}$  sea complejo. Esto es, de nuevo, la raíz más pequeña de  $\omega^2\mu_h^2 - 4\omega + 4 = 0$  que además verifica  $0 < \omega < 2$ , tal y como nos interesa.

Para estudiar el radio espectral de la matriz de iteración del método S.O.R. con el  $\omega$  óptimo que lo minimice, tengamos en cuenta que, para un  $\omega$  fijo,  $\bar{h}(\mu) = \omega^2\mu^2 - 4\omega + 4$  es la siguiente parábola en  $(-1, 1)$ :



Por tanto, para un  $\omega$  cercano a 2,  $\bar{h}(\mu_h)$  es negativo para todos los valores de  $\mu_h \in (-1, 1)$ . De aquí se sigue que todos los  $\lambda_h(\mu_h)$  son complejos con módulo  $(\omega - 1)$  y, por tanto, concluimos que el radio espectral es  $(\omega - 1)$ . Para hacer  $(\omega - 1)$  pequeño, vamos haciendo que  $\omega$  sea más pequeño hasta que llegamos a un  $\omega$  que hace real algún  $\lambda_h(\mu_h)$ . Esto ocurre cuando  $\mu_h$  es el autovalor con mayor módulo de la matriz de iteración de Jacobi. Más concretamente, cuando  $\omega^2\mu_{h,1,1}^2 - 4(\omega - 1) = 0$  que es el valor de  $\omega$  en (1.13), con  $\mu_h$  sustituido por  $\mu_{h,1,1}$ , más concretamente,

$$\omega_h^* = \frac{2}{1 + \sqrt{1 - \mu_{h,1,1}}}$$

Para valores más pequeños de  $\omega$ , el autovalor real más grande de  $\mathcal{L}_{\omega,h}$  seguirá siendo el correspondiente a  $\mu_{h,1,1}$  en (1.11) con signo positivo, y este será más grande teniendo en cuenta (1.12).

Como  $\omega_h^* > 1$ , siempre habrá algún autovalor complejo de  $\mathcal{L}_{\omega_h^*,h}$  que sabemos que toman módulo  $\omega_h^* - 1$ . Basta ahora darse cuenta de que, para ese valor de  $\omega_h^*$ ,

$$|\lambda_{h,1,1}| = \left( \frac{\omega_h^*}{2} \mu_{h,1,1} \right)^2 = \omega_h^* - 1,$$

con lo que queda probado que este  $\omega_h^*$  es el óptimo.

Teniendo en cuenta la expresión de  $\mu_{h,1,1} = \cos(\pi/N)$ , entonces

$$\omega_h^* = \frac{2}{1 + \sin(\frac{\pi}{N})}.$$

Para terminar, vemos que el radio espectral de la matriz de iteración del método S.O.R. para el  $\omega$  óptimo es

$$\rho(\mathcal{L}_{\omega_h^*,h}) = \omega_h^* - 1 = \frac{1 - \sin(\frac{\pi}{N})}{1 + \sin(\frac{\pi}{N})} \approx 1 - \frac{2\pi}{N} = 1 - 2h\pi + O(h^2).$$

De nuevo observamos, como en los otros métodos, que cuando  $h \rightarrow 0$ , la convergencia será más lenta. Pero, como es lógico, para  $h$  pequeño, el radio de convergencia es claramente menor que el correspondiente al método de Gauss-Seidel. Más concretamente, para dividir el error entre  $10^m$ , de acuerdo con la observación 1.2.1, se necesitan en general aproximadamente  $m/(2h\pi)$  iteraciones.



## Capítulo 2

# Método multigríd

### 2.1. Ideas principales

Utilizando el esquema del método de Gauss-Seidel con el orden lexicográfico (1.7) sobre el problema modelo ocurre que el error se hace cada vez más suave. Esto no quiere decir que se haga cada vez más pequeño, si no que la diferencia de los errores en nodos distintos será cada vez más pequeño. En la Figura de más abajo se pretende mostrar la idea del efecto del suavizado sobre el error. De hecho, si nos fijamos en la fórmula del error para Gauss-Seidel (1.8), podemos interpretarla como un promedio de los errores en los puntos ya aproximados. Esto nos lleva a uno de las ideas fundamentales con las que trabaja el método multigríd:

- **Principio de suavizado.** Este principio se basa en que la mayoría de los métodos clásicos iterativos aplicados sobre problemas elípticos discretos tienen un efecto de suavizado sobre el error de la aproximación.

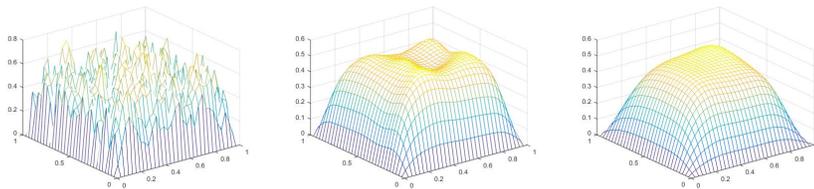


Figura 2.1: Representación del error utilizando el método de Gauss-Seidel con el orden lexicográfico sobre el problema modelo con  $h = (1/2)^5$ . Error con iterante inicial aleatorio, tras una iteración, tras 10 iteraciones y tras 20 iteraciones.

Una vez que tenemos una función que es suave en una cierta red de puntos podremos, sin una gran pérdida de información, aproximarla en una red más gruesa. Es decir, puesto que hay poca diferencia entre los valores de la función en

los distintos puntos de la red, podremos aproximar dicha función en un punto utilizando el valor de la función en puntos más alejados. Esto nos lleva a la segunda idea básica de los métodos multigrid:

- **Principio de engrosamiento de la red.** Los errores que sufren el efecto del suavizado tras varias iteraciones pueden ser aproximados en una red más gruesa.

Una ventaja clara del engrosamiento de la red es que, aproximar en una red más gruesa, tiene un coste menor que en la red más fina.

Otra de las ideas fundamentales del método multigrid es la que mostraremos a continuación observando los desarrollos de Fourier del error. En nuestro problema modelo, el error  $e_h = e_h^m(x, y)$ , con  $(x, y) \in \Omega_h$ , puede ser escrito de la forma

$$e_h(x, y) = \sum_{k,l=1}^{N-1} \alpha_{k,l} \sin(k\pi x) \sin(l\pi y), \quad (2.1)$$

para ciertas constantes  $\alpha_{k,l}$  ya que, como vimos en la sección 1.2.1, los vectores correspondientes a los valores nodales de las funciones

$$v_{k,l}(x, y) = \sin(k\pi x) \sin(l\pi y)$$

son linealmente independientes. De hecho, dichas funciones, aparte de ser autofunciones de la matriz de iteración de Jacobi  $B_h$ , también son autofunciones del operador discreto  $\Delta_h$  ya que

$$\Delta_h = \frac{4}{h^2}(B_h - I).$$

Claramente, los autovalores de  $\Delta_h$  se deducen inmediatamente de los de  $B_h$  en (1.6) y son

$$\frac{1}{h^2}(2 \cos(l\pi h) + 2 \cos(k\pi h) - 4).$$

A continuación definimos las componentes de alta y baja frecuencia del error.

**Definición 2.1.1** Para  $k, l \in \{1, \dots, N-1\}$ , llamaremos componentes de baja frecuencia a los términos  $v_{k,l}$  tales que  $\max(k, l) < N/2$ . De la misma forma, diremos que  $v_{k,l}$  es una componente de alta frecuencia si  $N/2 \leq \max(k, l) < N$ .

Esto nos permite tratar, de forma sucinta, el último principio en el que se basa el método multigrid:

- **Principio de corrección en red gruesa.** Basado en el hecho de que si el error se suaviza, esto significa que las componentes de alta frecuencia en (2.1) se hacen pequeñas tras varias iteraciones. Sin embargo, las componentes de baja frecuencia apenas se verán afectadas.

En nuestro problema modelo en la red  $\Omega_h$ , con  $h = 1/N$ , si  $N$  es un número par, consideramos la red más gruesa  $\Omega_H$ , donde  $H = 2h$ . Esta elección de engrosamiento de la red es conocida por engrosamiento estándar y será la que utilizaremos a lo largo de este trabajo.

Observamos que las autofunciones  $v_{k,l}$  verifican las siguiente igualdades en  $\Omega_{2h}$

$$v_{k,l}(x, y) = -v_{N-k,l}(x, y) = -v_{k,N-l}(x, y) = v_{N-k,N-l}(x, y), \quad (x, y) \in \Omega_{2h}.$$

Demostremos la igualdad  $v_{k,l}(x, y) = v_{N-k,N-l}(x, y)$ , con  $(x, y) \in \Omega_{2h}$ , para mostrar las ideas con las que se prueba el resto de igualdades. Primero, utilizando la fórmula trigonométrica del seno de la suma de ángulos tendremos:

$$\begin{aligned} v_{N-k,N-l}(x, y) &= (\sin(N\pi x) \cos(k\pi x) - \cos(N\pi x) \sin(k\pi x)) \\ &\quad (\sin(N\pi y) \cos(l\pi y) - \cos(N\pi y) \sin(l\pi y)). \end{aligned}$$

Ahora, teniendo en cuenta que  $h = 1/N$  y que los puntos  $(x, y) \in \Omega_{2h}$  son aquellos de la forma  $x = 2ih$ ,  $y = 2jh$ , para  $i, j \in \{1, \dots, (N-2)/2\}$ , entonces  $x = 2i/N$ ,  $y = 2j/N$ . Sustituyendo, y teniendo en cuenta que  $\sin(2\pi i) = \sin(2\pi j) = 0$  y  $\cos(2\pi i) = \cos(2\pi j) = 1$ , en la expresión de arriba llegamos a la igualdad que buscábamos.

Con esto lo que tenemos es que, en  $\Omega_{2h}$ , estas autofunciones son indistinguibles, y por lo tanto, tenemos que solo las componentes de baja frecuencia son visibles en  $\Omega_{2h}$ , ya que las de alta frecuencia coinciden con las de baja frecuencia (a esto se le conoce como aliasing de frecuencias).

Hasta ahora hemos descrito las bases del método multigrad en dos redes. Vamos a introducir las nociones del método en más de dos redes y, para ello, supondremos que  $N$  es una potencia de 2. Por lo tanto,  $h = 2^{-p}$  y denotaremos las redes de nodos de la forma

$$\Omega_h, \Omega_{2h}, \Omega_{4h}, \dots, \Omega_{h_0},$$

donde  $\Omega_h$  es la red más fina y  $\Omega_{h_0}$  es la red más gruesa formada por un único punto.

De la misma forma que hemos distinguido componentes de alta y baja frecuencia en las redes  $(\Omega_h, \Omega_{2h})$ , ahora haremos lo mismo con las redes  $(\Omega_{2h}, \Omega_{4h})$  y continuaremos con toda la secuencia de redes. Anteriormente comentamos que, usando la iteración de Gauss-Seidel en la red  $\Omega_h$ , generamos unas componentes de alta frecuencia  $(h, 2h)$  del error que se van haciendo cada vez más pequeñas. Las componentes de baja frecuencia  $(h, 2h)$  son visibles en  $\Omega_{2h}$  y pueden ser aproximadas en dicha red. La iteración de Gauss-Seidel afecta también al resto de redes y, por lo tanto, generamos componentes de alta frecuencia pequeñas  $(2h, 4h)$ , componentes de alta frecuencia  $(4h, 8h)$  pequeñas y así sucesivamente con el resto de redes.

En resumen, lo que se pretende destacar es la idea de que métodos iterativos adecuados, como puede ser Gauss-Seidel, darán una rápida reducción de las componentes de alta frecuencia de los errores en las diferentes mallas.

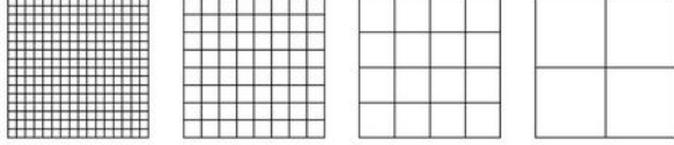


Figura 2.2: Secuencia de un engrosamiento estándar de mallas desde una malla de tamaño  $h = 1/16$ .

## 2.2. Procedimiento de suavizado del error

Como hemos mencionado anteriormente, los métodos iterativos clásicos tienen un efecto de suavizado sobre el error de una aproximación. Este efecto de suavizado es una de las bases del método multigrad y, en esta sección, analizaremos el efecto de suavizado de los métodos iterativos clásicos de Jacobi y Gauss-Seidel.

A los métodos clásicos iterativos, como el de Jacobi y Gauss-Seidel, se les conoce como métodos de suavizado o métodos de relajación cuando son utilizados con la intención de suavizar errores.

### 2.2.1. Iteración tipo Jacobi

Para el problema modelo, la fórmula de iteración del método de Jacobi es

$$\begin{aligned} z_h^{m+1}(x_i, y_j) &= \frac{1}{4} [-h^2 f_h(x_i, y_j) + u_h^m(x_i - h, y_j) + u_h^m(x_i + h, y_j) \\ &\quad + u_h^m(x_i, y_j - h) + u_h^m(x_i, y_j + h)] \\ u_h^{m+1}(x_i, y_j) &= z_h^{m+1}(x_i, y_j), \end{aligned}$$

con  $(x_i, y_j) \in \Omega_h$  y donde  $u_h^m, u_h^{m+1}$  denotan la antigua y la nueva aproximación respectivamente.

Podemos generalizar esta iteración introduciendo un parámetro de relajación  $\omega$ :

$$u_h^{m+1} = u_h^m + \omega(z_h^{m+1} - u_h^m).$$

A la relajación del método de Jacobi con el parámetro  $\omega$  lo denotaremos por  $\omega$ -JAC. Para  $\omega = 1$  tendremos que  $\omega$ -JAC es el método de Jacobi tradicional. El método iterativo  $\omega$ -JAC puede ser escrito como

$$u_h^{m+1} = S_h(\omega)u_h^m - \frac{h^2}{4}f_h,$$

donde el operador  $S_h(\omega)$  viene dado por

$$S_h(\omega) = I_h + \frac{\omega h^2}{4} \Delta_h. \quad (2.2)$$

Para un estudio de la convergencia del método  $\omega$ -JAC analizaremos los autovalores de las autofunciones del operador  $S_h(\omega)$ . Notemos que las autofunciones del operador  $S_h(\omega)$  son las mismas que las del operador  $\Delta_h$ , y teniendo en cuenta quiénes eran los autovalores de  $\Delta_h$ , los autovalores del operador de  $\omega$ -JAC son

$$\chi_h^{k,l}(\omega) = 1 - \frac{\omega}{2} [2 - \cos(k\pi h) - \cos(l\pi h)],$$

con  $k, l \in \{1, \dots, N-1\}$ . El radio espectral del operador,

$$\rho(S_h(\omega)) = \text{máx}\{|\chi_h^{k,l}(\omega)| : (k, l = 1, \dots, N-1)\},$$

vendrá dado por

$$\begin{aligned} \rho(S_h(\omega)) &= |\chi_h^{1,1}(\omega)| = |1 - \omega(1 - \cos(\pi h))| = 1 - |O(\omega h^2)|, \quad \omega \in (0, 1], \\ \rho(S_h(\omega)) &\geq 1, \quad \omega \in \mathbb{R} \setminus (0, 1], \text{ para } h \text{ suficientemente pequeño.} \end{aligned}$$

Por lo tanto, podemos concluir que, para tener una buena convergencia para el método  $\omega$ -JAC, tomar  $\omega = 1$  es la mejor opción.

Ahora bien, para las propiedades de suavizado la situación es muy diferente. Consideraremos las aproximaciones antes ( $w_h$ ) y después ( $\bar{w}_h$ ) de aplicar el método de suavizado  $\omega$ -JAC, para  $\omega \in (0, 1]$ , y los desarrollos de Fourier de sus errores, esto es

$$v_h = \sum_{k,l=1}^{N-1} \alpha_{k,l} v_{k,l}, \quad \bar{v}_h = \sum_{k,l=1}^{N-1} \chi_h^{k,l}(\omega) \alpha_{k,l} v_{k,l},$$

donde  $v_h := u_h - w_h$  y  $\bar{v}_h := u_h - \bar{w}_h$ . Recordemos que, en el análisis en el que estamos interesados ahora, estamos buscando que reduzcan muy rápido las componentes de alta frecuencia ( $h, 2h$ ). Sin embargo, cuando analizamos la convergencia del método, buscamos un  $\chi_h^{k,l}(\omega)$  de forma que todas las componentes de  $\bar{v}_h$  sean lo más pequeñas posibles tras cada iteración.

Para analizar las propiedades de suavizado del operador  $S_h(\omega)$  cuantitativamente, definimos el siguiente concepto.

**Definición 2.2.1** *Llamaremos factor de suavizado  $\mu(h; \omega)$  de  $S_h(\omega)$  al valor*

$$\mu(h; \omega) := \text{máx}\{|\chi_h^{k,l}(\omega)| : N/2 \leq \text{máx}(k, l) \leq N-1\}.$$

*Este valor representa el peor factor por el cual las componentes de alta frecuencia del error son reducidos tras una iteración. Denotaremos por  $\mu^*(\omega)$  al supremo de, entre todas las posibles elecciones de  $h$ , los factores de suavizado de  $S_h(\omega)$ , es decir*

$$\mu^*(\omega) := \sup_{h \in \mathbf{H}} \mu(h; \omega),$$

donde  $\mathbf{H}$  denota el conjunto de todos los tamaños de malla admisibles.

A continuación, mostraremos cómo el  $\omega$  que hace óptima la convergencia del método  $\omega$ -JAC no es la mejor elección para tener buenas propiedades de suavizado, es decir, no nos dará el factor de suavizado más pequeño.

Denotamos  $\mathbf{H} = \{h = 1/n : n \in \mathbb{N}, n \geq 2\}$ .

Por lo tanto, el factor de suavizado y su supremo sobre los posibles valores de los tamaños de la malla vienen dados por

$$\mu(h; \omega) = \max\left\{1 - \frac{\omega}{2}(2 - \cos(k\pi h) - \cos(l\pi h)) : N/2 \leq \max(k, l) \leq N - 1\right\},$$

$$\mu^*(\omega) = \max\left\{1 - \frac{\omega}{2}, |1 - 2\omega|\right\}.$$

La idea de cómo obtener  $\mu^*$  es que, tomando  $h$  suficientemente pequeño, los cosenos de la expresión anterior se acercarán al valor  $-1$  cuando tomemos  $k, l = N - 1$  y al valor  $0$  y  $1$  cuando tomemos, por ejemplo,  $k = N/2$  y  $l = 1$ .

Observemos que, para un  $h$  suficientemente pequeño,  $\mu(h; \omega) \geq 1$ , para  $\omega \leq 0$  y  $\omega \geq 1$ . Por lo tanto, no tendremos propiedades de suavizado en estos casos.

Para  $\omega \in (0, 1)$ , veamos que  $\omega = 4/5$  es la elección óptima para tener las mejores propiedades de suavizado para el método  $\omega$ -JAC en este modelo. Es decir, buscaremos  $\omega \in (0, 1)$  que haga lo más pequeña posible la expresión  $\max\{1 - \omega/2, |1 - 2\omega|\}$ . Viendo la figura 2.3 anteriores tenemos que la función

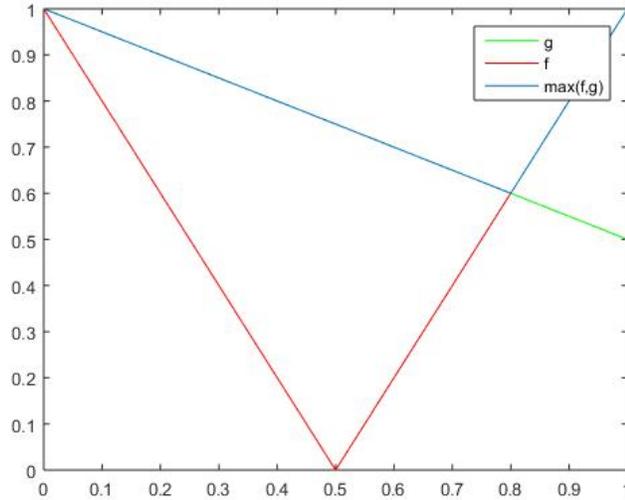


Figura 2.3: Representación gráfica de las funciones  $f(\omega) = |1 - 2\omega|$  (en rojo),  $g(\omega) = |1 - \omega/2|$  (en verde) y  $\max(f, g)$  (en azul) con  $\omega \in (0, 1)$ .

$\max\{1 - \omega/2, |1 - 2\omega|\}$  alcanza su valor más pequeño cuando  $1 - \omega/2 = 2\omega - 1$ . Despejando obtenemos fácilmente que  $\omega = 4/5$ .

Para el factor de suavizado  $\mu(h; \omega)$  veamos que

$$\inf\{\mu(h; \omega) : 0 \leq \omega \leq 1\} = \mu\left(h; \frac{4}{4 + \cos(\pi h)}\right) = \frac{3 \cos(\pi h)}{4 + \cos(\pi h)} = \frac{3}{5} - |O(h^2)|.$$

El ínfimo lo obtenemos razonando igual que hicimos en la gráfica anterior, salvo que ahora, con todas las funciones de la forma  $|1 - \omega(2 - \cos(k\pi h) - \cos(l\pi h))/2|$ , con  $N/2 \leq \max(k, l) \leq N - 1$ . Es decir, el ínfimo se alcanza para el  $\omega$  que verifique

$$1 - \omega \left(1 - \frac{\cos(\pi h)}{2}\right) = -1 + \omega(1 + \cos(\pi h)).$$

Todo esto lo que significa es que,  $\omega$ -JAC con  $\omega = 4/5$  reduce todas las componentes de alta frecuencia al menos por un factor de  $3/5$  independientemente del tamaño de la malla  $h$ .

### 2.2.2. Iteración tipo Gauss-Seidel

En el caso del método iterativo de Gauss-Seidel, si introducimos un parámetro de relajación obtendremos el ya mencionado método S.O.R. Para el problema modelo el esquema del método vendrá dado por

$$\begin{aligned} z_h^{m+1}(x_i, y_j) &= \frac{1}{4} [-h^2 f_h(x_i, y_j) + u_h^{m+1}(x_i - h, y_j) + u_h^m(x_i + h, y_j) \\ &\quad + u_h^{m+1}(x_i, y_j - h) + u_h^m(x_i, y_j + h)] \\ u_h^{m+1} &= u_h^m + \omega(z_h^{m+1} - u_h^m), \end{aligned}$$

con  $(x_i, y_j) \in \Omega_h$  y donde  $u_h^m, u_h^{m+1}$  denotan la antigua y la nueva aproximación respectivamente. Ya vimos en el capítulo anterior que  $\omega_h^* = 2/(1 + \sin(\pi/N))$  es el parámetro de relajación que mejor convergencia del método nos da. De hecho, obteníamos que el radio espectral de la matriz de iteración del método S.O.R. asociado es

$$\rho(\mathcal{L}_{\omega_h^*, h}) = 1 - O(h),$$

frente al radio espectral de la matriz de iteración método clásico de Gauss-Seidel que era

$$\rho(\mathcal{L}_h) = 1 - O(h^2).$$

Sin embargo, en el contexto de los métodos multigrid, las propiedades de suavizado de los métodos tipo Gauss-Seidel son más interesantes que sus propiedades de convergencia a la solución. A modo de resumen, antes de entrar a analizar dichas propiedades de suavizado, destaquemos que, en nuestro problema modelo con el orden lexicográfico, la introducción de un parámetro de suavizado no mejorará las propiedades de suavizado del método de Gauss-Seidel tradicional. Realizar un análisis del suavizado como el considerado para el método de  $\omega$ -JAC no es tan fácil y tenemos que recurrir al análisis local de Fourier, que considera una malla infinita para no tener en cuenta las condiciones en la frontera.

### 2.3. Análisis local de Fourier de las propiedades de suavizado

Consideraremos en este análisis las funciones definidas en una malla de puntos de la forma

$$\varphi_h(\theta, x) = e^{i\theta x/h},$$

donde  $x$  varía dentro de una malla de puntos infinita  $G_h$  y  $\theta$  es un parámetro continuo que caracteriza la frecuencia de la función.

Puesto que las funciones  $\varphi(\theta, \cdot)$  están definidas en una red infinita  $G_h$ , la frontera y las condiciones frontera de nuestro problema no se están teniendo en cuenta. Sin embargo, el objetivo del análisis local de Fourier, es determinar el comportamiento cuantitativo de la convergencia y eficiencia de los algoritmos multigrad si se incluye un tratamiento apropiado en la frontera. Ese tratamiento extra en la frontera es muy pequeño, en comparación con los puntos interiores de la malla, cuando el tamaño de la red tiende a 0.

#### 2.3.1. Primeros conceptos para el análisis del suavizado

Nos vamos a centrar en nuestro modelo con eje espaciado en dos dimensiones y el engrosamiento estándar de la malla. Denotemos por  $\mathbf{x} = (x_1, x_2)$  y por  $h$  a la magnitud que es el diámetro de la red en ambas direcciones. Para cada  $h$  tenemos una malla de infinitos puntos asociada

$$\mathbf{G}_h = \{\mathbf{x} = \mathbf{k}h := (k_1h, k_2h), \quad \mathbf{k} \in \mathbb{Z}^2\}.$$

Consideraremos un operador discreto  $L_h$  general sobre la malla de infinitos puntos  $\mathbf{G}_h$ , es decir,

$$L_h w_h(\mathbf{x}) = \sum_{\kappa \in V} s_\kappa w_h(\mathbf{x} + \kappa h) \quad (\kappa \in \mathbb{Z}^2), \quad (2.3)$$

donde  $s_\kappa$  son coeficientes constantes que estarán en  $\mathbb{R}$  o en  $\mathbb{C}$  y que generalmente dependerán de  $h$ . Además,  $V$  es un conjunto finito de índices.

En el análisis local de Fourier será muy importante el papel que juegan las siguientes funciones definidas en la malla de infinitos puntos

$$\varphi_h(\boldsymbol{\theta}, \mathbf{x}) = e^{i\boldsymbol{\theta}\mathbf{x}/h} := e^{i\theta_1 x_1/h} e^{i\theta_2 x_2/h}, \quad \mathbf{x} \in G_h.$$

Supondremos que  $\boldsymbol{\theta}$  varía continuamente en  $\mathbb{R}^2$ . Observemos que  $\varphi_h(\boldsymbol{\theta}, \mathbf{x}) \equiv \varphi_h(\bar{\boldsymbol{\theta}}, \mathbf{x})$  si y solo si

$$\theta_1 = \bar{\theta}_1 \quad \text{mód } 2\pi \quad \text{y} \quad \theta_2 = \bar{\theta}_2 \quad \text{mód } 2\pi.$$

Por lo tanto, es suficiente considerar  $\varphi_h(\boldsymbol{\theta}, \mathbf{x})$  con  $\boldsymbol{\theta} \in [-\pi, \pi]^2$ . Tendremos en cuenta el siguiente lema

**Lema 2.3.1** *Las funciones  $\varphi_h(\boldsymbol{\theta}, \cdot)$  con  $\boldsymbol{\theta} \in [-\pi, \pi]^2$  son linealmente independientes en la malla  $G_h$ .*

### 2.3. ANÁLISIS LOCAL DE FOURIER DE LAS PROPIEDADES DE SUAVIZADO 29

#### **Demostración.**

Si tenemos un conjunto finito de ángulos  $\{\theta_1, \dots, \theta_n\} \subset [-\pi, \pi)^2$  y hacemos una combinación lineal de las funciones  $\varphi_h(\theta_j, \cdot)$  igual a 0, esto es

$$\sum_{j=1}^n \alpha_j \varphi_h(\theta_j, \mathbf{x}) = 0, \quad \text{para todo } \mathbf{x} \in \mathbf{G}_h,$$

o equivalentemente,

$$\sum_{j=1}^n \alpha_j e^{i\theta_{j,1}l} e^{i\theta_{j,2}m} = 0, \quad \text{para todo } (l, m) \in \mathbb{Z}^2, \quad (2.4)$$

entonces necesariamente  $\alpha_j = 0$  para todo  $j \in \{1, \dots, n\}$ . Esto es sencillo de demostrar, en el caso en que los  $\theta_{j,2}$  sean todas distintas, tomando los puntos de la red de la forma  $(0, 0), (0, h), (0, 2h), \dots, (0, (n-1)h)$ , ya que con ellos obtenemos el sistema lineal:

$$\begin{cases} \sum_{j=1}^n \alpha_j & = 0 \\ \sum_{j=1}^n \alpha_j z_j & = 0 \\ \vdots & \\ \sum_{j=1}^n \alpha_j z_j^{n-1} & = 0 \end{cases}$$

donde  $z_j = e^{i\theta_{j,2}}$ . Este sistema tiene asociada la siguiente matriz

$$\begin{pmatrix} 1 & 1 & \cdots & 1 \\ z_1 & z_2 & \cdots & z_n \\ \vdots & \vdots & \vdots & \vdots \\ z_1^{n-1} & z_2^{n-1} & \cdots & z_n^{n-1} \end{pmatrix}.$$

Puesto que  $z_i \neq z_j$  para  $i \neq j$ , tenemos que la matriz es de Vandermonde con determinante distinto de 0 y, por lo tanto, el sistema tiene una única solución, que es la trivial.

Si no son todas las componentes  $\theta_{j,2}$  distintas, habría que escribir (2.4) como

$$\sum_{j:\theta_{j,2} \text{ distintas}} e^{i\theta_{j,2}m} \left( \sum_{j:\theta_{j,2} \text{ es la elegida y } \theta_{j,1} \text{ distintos}} \alpha_j e^{i\theta_{j,1}l} \right) = 0 \quad \forall (l, m) \in \mathbb{Z}^2.$$

De aquí, por lo anterior, los paréntesis deben anularse y, aplicando de nuevo el razonamiento anterior con  $\theta_{j,1}$  en lugar de  $\theta_{j,2}$  se obtendría que  $\alpha_j = 0$ .  $\square$

Veamos la siguiente propiedad de las funciones  $\varphi_h(\theta, \mathbf{x})$  que recogeremos en el siguiente lema, cuya prueba es prácticamente inmediata.

**Lema 2.3.2** *Para  $-\pi \leq \theta < \pi$ , todas las funciones  $\varphi_h(\theta, \mathbf{x})$  son autofunciones de cualquier operador discreto que pueda ser descrito como en (2.3). Tendremos la siguiente relación*

$$L_h \varphi_h(\theta, \mathbf{x}) = \tilde{L}_h(\theta) \varphi_h(\theta, \mathbf{x}) \quad (\mathbf{x} \in \mathbf{G}_h),$$

con

$$\tilde{L}_h(\boldsymbol{\theta}) = \sum_{\boldsymbol{\kappa}} s_{\boldsymbol{\kappa}} e^{i\boldsymbol{\theta}\boldsymbol{\kappa}},$$

el autovalor formal de  $L_h$ .

Una vez vistas estas propiedades de las funciones  $\varphi_h(\boldsymbol{\theta}, \mathbf{x})$  volvemos a centrarnos en las propiedades de suavizado. Para el análisis de dos mallas y su suavizado tendremos de nuevo que distinguir entre componentes de alta y baja frecuencia en  $\mathbf{G}_h$ , con respecto a  $\mathbf{G}_{2h}$ . La idea de esta distinción vuelve a estar basada en el hecho de que solo las funciones de la forma

$$\varphi_h(\boldsymbol{\theta}, \cdot) \quad \text{con} \quad \frac{-\pi}{2} \leq \boldsymbol{\theta} < \frac{\pi}{2}$$

son distinguibles en  $G_{2h}$ . Más concretamente, para cada  $\bar{\boldsymbol{\theta}} \in [-\pi/2, \pi/2)^2$ , tenemos otras tres funciones de alta frecuencia,  $\varphi_h(\boldsymbol{\theta}, \cdot)$ , que coinciden en  $G_H$  con  $\varphi_h(\bar{\boldsymbol{\theta}}, \cdot)$  y, por lo tanto, no son distinguibles en  $G_{2h}$ . De hecho, tenemos

$$\varphi_h(\boldsymbol{\theta}, \mathbf{x}) = \varphi_h(\bar{\boldsymbol{\theta}}, \mathbf{x}) \quad \text{para cada } \mathbf{x} \in G_{2h} \text{ si y solo si } \boldsymbol{\theta} = \bar{\boldsymbol{\theta}} \text{ mód } \pi.$$

Definimos rigurosamente las funciones de alta y baja frecuencia.

**Definición 2.3.1** (*Funciones  $\varphi_h$  de alta y baja frecuencia*)

- Diremos que  $\varphi_h(\boldsymbol{\theta}, \cdot)$  es de baja frecuencia si y solo si  $\boldsymbol{\theta} \in T^{low} := [-\frac{\pi}{2}, \frac{\pi}{2})^2$ .
- Diremos que  $\varphi_h(\boldsymbol{\theta}, \cdot)$  es una componente de alta frecuencia si y solo si  $\boldsymbol{\theta} \in T^{high} := [-\pi, \pi)^2 \setminus [-\frac{\pi}{2}, \frac{\pi}{2})^2$ .

Con la introducción de estos conceptos desarrollaremos a continuación el análisis del suavizado.

### 2.3.2. Análisis del suavizado

El método de suavizado para  $\Delta_h u_h = f_h$  puede escribirse de la forma

$$L_h^+ \bar{w}_h + L_h^- w_h = f_h, \tag{2.5}$$

donde  $w_h$  es la aproximación anterior y  $\bar{w}_h$  es la nueva aproximación de  $u_h$ . Entonces, podemos caracterizar al método de suavizado por la siguiente escisión

$$\Delta_h = L_h^+ + L_h^-.$$

En el caso concreto de GS-LEX tendremos que la escisión viene dada por

$$L_h^+ = \frac{1}{h^2} \begin{bmatrix} 0 & & \\ 1 & -4 & 0 \\ & 1 & \end{bmatrix}_h, \quad L_h^- = \frac{1}{h^2} \begin{bmatrix} 1 & & \\ 0 & 0 & 1 \\ & 0 & \end{bmatrix}_h.$$

### 2.3. ANÁLISIS LOCAL DE FOURIER DE LAS PROPIEDADES DE SUAVIZADO 31

Hemos escrito los operadores con la notación de la molécula computacional para aligerar la notación. Una breve explicación de cómo manejar esta notación es que funciona de forma similar a una matriz de coeficientes de un sistema lineal, es decir, los coeficientes que aparecen en los corchetes son los coeficientes que tendrá la función (evaluada en un cierto punto) a la que se le aplica el operador (donde no hay coeficiente entonces es 0). Entonces nos preguntamos, ¿en qué punto se evalúa la función que está siendo multiplicada por uno de los coeficientes? La respuesta nos la da la posición que ocupa cada coeficiente en el corchete. En el punto central estará evaluada en el punto  $(x_i, y_j)$ , si nos movemos a la izquierda entonces en  $(x_i - h, y_j)$  y si nos movemos a la derecha en  $(x_i + h, y_j)$ . De la misma forma, si nos movemos hacia arriba entonces  $(x_i, y_j + h)$  y si nos movemos hacia abajo  $(x_i, y_j - h)$ . Con esto hacemos el resto de posiciones de dentro de los corchetes y, notemos que, el salto en los puntos de la malla nos lo da el subíndice en el corchete derecho.

También podemos hacer la escisión para  $\omega$ -JAC

$$L_h^+ = \frac{1}{h^2} \begin{bmatrix} 0 & & \\ 0 & -4/\omega & 0 \\ & 0 & \end{bmatrix}_h, \quad L_h^- = \frac{1}{h^2} \begin{bmatrix} & & 1 \\ 1 & -4(1 - 1/\omega) & 1 \\ & & 1 \end{bmatrix}_h.$$

Todo esto será válido para cualquier operador discreto que admita localmente la escisión (2.5).

Puesto que  $\Delta_h u_h = f_h$ , restando en (2.5) tendremos

$$L_h^+ \bar{v}_h + L_h^- v_h = 0,$$

donde  $\bar{v}_h = u_h - \bar{w}_h$ ,  $v_h = u_h - w_h$ . También tenemos, de forma equivalente, que  $\bar{v}_h = S_h v_h$ , donde  $S_h$  es el operador de suavizado resultante.

Aplicando a  $L_h^+$  y a  $L_h^-$  las autofunciones  $\varphi(\boldsymbol{\theta}, \mathbf{x})$ , obtenemos

$$\begin{aligned} L_h^+ e^{i\boldsymbol{\theta}\mathbf{x}/\mathbf{h}} &= \tilde{L}_h^+(\boldsymbol{\theta}) e^{i\boldsymbol{\theta}\mathbf{x}/\mathbf{h}} \\ L_h^- e^{i\boldsymbol{\theta}\mathbf{x}/\mathbf{h}} &= \tilde{L}_h^-(\boldsymbol{\theta}) e^{i\boldsymbol{\theta}\mathbf{x}/\mathbf{h}}, \end{aligned}$$

donde  $\tilde{L}_h^+(\boldsymbol{\theta})$  y  $\tilde{L}_h^-(\boldsymbol{\theta})$  denotan los autovalores formales de los operadores  $L_h^+$  y  $L_h^-$ , respectivamente. De aquí se deduce de manera inmediata el siguiente lema:

**Lema 2.3.3** *Todas las funciones  $\varphi_h(\boldsymbol{\theta}, \cdot)$  con  $\tilde{L}_h^+(\boldsymbol{\theta}) \neq 0$  son autofunciones de  $S_h$ :*

$$S_h \varphi(\boldsymbol{\theta}, \mathbf{x}) = \tilde{S}_h(\boldsymbol{\theta}) \varphi(\boldsymbol{\theta}, \mathbf{x}) \quad (-\pi \leq \boldsymbol{\theta} < \pi)$$

con  $\tilde{S}_h(\boldsymbol{\theta}) := \tilde{L}_h^-(\boldsymbol{\theta}) / \tilde{L}_h^+(\boldsymbol{\theta})$ .

Basándonos en este lema tenemos la siguiente definición:

**Definición 2.3.2** *El factor de suavizado ( $\mu_{loc}$ ) es la magnitud dada por*

$$\mu_{loc} = \mu_{loc}(S_h) := \sup\{|\tilde{S}_h(\boldsymbol{\theta})| : \boldsymbol{\theta} \in T^{high}\}.$$

**GS-LEX**,  $\mu_{loc}(S_h) = 0,5$

En el caso GS-LEX, tendremos que

$$\begin{aligned} L_h^+ e^{i\theta \mathbf{x}/h} &= \frac{1}{h^2} \begin{bmatrix} 0 & & \\ 1 & -4 & 0 \\ & 1 & \end{bmatrix}_h e^{i\theta \mathbf{x}/h} \\ &= \frac{1}{h^2} (-4 + e^{-i\theta_1} + e^{-i\theta_2}) e^{i\theta \mathbf{x}/h}, \\ L_h^- e^{i\theta \mathbf{x}/h} &= \frac{1}{h^2} \begin{bmatrix} & 1 & \\ 0 & 0 & 1 \\ & 0 & \end{bmatrix}_h e^{i\theta \mathbf{x}/h} \\ &= \frac{1}{h^2} (e^{-i\theta_1} + e^{-i\theta_2}) e^{i\theta \mathbf{x}/h}. \end{aligned}$$

Y, por lo tanto, aplicando el lema 2.3.3 tenemos que

$$\tilde{S}_h(\boldsymbol{\theta}) := -\frac{\tilde{L}_h^-(\boldsymbol{\theta})}{\tilde{L}_h^+(\boldsymbol{\theta})} = \frac{e^{-i\theta_1} + e^{-i\theta_2}}{4 - e^{-i\theta_1} - e^{-i\theta_2}}.$$

Por tanto, el factor de suavizado para GS-LEX, resulta ser

$$\mu_{loc}(S_h) = \sup \left\{ \left| \frac{e^{-i\theta_1} + e^{-i\theta_2}}{4 - e^{-i\theta_1} - e^{-i\theta_2}} \right| : \boldsymbol{\theta} \in T^{high} \right\}.$$

Para buscar dicho superior, notemos que buscar el superior de la función

$$g(\boldsymbol{\theta}) = \left| \frac{e^{-i\theta_1} + e^{-i\theta_2}}{4 - e^{-i\theta_1} - e^{-i\theta_2}} \right|$$

es equivalente a buscarlo de la función  $f(\boldsymbol{\theta}) = g(\boldsymbol{\theta})^2$ , que es más sencilla de manejar. Además haciendo algunas operaciones tenemos

$$\begin{aligned} f(\boldsymbol{\theta}) := g(\boldsymbol{\theta})^2 &= \left| \frac{(\cos(\theta_1) + \cos(\theta_2)) - i(\sin(\theta_1) + \sin(\theta_2))}{(\cos(\theta_1) + \cos(\theta_2) - 4) - i(\sin(\theta_1) + \sin(\theta_2))} \right|^2 \\ &= \frac{(\cos(\theta_1) + \cos(\theta_2))^2 + (\sin(\theta_1) + \sin(\theta_2))^2}{(\cos(\theta_1) + \cos(\theta_2) - 4)^2 + (\sin(\theta_1) + \sin(\theta_2))^2} \\ &= \frac{1 + \cos(\theta_1 - \theta_2)}{9 + \cos(\theta_1 - \theta_2) - 4(\cos(\theta_1) + \cos(\theta_2))}. \end{aligned}$$

Empezamos estudiando los puntos críticos en el interior de  $T^{high}$ . Para ello, tendremos que ver los puntos en  $T^{high}$  donde se anule  $\nabla f$ , que toma la forma

$$\begin{pmatrix} \frac{-\sin(\theta_1 - \theta_2)(9 + \cos(\theta_1 - \theta_2) - 4(\cos(\theta_1) + \cos(\theta_2))) - (1 + \cos(\theta_1 - \theta_2))(-\sin(\theta_1 - \theta_2) + 4\sin(\theta_1))}{(9 + \cos(\theta_1 - \theta_2) - 4(\cos(\theta_1) + \cos(\theta_2)))^2} \\ \frac{\sin(\theta_1 - \theta_2)(9 + \cos(\theta_1 - \theta_2) - 4(\cos(\theta_1) + \cos(\theta_2))) - (1 + \cos(\theta_1 - \theta_2))(\sin(\theta_1 - \theta_2) + 4\sin(\theta_2))}{(9 + \cos(\theta_1 - \theta_2) - 4(\cos(\theta_1) + \cos(\theta_2)))^2} \end{pmatrix}.$$

### 2.3. ANÁLISIS LOCAL DE FOURIER DE LAS PROPIEDADES DE SUAVIZADO33

Anulamos ambas componentes y sumándolas, se obtiene que debe cumplirse

$$(\cos(\theta_2 - \theta_1) + 1)(\sin(\theta_1) + \sin(\theta_2)) = 0.$$

Por tanto, o bien  $\cos(\theta_2 - \theta_1) = -1$  o bien  $\sin(\theta_1) = -\sin(\theta_2)$ . En el primer caso,  $f(\theta) = 0$ , con lo cual es un mínimo de la función y no un máximo. Si no ocurre eso, pero ocurre lo segundo, para que  $(\theta_1, \theta_2) \in T^{high}$  debe ocurrir  $\theta_1 = -\theta_2$  con  $|\theta_1| \geq \pi/2$ . Estudiando la función  $\bar{f}(\theta_1) = f(\theta_1, -\theta_1)$  en el intervalo  $[-\pi/2, \pi/2]$ , puede observarse que es positiva y el máximo se toma en  $\theta_1 = \pi$  y  $\bar{f}(\pi) = 1/5$ . Se ha obtenido en cualquier caso, un punto de la frontera de  $T^{high}$  y no del interior, con lo que concluimos que no hay puntos críticos en el interior de  $T^{high}$ .

Estudiemos ahora el supremo en la frontera de  $T^{high}$ . Se puede ver, haciendo la gráfica de cada una de las funciones que obtendremos a continuación, que ninguna de las siguientes funciones alcanza su máximo en los extremos del intervalo.

- $\theta_1 = \pi/2, \quad -\pi/2 \leq \theta_2 \leq \pi/2:$

$$f(\pi/2, \theta_2) = \frac{1 + \cos(\pi/2 - \theta_2)}{9 + \cos(\pi/2 - \theta_2) - 4 \cos(\theta_2)} = \frac{1 + \sin(\theta_2)}{9 + \sin(\theta_2) - 4 \cos(\theta_2)}.$$

Se trata de una función en una variable. Estudiamos sus puntos críticos:

$$0 = \frac{\partial}{\partial \theta_2} f(\pi/2, \theta_2) \iff 0 = 2 \cos(\theta_2) - \sin(\theta_2) - 1.$$

Despejando  $\cos(\theta_2)$ , elevando al cuadrado y utilizando  $\sin^2(\theta_2) + \cos^2(\theta_2) = 1$ , se obtiene la ecuación  $5 \sin^2(\theta_2) + 2 \sin(\theta_2) - 3 = 0$ , con lo que tendríamos dos puntos críticos:  $\theta_2 = \arccos(4/5), -\pi/2$ . Evaluando la función en cada uno de los puntos críticos y en el otro extremo tendremos que:

$$f(\pi/2, -\pi/2) = 0, \quad f(\pi/2, \arccos(4/5)) = 1/4, \quad f(\pi/2, \pi/2) = 1/5.$$

- $\theta_2 = -\pi/2, \quad -\pi/2 \leq \theta_1 \leq \pi/2:$

$$f(\theta_1, -\pi/2) = \frac{1 - \sin(\theta_1)}{9 - \sin(\theta_1) - 4 \cos(\theta_1)}.$$

De nuevo estudiamos los puntos críticos de esta función de una variable:

$$0 = \frac{\partial}{\partial \theta_1} f(\theta_1, -\pi/2) \iff 0 = 1 - 2 \cos(\theta_1) - \sin(\theta_1).$$

Las soluciones de esta ecuación son  $\theta_1 = -\arccos(4/5), \pi/2$  y, evaluando, en estos puntos críticos y en el otro extremo tendremos

$$f(\pi/2, -\pi/2) = 0, \quad f(\arccos(4/5), -\pi/2) = 1/4, \quad f(-\pi/2, -\pi/2) = 1/5$$

- $\theta_2 = \pi/2, \quad -\pi/2 \leq \theta_1 \leq \pi/2:$

$$f(\theta_1, \pi/2) = \frac{1 + \sin(\theta_1)}{9 + \sin(\theta_1) - 4 \cos(\theta_1)}.$$

Esto es la misma ecuación que en  $\theta_1 = \pi/2$  y, por lo tanto, el máximo se alcanza en  $\theta_1 = \arccos(4/5)$ .

- $\theta_1 = -\pi/2, \quad -\pi/2 \leq \theta_2 \leq \pi/2:$

$$f(-\pi/2, \theta_2) = \frac{1 - \sin(\theta_2)}{9 - \sin(\theta_2) - 4 \cos(\theta_2)}.$$

Misma ecuación que en el caso  $\theta_2 = -\pi/2$ , entonces el máximo está en  $\theta_2 = -\arccos(4/5)$ .

- $\theta_1 = -\pi, \quad -\pi \leq \theta_2 \leq \pi:$

$$f(-\pi, \theta_2) = \frac{1 - \cos(\theta_2)}{13 - 5 \cos(\theta_2)}.$$

Derivamos y tenemos

$$0 = \frac{\partial}{\partial \theta_2} f(-\pi, \theta_2) \iff \sin(\theta_2) = 0.$$

Concluimos que los puntos críticos son  $\theta_2 = 0, -\pi$ , y evaluando tendremos que

$$f(-\pi, 0) = 0, \quad f(-\pi, -\pi) = 1/9.$$

- $\theta_2 = -\pi, \quad -\pi \leq \theta_1 \leq \pi:$

$$f(\theta_1, -\pi) = \frac{1 - \cos(\theta_1)}{13 - 5 \cos(\theta_1)}.$$

Tenemos las mismas ecuaciones de antes. Por lo tanto, tendremos los mismos resultados solo que ahora para  $\theta_1$ .

- $\theta_1 = \pi, \quad -\pi \leq \theta_2 \leq \pi:$

$$f(\pi, \theta_2) = \frac{1 - \cos(\theta_2)}{13 - 5 \cos(\theta_2)}.$$

Mismos resultados que en los casos anteriores.

- $\theta_2 = \pi, \quad -\pi \leq \theta_1 \leq \pi:$

$$f(\theta_1, \pi) = \frac{1 - \cos(\theta_1)}{13 - 5 \cos(\theta_1)}.$$

Mismos resultados que en los casos anteriores.

El Teorema de Weirstrass nos aseguraba que la función  $f$  alcanza su máximo en  $\overline{T^{high}}$  pero, de hecho, hemos visto que lo alcanza en  $T^{high}$ . Comparando las evaluaciones de  $f$  en los puntos críticos que hemos obtenido en la frontera tenemos que el máximo se alcanza en  $(\theta_1, \theta_2) = (\pi/2, \arccos(4/5))$  donde  $f(\pi/2, \arccos(4/5)) = 1/4$ . Con lo que se demuestra que  $\mu_{loc} = 0,5$ , para GS-LEX con el engrosamiento estándar.

### 2.3. ANÁLISIS LOCAL DE FOURIER DE LAS PROPIEDADES DE SUAVIZADO 35

#### S.O.R. o $\omega$ GS-LEX

El método  $\omega$  GS-LEX viene dado por

$$\begin{aligned} z_h^{m+1}(x_i, y_j) &= \frac{1}{4}[-h^2 f_h(x_i, y_j) + u_h^{m+1}(x_i - h, y_j) + u_h^m(x_i + h, y_j) \\ &\quad + u_h^{m+1}(x_i, y_j - h) + u_h^m(x_i, y_j + h)] \\ u_h^{m+1} &= u_h^m + \omega(z_h^{m+1} - u_h^m), \end{aligned}$$

de donde podemos deducir la escisión  $L_h^-, L_h^+$  teniendo en cuenta que, despejando  $z_h^{n+1}$  en la segunda ecuación, sustituyendo en la primera, y despejando  $f_h$ , se tiene

$$\begin{aligned} f_h(x_i, y_j) &= \frac{4}{\omega h^2} \left[ -u_h^{m+1}(x_i, y_j) - \frac{\omega}{4}(u_h^{m+1}(x_i - h, y_j) + u_h^{m+1}(x_i, y_j - h) \right. \\ &\quad \left. + u_h^m(x_i + h, y_j) + u_h^m(x_i, y_j + h)) - (\omega - 1)u_h^m(x_i, y_j) \right], \end{aligned}$$

con lo que finalmente las matrices  $L_h^+, L_h^-$  son

$$L_h^+ = \frac{1}{h^2} \begin{bmatrix} 0 & & \\ 1 & -\frac{4}{\omega} & 0 \\ & 1 & \end{bmatrix}_h, \quad L_h^- = \frac{1}{h^2} \begin{bmatrix} & 1 & \\ 0 & -\frac{4(\omega-1)}{\omega} & 1 \\ & 0 & \end{bmatrix}_h.$$

Para obtener el factor de suavizado notaremos entonces que

$$\begin{aligned} L_h^+ e^{i\theta \mathbf{x}/h} &= \frac{1}{h^2} \left( -\frac{4}{\omega} + e^{-i\theta_1} + e^{-i\theta_2} \right) e^{i\theta \mathbf{x}/h}, \\ L_h^- e^{i\theta \mathbf{x}/h} &= \frac{1}{h^2} \left( -\frac{4(\omega-1)}{\omega} + e^{i\theta_1} + e^{i\theta_2} \right) e^{i\theta \mathbf{x}/h}, \end{aligned}$$

y por lo tanto,

$$\tilde{S}_h(\boldsymbol{\theta}) := -\frac{\tilde{L}_h^-(\boldsymbol{\theta})}{\tilde{L}_h^+(\boldsymbol{\theta})} = \frac{1 - \omega + \frac{\omega e^{i\theta_1}}{4} + \frac{\omega e^{i\theta_2}}{4}}{\frac{\omega e^{-i\theta_1}}{4} + \frac{\omega e^{-i\theta_2}}{4} - 1}.$$

Estudiamos el superior de  $|\tilde{S}_h(\boldsymbol{\theta})|$  en  $T^{high}$  y para ello, como en el caso de GS-LEX, estudiamos la función  $|\tilde{S}_h(\boldsymbol{\theta})|^2$ . Escribiendo las exponenciales complejas en función de senos y cosenos tendremos

$$|\tilde{S}_h(\boldsymbol{\theta})|^2 = \frac{\left( (\omega - 1) - \frac{\omega}{4}(\cos(\theta_1) + \cos(\theta_2)) \right)^2 + \left( \frac{\omega}{4}(\sin(\theta_1) + \sin(\theta_2)) \right)^2}{\left( 1 - \frac{\omega}{4}(\cos(\theta_1) + \cos(\theta_2)) \right)^2 + \left( \frac{\omega}{4}(\sin(\theta_1) + \sin(\theta_2)) \right)^2}.$$

Desarrollando los cuadrados y haciendo las simplificaciones oportunas, se llega a la expresión

$$|\tilde{S}_h(\boldsymbol{\theta})|^2 = \frac{(\omega - 1)^2 + \frac{\omega^2}{8}(1 + \cos(\theta_1 - \theta_2)) - (\omega - 1)\frac{\omega}{2}(\cos(\theta_1) + \cos(\theta_2))}{1 - \frac{\omega}{2}(\cos(\theta_1) + \cos(\theta_2)) + \frac{\omega^2}{8}(1 + \cos(\theta_1 - \theta_2))}.$$

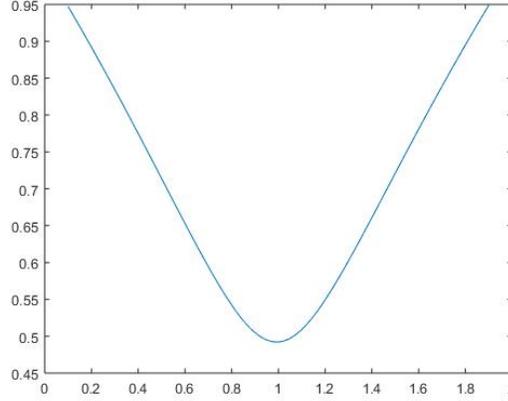


Figura 2.4:  $|\tilde{S}_h(\theta)|$  frente a  $\omega$  para  $0 < \omega < 2$ .

Buscar analíticamente el valor superior de la función en  $T^{high}$  es más complicado que en el caso GSLEX. Por lo tanto, hemos realizado con MATLAB una representación de la función para valores  $\omega$  equiespaciados 0,01 en el intervalo  $(0, 2)$  que se muestra en el figura 2.4. En esta representación gráfica vemos que el  $\omega_{opt}$ , que más pequeño hace el factor de suavizado, está muy cerca de 1. Los beneficios que proporciona sobre las propiedades de suavizado el  $\omega_{opt}$  son muy pequeñas en comparación con el trabajo adicional que supone encontrarlo y utilizarlo en el proceso de suavizado (2 operaciones más por cada punto de la malla).

## 2.4. Introducción al ciclo de 2 mallas

Como ya dijimos al principio del capítulo, el método multigrid se basa en dos principales ideas: el suavizado y la corrección en malla grosera. En esta sección vamos a presentar estas ideas para dos mallas con la intención de generalizarlo, más adelante, en múltiples mallas.

### 2.4.1. Aproximación de la solución de la ecuación del defecto

Para el problema discreto dado por  $L_h u_h = f_h$ , si  $u_h^m$  es una aproximación de  $u_h$ , denotaremos en esta sección el error por

$$v_h^m := u_h - u_h^m.$$

Ahora, denotaremos por el defecto o residuo a

$$d_h^m := f_h - L_h u_h^m. \quad (2.6)$$

Obsérvese que el error es la diferencia entre la solución y la aproximación de la solución de nuestro problema discreto. Sin embargo, el defecto es la diferencia entre la imagen de la solución y de la aproximación de la solución de nuestro problema discreto a través del operador  $L_h$ . Es claro por (2.6) que se tiene la ecuación del defecto

$$L_h v_h^m = d_h^m, \quad (2.7)$$

y se puede considerar el siguiente proceso

$$u_h^m \longrightarrow d_h^m = f_h - L_h u_h^m \longrightarrow L_h v_h^m = d_h^m \longrightarrow u_h = u_h^m + v_h^m.$$

En sí, este proceso no supone ninguna mejora o diferencia a lo que ya teníamos. Sin embargo, si en lugar de considerar  $L_h$ , tomamos una aproximación más simple  $\hat{L}_h$  de forma que exista su inverso  $\hat{L}_h^{-1}$ , la solución de

$$\hat{L}_h \hat{v}_h^m = d_h^m,$$

nos dará una nueva aproximación de nuestro problema discreto dada por

$$u_h^{m+1} := u_h^m + \hat{v}_h^m.$$

Por tanto, el proceso anterior se convierte en lo siguiente

$$u_h^m \longrightarrow d_h^m = f_h - L_h u_h^m \longrightarrow \hat{L}_h \hat{v}_h^m = d_h^m \longrightarrow u_h^{m+1} = u_h^m + \hat{v}_h^m.$$

Entonces, empezando con un  $u_h^0$  y aplicando el proceso sucesivas veces, tendremos un proceso iterativo. El operador de iteración del proceso vendrá dado por

$$M_h = I_h - C_h L_h : \mathcal{G}(\Omega_h) \longrightarrow \mathcal{G}(\Omega_h),$$

donde  $C_h := (\hat{L}_h)^{-1}$  e  $I_h$  denota el operador identidad en  $\mathcal{G}(\Omega_h)$ , que es el espacio de funciones definidas en la malla  $\Omega_h$ . Es decir, el proceso iterativo es de la forma

$$u_h^{m+1} = M_h u_h^m + s_h, \quad \text{donde } s_h = (\hat{L}_h)^{-1} f_h \quad (m = 0, 1, 2, \dots).$$

Para el error tendremos que

$$v_h^{m+1} = M_h v_h^m = (I_h - C_h L_h) v_h^m \quad (m = 0, 1, 2, \dots),$$

y para el defecto se tiene que

$$d_h^{m+1} = L_h M_h L_h^{-1} d_h^m = (I_h - L_h C_h) d_h^m \quad (m = 0, 1, 2, \dots).$$

Comenzando con el iterante inicial  $u^0 = 0$ , tendremos

$$\begin{aligned} u_h^m &= (I_h + M_h + M_h^2 + \dots + M_h^{m-1})(\hat{L}_h)^{-1} f_h \\ &= (I_h - M_h^m)(I_h - M_h)^{-1}(\hat{L}_h)^{-1} f_h \\ &= (I_h - M_h^m)L_h^{-1} f_h. \end{aligned} \quad (2.8)$$

Las propiedades de convergencia del proceso arriba descrito están caracterizadas por el radio espectral del operador de iteración, dado por

$$\rho(M_h) = \rho(I_h - C_h L_h) = \rho(I_h - L_h C_h).$$

Para cualquier norma matricial asociada a los valores en la red  $\Omega_h$  entonces, los valores  $\|I_h - C_h L_h\|$  y  $\|I_h - L_h C_h\|$  son cotas superiores para el factor de reducción del error y para el factor de reducción del defecto, respectivamente, por cada iteración.

**Observación 2.4.1** *Aplicar un método clasico iterativo, como Gauss-Seidel o Jacobi a la ecuación  $L_h u_h = f_h$ , se puede interpretar como aproximar la solución de la ecuación del defecto (2.7) en el siguiente sentido:*

- Si consideramos como aproximación  $\hat{v}_h^m$  la obtenida tras una iteración del método Gauss-Seidel (con iterante anterior nula) a la ecuación del defecto, esto queda

$$\begin{aligned} (\hat{v}_h^m)_i &= \frac{(d_h^m)_i - \sum_{j<i} a_{i,j} (\hat{v}_h^m)_j}{a_{i,i}} \\ &= \frac{(f_h)_i - (L_h u_h^m)_i - \sum_{j<i} a_{i,j} (u_h^{m+1} - u_h^m)_j}{a_{i,i}}, \end{aligned}$$

donde  $a_{i,j}$  son los coeficientes de la matriz del operador. Por lo tanto, aplicando que  $u_h^{m+1} = u_h^m + \hat{v}_h^m$  en el lado izquierdo, tendremos

$$(u_h^{m+1})_i = \frac{(f_h)_i - \sum_{j<i} a_{i,j} (u_h^{m+1})_j - \sum_{j>i} a_{i,j} (u_h^m)_j}{a_{i,i}}.$$

Esto es el método iterativo Gauss-Seidel aplicado a  $L_h u_h = f_h$ .

- Veamos ahora que tomando  $\hat{L}_h = \frac{1}{\omega} D_h$ , donde  $D_h$  es la diagonal de la matriz del operador  $L_h$ , obtendremos en el proceso el método iterativo de  $\omega$ -JAC. Es obvio que

$$(\hat{v}_h^m)_i = \omega \frac{(d_h^m)_i}{d_i},$$

y por lo tanto,

$$(u_h^{m+1})_i = (u_h^m)_i + \frac{\omega}{d_i} ((f_h)_i - (L_h u_h^m)_i).$$

Esto último es exactamente el método iterativo de  $\omega$ -JAC aplicado al problema  $L_h u_h = f_h$ .

## 2.4.2. Corrección en malla grosera

La idea es resolver la ecuación del defecto (2.7) en otra malla más gruesa  $\Omega_H$ . Esto es, la ecuación del defecto es remplazada por

$$L_H \hat{v}_H^m = d_H^m.$$

Supondremos que  $L_H : \mathcal{G}(\Omega_H) \rightarrow \mathcal{G}(\Omega_h)$ , con  $\dim(\mathcal{G}(\Omega_H)) < \dim(\mathcal{G}(\Omega_h))$ , y que existe  $(L_H)^{-1}$ . Puesto que  $d_H^m$  y  $\hat{v}_H^m$  son funciones en la red grosera  $\Omega_H$ , definimos dos operadores que nos lleven desde la malla fina  $\Omega_h$  a la grosera  $\Omega_H$  y viceversa. Estos son

$$I_h^H : \mathcal{G}(\Omega_h) \rightarrow \mathcal{G}(\Omega_H), \quad I_H^h : \mathcal{G}(\Omega_H) \rightarrow \mathcal{G}(\Omega_h).$$

Usaremos como  $I_h^H$  la restricción de una función en la malla  $\Omega_h$  a la malla  $\Omega_H$  y, por lo tanto,  $d_H^m$  vendrá dado por la expresión

$$d_H^m := I_h^H d_h^m,$$

y usaremos  $I_H^h$  un interpolador, es decir, un operador que a partir de los valores de una función en  $\Omega_H$  la extienda a  $\Omega_h$ , y por lo tanto tendremos

$$\hat{v}_h^m := I_H^h \hat{v}_H^m.$$

En resumen, llamaremos corrección de malla grosera al siguiente proceso:

- Partimos de la ecuación del defecto  $d_h^m = f_h - L_h u_h^m$ .
- Restringimos el defecto a la malla grosera  $d_H^m = I_h^H d_h^m$ .
- Resolvemos en la malla grosera el problema  $L_H \hat{v}_H^m = d_H^m$ .
- Interpolamos para volver a la malla fina  $\hat{v}_h^m = I_H^h \hat{v}_H^m$ .
- Obtenemos la nueva aproximación  $u_h^{m+1} = u_h^m + \hat{v}_h^m$ .

El operador asociado viene dado por  $I_h - C_h L_h$ , donde  $C_h = I_H^h L_H^{-1} I_h^H$ . Pero debemos destacar que el proceso como tal no es convergente ya que tenemos que

$$\rho(I_h - I_H^h L_H^{-1} I_h^H L_h) \geq 1 \quad (2.9)$$

Esto ocurre debido a que el operador  $I_h^H$  parte de un espacio de dimensión mayor a uno de dimensión menor y por lo tanto el operador  $C_h = I_H^h L_H^{-1} I_h^H$  es no invertible. En particular, esto implica que  $C_h L_h v_h = 0$ , para cierto  $v_h \neq 0$ , y por tanto 1 es siempre autovalor de la matriz en (2.9).

**Ejemplo 2.4.1** Para el operador inyección  $I_h^H$ , cualquier error  $v_h \in \mathcal{G}(\Omega_h)$  con

$$L_h v_h(P) = \begin{cases} 0 & \text{para } P \in \Omega_H \\ \text{arbitrario} & \text{para } P \notin \Omega_H \end{cases}$$

se convierte en 0 tras aplicarle  $I_h^H$  y, por lo tanto, al aplicarle  $C_h$ . Esto es, el iterante no cambia tras ser aplicada la corrección en malla grosera.

Más concretamente, consideramos nuestro problema modelo con  $h = 1/N$  y

$$v_h(x, y) = \sin\left(\frac{N}{2}\pi x\right) \sin\left(\frac{N}{2}\pi y\right).$$

Esta función es de alta frecuencia y, con el engrosamiento estándar, tendremos

$$v_h(P) = L_h v_h(P) = I_h^{2h} L_h v_h(P) = 0 \quad \forall P \in \Omega_{2h}.$$

Con este ejemplo pretendemos mostrar que es una función de alta frecuencia la que se va a 0. La forma de hacer el método convergente será suavizar, como veremos a continuación.

### 2.4.3. Estructura del operador de dos mallas

Hasta este punto lo que tenemos es que el ciclo de dos mallas se basa en la aproximación del error utilizando un operador más simple que  $(L_h)^{-1}$ . Dicha aproximación será  $I_H^h L_H^{-1} I_h^H$ . Sin embargo, esto no nos da un método convergente y para arreglarlo se introduce un pre y un post suavizado.

A continuación, mostraremos cómo quedaría el proceso de un ciclo de dos mallas.

**Ciclo de dos mallas**  $u_h^{m+1} = TGCYC(u_h^m, L_h, f_h, \nu_1, \nu_2)$

1. Presuavizado: Calculamos  $\bar{u}_h^m$  aplicándole  $\nu_1$  iteraciones de un método de suavizado a  $u_h^m$ :

$$\bar{u}_h^m = SMOOTH^{\nu_1}(u_h^m, L_h, f_h).$$

2. Corrección en malla grosera

- Partimos de la ecuación del defecto  $\bar{d}_h^m = f_h - L_h \bar{u}_h^m$ .
- Restringimos el defecto a la malla grosera  $\bar{d}_H^m = I_H^H \bar{d}_h^m$ .
- Resolvemos en la malla grosera el problema  $L_H \hat{v}_H^m = \bar{d}_H^m$ .
- Interpolamos para volver a la malla fina  $\hat{v}_h^m = I_H^h \hat{v}_H^m$ .
- Obtenemos la nueva aproximación  $u_h^{m,C} = \bar{u}_h^m + \hat{v}_h^m$ .

3. Postsuavizado: Calculamos  $u_h^{m+1}$  aplicándole  $\nu_2$  iteraciones del método de suavizado dado a  $u_h^{m,C}$ :

$$u_h^{m+1} = SMOOTH^{\nu_2}(u_h^{m,C}, L_h, f_h).$$

Observamos que el operador que obtenemos ahora es

$$M_h^H = S_h^{\nu_2} K_h^H S_h^{\nu_1} \text{ donde } K_h^H := I_h - I_H^h L_H^{-1} I_h^H L_h.$$

La experiencia con los métodos multigrid muestra que la elección de cada componente del método tiene una gran influencia en la eficiencia del método. Sin embargo, no hay reglas que nos indiquen cuáles son las mejores elecciones de los componentes para construir algoritmos óptimos en las distintas aplicaciones del método. Una forma de elegir una buena componente del método puede ser estudiando las propiedades de cada componente en cada problema concreto, de forma similar a como hicimos con las propiedades de suavizado de los métodos clásicos para nuestro problema modelo.

## 2.5. Componentes del método multigrid

A continuación daremos algunos ejemplos de componentes del método multigrid que serán utilizados en el desarrollo de este trabajo.

### 2.5.1. Elección del operador en la red grosera $L_H$

Consideramos el engrosamiento estándar de  $\Omega_h$  a  $\Omega_{2h}$ . La elección natural que haremos será la del análogo a  $L_h$  en la red  $\Omega_H$ . Es decir, para el problema modelo tendremos que

$$L_H = \frac{1}{H^2} \begin{bmatrix} 1 & 1 & \\ 1 & -4 & 1 \\ & 1 & \end{bmatrix}_H.$$

### 2.5.2. Elección del operador restricción

El operador restricción  $I_h^{2h}$  más sencillo sería la inyección. Sin embargo, en este trabajo el que más utilizaremos será el operador *full weighting* (FW) que en notación de molécula computacional viene dado por

$$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}_h^{2h},$$

y aplicado a la función  $d_h$ , tendremos que, para  $(x, y) \in \Omega_{2h}$ ,

$$\begin{aligned} d_{2h}(x, y) &= I_h^{2h} d_h(x, y) \\ &= \frac{1}{16} [4d_h(x, y) + 2d_h(x + h, y) + 2d_h(x - h, y) + 2d_h(x, y + h) \\ &\quad + 2d_h(x, y - h) + d_h(x + h, y + h) + d_h(x + h, y - h) \\ &\quad + d_h(x - h, y + h) + d_h(x - h, y - h)]. \end{aligned}$$

**Observación 2.5.1** *El operador FW se puede obtener a partir de la versión discreta de la condición*

$$\int_{\Omega_{x,y}} \omega(\bar{x}, \bar{y}) d\bar{x}d\bar{y} = \int_{\Omega_{x,y}} (I_h^{2h}\omega)(\bar{x}, \bar{y}) d\bar{x}d\bar{y},$$

donde  $\Omega_{x,y} = [x - h, x + h] \times [y - h, y + h]$ . Indicaremos brevemente cómo verlo: Aplicando la regla de cuadratura del punto medio a la derecha, tendremos

$$\int_{\Omega_{x,y}} (I_h^{2h}\omega)(\bar{x}, \bar{y}) d\bar{x}d\bar{y} \simeq (2h)^2 (I_h^{2h}\omega)(x, y).$$

En el lado izquierdo, escribiremos la integral de la forma

$$\begin{aligned} \int_{\Omega_{x,y}} \omega(\bar{x}, \bar{y}) d\bar{x}d\bar{y} &= \int_y^{y+h} \int_{x-h}^x \omega(\bar{x}, \bar{y}) d\bar{x}d\bar{y} + \int_y^{y+h} \int_x^{x+h} \omega(\bar{x}, \bar{y}) d\bar{x}d\bar{y} \\ &\quad + \int_{y-h}^y \int_x^{x+h} \omega(\bar{x}, \bar{y}) d\bar{x}d\bar{y} + \int_{y-h}^y \int_{x-h}^x \omega(\bar{x}, \bar{y}) d\bar{x}d\bar{y}. \end{aligned}$$

A cada una de estas integrales le aplicamos la fórmula de cuadratura del trapecio. Igualando ambas expresiones llegamos al operador FW.

### 2.5.3. Elección del operador interpolador

Consideraremos a lo largo del trabajo el interpolador bilineal, que viene dado por

$$I_{2h}^h \hat{v}_{2h}(x, y) = \begin{cases} \hat{v}_{2h}(x, y), & \text{para } \frac{x}{h}, \frac{y}{h} \text{ par,} \\ \frac{1}{2}[\hat{v}_{2h}(x, y+h) + \hat{v}_{2h}(x, y-h)], & \text{para } \frac{x}{h} \text{ par e } \frac{y}{h} \text{ impar,} \\ \frac{1}{2}[\hat{v}_{2h}(x+h, y) + \hat{v}_{2h}(x-h, y)], & \text{para } \frac{x}{h} \text{ impar e } \frac{y}{h} \text{ par,} \\ \frac{1}{4}[\hat{v}_{2h}(x+h, y+h) + \hat{v}_{2h}(x+h, y-h) \\ + \hat{v}_{2h}(x-h, y+h) + \hat{v}_{2h}(x-h, y-h)], & \text{para } \frac{x}{h}, \frac{y}{h} \text{ impar,} \end{cases}$$

que en notación de la molécula computacional lo escribimos de la forma

$$I_{2h}^h = \frac{1}{4} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}_{2h}^h.$$

## 2.6. Análisis local de Fourier del ciclo de 2 mallas

En esta sección aplicaremos el análisis local de Fourier al operador en dos mallas  $M_h^H$  definido anteriormente, donde estaremos considerando el engrosamiento estándar  $H = 2h$ . Para el cálculo de su factor de convergencia analizaremos cómo actúan los operadores  $L_h, I_h^H, L_H, I_H^h$  y  $S_h$  sobre las funciones  $\varphi_h(\boldsymbol{\theta}, \cdot)$ .

Para toda baja frecuencia  $\boldsymbol{\theta} = (\theta_1, \theta_2) \in T^{low} = [-\pi/2, \pi/2]^2$ , consideramos las frecuencias

$$\begin{aligned} \boldsymbol{\theta}^{(0,0)} &:= (\theta_1, \theta_2), & \boldsymbol{\theta}^{(1,1)} &:= (\bar{\theta}_1, \bar{\theta}_2), \\ \boldsymbol{\theta}^{(1,0)} &:= (\bar{\theta}_1, \theta_2), & \boldsymbol{\theta}^{(0,1)} &:= (\theta_1, \bar{\theta}_2), \end{aligned}$$

donde

$$\bar{\theta}_i := \begin{cases} \theta_i + \pi & \text{si } \theta_i < 0, \\ \theta_i - \pi & \text{si } \theta_i \geq 0, \end{cases}.$$

Con esto damos el siguiente lema

**Lema 2.6.1** 1. Para toda frecuencia  $\boldsymbol{\theta}^{(0,0)} \in T^{low}$  tenemos que, para todo  $\mathbf{x} \in \Omega_{2h}$ , se cumple:

$$\varphi_h(\boldsymbol{\theta}^{(0,0)}, \mathbf{x}) \equiv \varphi_h(\boldsymbol{\theta}^{(1,1)}, \mathbf{x}) \equiv \varphi_h(\boldsymbol{\theta}^{(1,0)}, \mathbf{x}) \equiv \varphi_h(\boldsymbol{\theta}^{(0,1)}, \mathbf{x}).$$

2. Cada una de estas cuatro funciones  $\varphi_h(\boldsymbol{\theta}^\alpha, \cdot)$  con  $\boldsymbol{\alpha} \in \{(0,0), (1,0), (0,1), (1,1)\}$  coinciden en  $\mathbf{G}_{2h}$  (malla de infinitos puntos) con la función  $\varphi_{2h}(2\boldsymbol{\theta}^{(0,0)}, \cdot)$ , esto es

$$\varphi_h(\boldsymbol{\theta}^\alpha, \mathbf{x}) \equiv \varphi_{2h}(2\boldsymbol{\theta}^{(0,0)}, \mathbf{x}), \quad \text{para todo } \mathbf{x} \in \mathbf{G}_{2h}$$

La demostración de la primera parte del lema se obtiene de forma directa sin más que sustituir por cada valor de los ángulos  $\theta^\alpha$ . Para la segunda parte, veamos la igualdad para  $\theta^{(1,1)}$  y el resto será de forma análoga.

$$\begin{aligned}\varphi_h(\theta^{(1,1)}, \mathbf{x}) &= e^{i(\theta_1 \pm \pi)x_1/h} e^{i(\theta_2 \pm \pi)x_2/h} \\ &= e^{i\theta_1 x_1/h} e^{\pm i\pi 2k_1} e^{i\theta_2 x_2/h} e^{\pm i2k_2\pi} \quad (x_i = 2k_i h, k_i \in \mathbb{Z}) \\ &= e^{i\theta_1 2x_1/(2h)} e^{i\theta_2 2x_2/(2h)} = \varphi_{2h}(2\theta^{(0,0)}, \mathbf{x}).\end{aligned}$$

Lo principal de este resultado es el hecho, que vamos a utilizar a lo largo de este análisis, de que las cuádruplas  $\varphi_h(\theta^\alpha, \cdot)$  coinciden en  $G_h$ . Todo esto nos lleva a la siguiente definición:

**Definición 2.6.1** *Llamaremos a las cuatro funciones  $\varphi_h(\theta^\alpha, \cdot)$  armónicos para el engrosamiento estándar. Para un  $\theta = \theta^{(0,0)} \in T^{\text{low}}$  definimos el espacio vectorial de dimensión 4 generado por los cuatro armónicos, es decir*

$$E_h^\theta := \langle \varphi_h(\theta^\alpha, \cdot) : \alpha \in \{(0,0), (1,0), (0,1), (1,1)\} \rangle.$$

Tomemos  $\psi \in E_h^\theta$  que lo representaremos de la forma

$$\psi = A^{(0,0)}\varphi_h(\theta^{(0,0)}, \cdot) + A^{(1,1)}\varphi_h(\theta^{(1,1)}, \cdot) + A^{(1,0)}\varphi_h(\theta^{(1,0)}, \cdot) + A^{(0,1)}\varphi_h(\theta^{(0,1)}, \cdot).$$

Vamos a analizar cómo son transformados los coeficientes  $A^\alpha$  tras aplicar a  $\psi$  el operador de dos redes  $M_h^H = S_h^{\nu_2} K_h^H S_h^{\nu_1}$  para nuestro problema modelo con engrosamiento estándar. Para ello daremos el siguiente teorema, donde  $I_h^{2h}$  será el operador FW,  $I_{2h}^h$  será el operador interpolador bilineal y  $L_{2h}$  será la elección natural vista en la sección 2.5.1, para el cual supondremos la existencia de  $(L_{2h})^{-1}$  que será discutido después. La demostración del teorema se irá dando de forma constructiva con las explicaciones posteriores a su enunciado.

**Teorema 2.6.1** *Bajo las hipótesis antes mencionadas tendremos*

1. *El espacio  $E_h^\theta$  es invariante por el operador de dos mallas  $M_h^H$ .*
2. *El operador de corrección en la malla grosera  $K_h^{2h}$  es representado en  $E_h^\theta$  por una matriz cuadrada de dimensión cuatro que denotaremos por  $\hat{K}_h^{2h}(\theta)$ . Dicha matriz viene dada por*

$$\hat{K}_h^{2h}(\theta) = \hat{I}_h - \hat{I}_{2h}^h(\theta)(\hat{L}_{2h}(2\theta))^{-1}\hat{I}_h^{2h}(\theta)\hat{L}_h(\theta),$$

para cada  $\theta \in T^{\text{low}}$ , donde  $\hat{I}_h, \hat{L}_h(\theta)$  son matrices  $(4 \times 4)$  que representan los operadores  $I_h, L_h(\theta)$ ,  $\hat{I}_h^{2h}(\theta)$  es una matriz  $(1 \times 4)$  que representa al operador  $I_h^{2h}(\theta)$ ,  $(\hat{L}_{2h}(2\theta))^{-1}$  es una matriz  $(1 \times 1)$  que representa el operador  $(L_{2h}(2\theta))^{-1}$  y, por último,  $\hat{I}_{2h}^h(\theta)$  es una matriz  $(4 \times 1)$  que representa el operador  $I_{2h}^h(\theta)$ .

3. Como los espacios  $E_h^\theta$  son invariantes por el operador de suavizado  $S_h$  (ver lema 2.3.3), es decir,  $S_h : E_h^\theta \rightarrow E_h^\theta$  para todo  $\theta \in T^{\text{low}}$ , tendremos una matriz  $(4 \times 4)$ , que denotaremos por  $\hat{M}_h^{2h}(\theta)$ , que representa al operador  $M_h^{2h}$  en  $E_h^\theta$  y que viene dada por

$$\hat{M}_h^{2h}(\theta) = \hat{S}_h(\theta)^{\nu_2} \hat{K}_h^{2h}(\theta) \hat{S}_h(\theta)^{\nu_1},$$

donde  $\hat{K}_h^{2h}(\theta)$  viene dado en la primera parte de este teorema y  $\hat{S}_h(\theta)$  es una matriz  $(4 \times 4)$  que representa el operador de suavizado  $S_h(\theta)$ .

Observamos que además de suponer la existencia de  $(L_{2h})^{-1}$  también estamos suponiendo la existencia de  $(L_h)^{-1}$  implícitamente. Es necesario suponer esto aunque, como veremos en el ejemplo siguiente, estas condiciones no se tienen por qué cumplir formalmente, ni si quiera en los casos más simples.

**Ejemplo 2.6.1** Tomando  $L_h = \Delta_h$  y  $L_{2h} = \Delta_{2h}$  en una red infinita, tendremos que para  $\theta = 0$ ,  $\varphi_h(\theta, \mathbf{x}) = 1$ , para todo  $h$ , se tiene

$$L_h \varphi_h(\theta, \mathbf{x}) = L_{2h} \varphi_{2h}(\theta, \mathbf{x}) = 0.$$

Entonces, los operadores  $L_h, L_{2h}$  tienen autovalores formales nulos, lo cual implica que no existe su inversa.

Con la intención de evitar el problema mostrado en el ejemplo anterior, razonaremos quitando el conjunto

$$\Lambda = \{\theta \in [-\pi/2, \pi/2]^2 : \tilde{L}_h = 0 \text{ o } \tilde{L}_{2h} = 0\}.$$

**Observación 2.6.1** El teorema 2.6.1 nos permite reducir el análisis de la convergencia del método en dos mallas al estudio del radio espectral y de las normas de matrices cuadradas de dimensión 4. Más adelante, cuando definamos el factor de convergencia asintótica, se verá esto aún más palpable.

A continuación obtendremos la representación matricial de  $M_h^{2h}$  estudiando la representación de cada una de las componentes del ciclo de dos mallas.

■ Operadores  $L_h, I_h$ .

Puesto que las funciones  $\varphi_h$  son autofunciones de  $L_h$ , tendremos que, por el lema 2.3.2, los coeficientes  $A^\alpha$  son transformados de la forma:

$$\begin{bmatrix} A^{(0,0)} \\ A^{(1,1)} \\ A^{(1,0)} \\ A^{(0,1)} \end{bmatrix} \longrightarrow \begin{bmatrix} \tilde{L}_h(\theta^{(0,0)}) & & & \\ & \tilde{L}_h(\theta^{(1,1)}) & & \\ & & \tilde{L}_h(\theta^{(1,0)}) & \\ & & & \tilde{L}_h(\theta^{(0,1)}) \end{bmatrix} \begin{bmatrix} A^{(0,0)} \\ A^{(1,1)} \\ A^{(1,0)} \\ A^{(0,1)} \end{bmatrix}.$$

Esta matriz diagonal será la matriz que denotamos en el teorema por  $\hat{L}_h(\theta)$ . Por otro lado, es trivial que la matriz del operador  $I_h$  es la matriz identidad  $(4 \times 4)$  en el espacio  $E_h^\theta$ .

- Operador restricción  $I_h^{2h}$ .

Observemos cómo el operador  $I_h^{2h}$  transforma las funciones  $\varphi(\boldsymbol{\theta}^\alpha, \cdot)$ .

$$\begin{aligned} I_h^{2h} \varphi_h(\boldsymbol{\theta}^\alpha, (x, y)) &= \frac{1}{16} [4 + 2(e^{i\theta_1} + e^{-i\theta_1} + e^{i\theta_2} + e^{-i\theta_2}) + e^{i\theta_1} e^{i\theta_2} \\ &\quad + e^{i\theta_1} e^{-i\theta_2} + e^{-i\theta_1} e^{i\theta_2} + e^{-i\theta_1} e^{-i\theta_2}] \varphi_h(\boldsymbol{\theta}^\alpha, (x, y)) \\ &= \frac{1}{4} (1 + \cos \theta_1)(1 + \cos \theta_2) \varphi_h(\boldsymbol{\theta}^\alpha, (x, y)), \end{aligned}$$

y puesto que en la red grosera  $\varphi_h(\boldsymbol{\theta}^\alpha, (x, y)) = \varphi_{2h}(2\boldsymbol{\theta}^{(0,0)}, (x, y))$ , tenemos que de  $E_h^\theta$  pasamos al espacio unidimensional  $\langle \varphi_{2h}(2\boldsymbol{\theta}^{(0,0)}, \cdot) \rangle$  y que la matriz que representa dicha aplicación en las correspondientes bases es

$$\hat{I}_h^{2h}(\boldsymbol{\theta}) = \frac{1}{4} \begin{bmatrix} (1 + \cos \theta_1)(1 + \cos \theta_2) \\ (1 + \cos(\theta_1 \pm \pi))(1 + \cos(\theta_2 \pm \pi)) \\ (1 + \cos(\theta_1 \pm \pi))(1 + \cos \theta_2) \\ (1 + \cos \theta_1)(1 + \cos(\theta_2 \pm \pi)) \end{bmatrix}^T,$$

siendo  $+$  cuando  $\theta_i < 0$  y siendo  $-$  cuando  $\theta_i \leq 0$ .

- Operador en la malla  $2h$ .

La función  $\varphi_{2h}(2\boldsymbol{\theta}^{(0,0)}, \cdot)$  es una autofunción de la elección natural, que hicimos en la sección 2.5.1, del operador  $L_{2h}$ . Más concretamente

$$L_{2h} \varphi_{2h}(2\boldsymbol{\theta}^{(0,0)}, (x, y)) = \frac{1}{2h^2} (-2 + \cos(2\theta_1) + \cos(2\theta_2)) \varphi_{2h}(2\boldsymbol{\theta}^{(0,0)}, (x, y)),$$

y, por tanto,

$$(\hat{L}_{2h}(2\boldsymbol{\theta}))^{-1} = \frac{2h^2}{(-2 + \cos(2\theta_1) + \cos(2\theta_2))}.$$

- Operador interpolador  $I_{2h}^h$ .

Aplicando el operador interpolador bilineal a la función  $\varphi_{2h}(2\boldsymbol{\theta}^{(0,0)}, \cdot)$  tendremos:

$$I_{2h}^h \varphi_{2h}(2\boldsymbol{\theta}^{(0,0)}, (x, y)) = \varphi_{2h}(2\boldsymbol{\theta}^{(0,0)}, (x, y)) \cdot \begin{cases} 1, & \text{si } \frac{x}{h}, \frac{y}{h} \text{ par,} \\ \cos(\theta_2), & \text{si } \frac{x}{h} \text{ par e } \frac{y}{h} \text{ impar,} \\ \cos(\theta_1), & \text{si } \frac{x}{h} \text{ impar e } \frac{y}{h} \text{ par,} \\ \cos(\theta_1) \cos(\theta_2), & \text{si } \frac{x}{h}, \frac{y}{h} \text{ impar,} \end{cases}$$

Puesto que toda función en  $E_h^\theta$  puede ser representada como una combinación lineal de la forma

$$a^{(0,0)} e^{i\boldsymbol{\theta}^{(0,0)} \mathbf{x}/h} + a^{(1,1)} e^{i\boldsymbol{\theta}^{(1,1)} \mathbf{x}/h} + a^{(1,0)} e^{i\boldsymbol{\theta}^{(1,0)} \mathbf{x}/h} + a^{(0,1)} e^{i\boldsymbol{\theta}^{(0,1)} \mathbf{x}/h},$$

los coeficientes  $a^{(0,0)}, a^{(1,1)}, a^{(1,0)}, a^{(0,1)}$  los obtendremos de evaluar esta última expresión y la de más arriba en los nodos de la forma  $(2kh, 2mh), (2kh, (2m+$

$1)h), ((2k+1)h, 2mh), ((2k+1)h, (2m+1)h)$ , igualando ambas expresiones, llegamos al sistema:

$$\begin{bmatrix} 1 \\ \cos \theta_2 \\ \cos \theta_1 \\ \cos \theta_1 \cos \theta_2 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & -1 & -1 & 1 \\ 1 & 1 & -1 & -1 \end{bmatrix} \begin{bmatrix} a^{(0,0)} \\ a^{(1,1)} \\ a^{(1,0)} \\ a^{(0,1)} \end{bmatrix}.$$

Por tanto, considerando la inversa de la matriz

$$\frac{1}{4} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & -1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & 1 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ \cos \theta_2 \\ \cos \theta_1 \\ \cos \theta_1 \cos \theta_2 \end{bmatrix} = \begin{bmatrix} a^{(0,0)} \\ a^{(1,1)} \\ a^{(1,0)} \\ a^{(0,1)} \end{bmatrix}.$$

Finalmente deducimos de todo lo anterior que la matriz que representa al operador interpolador bilineal  $\hat{I}_{2h}^h$  viene dada por

$$\hat{I}_{2h}^h(\boldsymbol{\theta}) = \frac{1}{4} \begin{bmatrix} (1 + \cos \theta_1)(1 + \cos \theta_2) \\ (1 - \cos \theta_1)(1 - \cos \theta_2) \\ (1 - \cos \theta_1)(1 + \cos \theta_2) \\ (1 + \cos \theta_1)(1 - \cos \theta_2) \end{bmatrix}.$$

- Operador del suavizado  $S_h$

Puesto que las funciones de Fourier  $\varphi_h(\boldsymbol{\theta}, \cdot)$  son autofunciones formales del operador  $S_h$ , donde  $S_h$  es el operador iteración de algún método de suavizado visto en este trabajo, es claro que  $S_h$  puede ser representado por la matriz  $(4 \times 4)$  diagonal

$$\hat{S}_h(\boldsymbol{\theta}) = \begin{bmatrix} \tilde{S}_h(\boldsymbol{\theta}^{(0,0)}) & & & \\ & \tilde{S}_h(\boldsymbol{\theta}^{(1,1)}) & & \\ & & \tilde{S}_h(\boldsymbol{\theta}^{(1,0)}) & \\ & & & \tilde{S}_h(\boldsymbol{\theta}^{(0,1)}) \end{bmatrix},$$

para cada  $\boldsymbol{\theta} \in T^{low}$ , donde  $\tilde{S}_h(\boldsymbol{\theta}^\alpha)$  es el autovalor formal asociado a la autofunción formal  $\varphi_h(\boldsymbol{\theta}^\alpha, \cdot)$ . Por ejemplo para GSLEX, como ya vimos, tendremos que  $\tilde{S}_h(\boldsymbol{\theta}) = (e^{-i\theta_1} + e^{-i\theta_2}) / (4 - e^{-i\theta_1} - e^{-i\theta_2})$ .

Con esto habríamos demostrado, de forma constructiva (pues damos todas las matrices), el teorema 2.6.1.

### 2.6.1. Factor de convergencia asintótico

Empecemos definiendo un concepto mencionado anteriormente

**Definición 2.6.2** *El factor de convergencia asintótico  $\rho_{loc}(M_h^{2h})$  es el valor que verifica*

$$\rho_{loc}(M_h^{2h}) = \sup\{\rho(\hat{M}_h^{2h}(\boldsymbol{\theta})) : \boldsymbol{\theta} \in T^{low} \setminus \Lambda\}.$$

Considerando nuestro problema modelo y el método multigrid en dos redes con GSLEX como método de suavizado, FW como operador restricción y el interpolador bilineal como operador de interpolación, hemos aproximado con MATLAB el factor de convergencia asintótico, que notemos que es independiente de  $h$  porque las matrices resultantes  $M_h^{2h}$  son independientes de  $h$ . Así, hemos obtenido la tabla 2.1, en la que se observa que el aumento de los valores  $\nu_1, \nu_2$  reduce ambos factores, sin embargo, el factor de suavizado  $\mu_{loc}$  se ve reducido más rápido que el factor de convergencia asintótica  $\rho_{loc}$ . Por lo tanto, si tenemos

$\nu_1, \nu_2$	$\mu_{loc}^{\nu_1+\nu_2}$	$\rho_{loc}$
1,0	0.500	0.400
1,1	0.250	0.193
2,1	0.125	0.119
2,2	0.063	0.084

Tabla 2.1: Comparación del método de suavizado y del factor de convergencia asintótico cuando se aplica un número distinto de veces el suavizado.

demasiadas iteraciones del suavizado entonces puede ocurrir que la aproximación en la red más gruesa no haga frente a los efectos del suavizado. Para nuestro problema modelo no es recomendable tomar  $\nu_1 + \nu_2 \geq 4$ , ya que la mejora en el número de iteraciones para obtener convergencia no compensa el coste computacional mayor que supone.

## 2.7. El ciclo multigrid

Hasta ahora hemos descrito los principios del ciclo multigrid solo para dos mallas de puntos. En el contexto de la teoría del método multigrid esto es poco práctico, debido a la complejidad que aún supone el problema en la red gruesa.

A continuación, pasaremos de dos mallas a varias mallas. La idea del multigrid reside en obtener una buena convergencia del método en dos redes si damos una aproximación de la ecuación del defecto (2.7) en lugar de resolverla. Una manera de obtener dicha aproximación, sin tener una pérdida en velocidad de convergencia, puede venir de aplicar la idea en dos redes de nuevo, implicando mallas más gruesas que  $\Omega_H$ . Esta idea, aplicada recursivamente, nos terminará llevando a la red más gruesa, que puede estar formada por un único punto, y donde resolver la ecuación del defecto será más sencillo. Denotaremos por  $\gamma$  al índice del ciclo, que hará referencia al número de veces que se repite el ciclo multigrid. Para  $\gamma = 1$  hablamos de V-ciclos y para  $\gamma = 2$  hablamos de W-ciclos.

### 2.7.1. Definición recursiva del método multigrid

Consideramos una secuencia de mallas cada vez más gruesas  $\Omega_{h_k}$ , cuyo tamaño de malla es  $h_k$  y que para aligerar la notación denotaremos simplemente

por  $\Omega_k$ . La malla más gruesa será denotada por  $\Omega_0$  y la más fina por  $\Omega_l$ . Consecuentemente, denotaremos las componentes del método multigrid en la sección 2.5 por

$$\begin{aligned} L_k : \mathcal{G}(\Omega_k) &\longrightarrow \mathcal{G}(\Omega_k), & S_k : \mathcal{G}(\Omega_k) &\longrightarrow \mathcal{G}(\Omega_k), \\ I_k^{k-1} : \mathcal{G}(\Omega_k) &\longrightarrow \mathcal{G}(\Omega_k), & I_{k-1}^k : \mathcal{G}(\Omega_k) &\longrightarrow \mathcal{G}(\Omega_k), \end{aligned}$$

donde  $L_k$  es la discretización de  $L$  ( $= \Delta$ , en nuestro problema modelo) en la malla  $\Omega_k$  para  $k = l, \dots, 0$ , y donde la ecuación original se discretiza por

$$L_l u_l = f_l \quad (\Omega_l), \quad (2.10)$$

que es el problema a resolver. El operador  $S_k$  denota el operador lineal correspondiente al método de suavizado que esta siendo aplicado en la malla  $\Omega_k$ . De la misma forma que en el ciclo de dos mallas,  $\nu$  será el número de veces que aplicamos el suavizado al problema discreto de la forma  $L_k u_k = f_k$  para alguna aproximación inicial  $w_k$ . El resultado de aplicar el suavizado será  $\bar{w}_k$  y denotaremos el proceso por

$$\bar{w}_k = SMOOTH^\nu(w_k, L_k, f_k).$$

Con estas distinciones hechas vamos a describir el ciclo multigrid o, más concretamente, el ciclo de  $(l + 1)$  mallas para resolver el problema (2.10) para un  $l$  fijo.

**Ciclo multigrid**  $u_k^{m+1} = MGCCY(k, \gamma, u_k^m, L_k, f_k, \nu_1, \nu_2)$

1. Presuavizado.

- Calculamos  $\bar{u}_k^m$  aplicando el suavizado  $\nu_1$  veces con iterante inicial  $u_k^m$ :

$$\bar{u}_k^m = SMOOTH^{\nu_1}(u_k^m, L_k, f_k).$$

2. Corrección en malla gruesa.

- Obtenemos la ecuación del defecto  $\bar{d}_k^m = f_k - L_k \bar{u}_k^m$ .
- Le aplicamos el operador restricción  $\bar{d}_{k-1}^m = I_k^{k-1} \bar{d}_k^m$ .
- Calculamos la solución aproximada  $\hat{v}_{k-1}^m$  de la ecuación del defecto en la malla  $\Omega_{k-1}$

$$L_{k-1} \hat{v}_{k-1}^m = \bar{d}_{k-1}^m, \quad (2.11)$$

utilizando el siguiente procedimiento:

- Si  $k = 1$ , entonces resolvemos exactamente la ecuación (2.11).
- Si  $k > 1$ , entonces obtendremos la aproximación de (2.11) aplicando  $\gamma$  veces el ciclo de  $k$  mallas utilizando la función nula en la malla  $\Omega_{k-1}$  como primera aproximación

$$\hat{v}_{k-1}^m = MGCCY^\gamma(k-1, \gamma, 0, L_{k-1}, \bar{d}_{k-1}^m, \nu_1, \nu_2).$$

- Interpolamos la corrección  $\hat{v}_k^m = I_{k-1}^k \hat{v}_{k-1}^m$ .
- Calculamos la aproximación en  $\Omega_k$   $u^{m,DCMG} = \bar{u}_k^m + \hat{v}_k^m$ .

3. Postsuavizado.

- Calculamos  $u_k^{m+1}$  aplicando el suavizado  $\nu_2$  veces con iterante inicial  $u^{m,DCMG}$ :

$$u_k^{m+1} = SMOOTH^{\nu_2}(u^{m,DCMG}, L_k, f_k).$$

Denotaremos por  $M_k$  al operador de iteración del método multigrid descrito anteriormente. Con todo esto tendremos el siguiente teorema

**Teorema 2.7.1** *El operador de iteración del método multigrid  $M_l$  viene dado por la siguiente recursión:*

$$\begin{aligned} M_0 &= 0 \\ M_k &= S_k^{\nu_2} (I_k - I_{k-1}^k (I_{k-1} - M_{k-1}^\gamma) L_{k-1}^{-1} I_k^{k-1} L_k) S_k^{\nu_1} \\ &\quad (k = 1, \dots, l). \end{aligned} \quad (2.12)$$

**Demostración.**

Es claro que para el caso en el que  $l = 1$  estamos en el ciclo de dos mallas y la fórmula (2.12) se convierte en

$$M_1 = S_1^{\nu_2} (I_1 - I_0^1 L_0^{-1} I_1^0 L_1) S_1^{\nu_1}.$$

Este operador es el mismo que vimos en la sección de la corrección en malla grosera, al cual le hemos aplicado un presuavizado y un postsuavizado. Por lo tanto, para el caso de dos mallas ya tenemos demostrado este resultado.

Ahora bien, para la multimalla tenemos, por el proceso descrito anteriormente, que el operador viene dado por

$$M_{k+1} = S_{k+1}^{\nu_2} [I_{k+1} - I_k^{k+1} \hat{L}_k^{-1} I_{k+1}^k L_{k+1}] S_{k+1}^{\nu_1},$$

donde  $\hat{L}_k$  denota el operador que viene de la aproximación  $L_k \bar{v}_k^m = \bar{d}_k^m$ , más concretamente el operador que viene de la ecuación  $\hat{L}_k \hat{v}_k^m = \bar{d}_k^m$ . Por lo tanto, tenemos que ver quién es el operador  $\hat{L}_k^{-1}$ . Dicho operador es  $M_k$  aplicado  $\gamma$  veces con aproximación inicial la nula y término independiente el defecto. Por lo tanto, razonando como en (2.8) tendremos que  $\hat{L}_k^{-1} = (I_k - M_k^\gamma) L_k^{-1}$  y sustituyendo en la ecuación anterior tenemos

$$M_{k+1} = S_{k+1}^{\nu_2} [I_{k+1} - I_k^{k+1} ((I_k - M_k^\gamma) L_k^{-1}) I_{k+1}^k L_{k+1}] S_{k+1}^{\nu_1},$$

con lo que hemos demostrado el resultado.  $\square$

La representación gráfica de los métodos multigrid para diferentes valores de  $l$  suele hacerse como en la figura 2.5 donde queda claro por qué al ciclo multigrid correspondiente a  $\gamma = 1$  se la llama V-ciclos y al correspondiente a  $\gamma = 2$ , se le denomina W-ciclo.

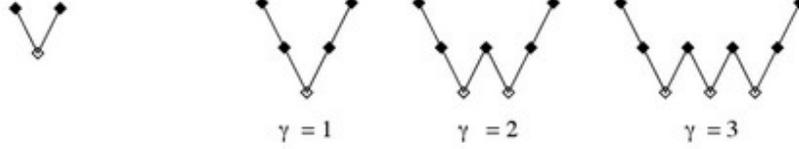


Figura 2.5: Estructura del ciclo multigrad para dos mallas (primera figura) y para tres mallas (el resto) y disntitos valores de  $\gamma$  .

### 2.7.2. Coste computacional

En esta sección vamos a estimar el coste computacional del método multigrad descrito en este trabajo. Por coste computacional entenderemos el número de operaciones aritméticas necesarias que se producen al aplicar el método. Siguiendo con la notación de la sección anterior denotaremos por  $W_l$  al coste computacional tras un ciclo multigrad en  $\Omega_l$ . Podemos convencernos fácilmente de que el coste computacional viene dado recursivamente por

$$\begin{aligned} W_1 &= W_1^0 + W_0, \\ W_{k+1} &= W_{k+1}^k + \gamma W_k, \quad (k = 1, \dots, l-1) \end{aligned} \quad (2.13)$$

donde  $W_{k+1}^k$  denota el coste computacional del ciclo de dos mallas  $(h_{k+1}, h_k)$  quitando el trabajo necesario para resolver la ecuación del defecto en  $\Omega_k$  y  $W_0$  denota el trabajo necesario para resolver exactamente la ecuación del defecto en la red más gruesa  $\Omega_0$ . Esta formulación recursiva simplemente nos dice que el coste computacional del método multigrad es la suma del coste computacional de cambiar a una malla más gruesa más el coste computacional de volver aplicar el método, un número  $\gamma$  de veces, en la red más gruesa.

Si suponemos que  $\gamma$  es independiente de la malla, es decir de  $k$ , tendremos inductivamente la siguiente fórmula para el trabajo computacional para un ciclo multigrad:

$$W_l = \sum_{k=1}^l \gamma^{l-k} W_k^{k-1} + \gamma^{l-1} W_0 \quad (l \geq 1). \quad (2.14)$$

Denotaremos por  $N_k$  al número de puntos en la red  $\Omega_k$ . Por lo tanto, para el engrosamiento estándar en dos dimensiones tendremos que  $N_k \doteq 4N_{k-1}$ , para  $k = 1, \dots, l$  y donde  $\doteq$  quiere decir igualdad salvo términos de orden inferior (por ejemplo, los de la frontera de la malla). Supondremos que las componentes del método multigrad (relajación, cálculo del defecto, operaciones de restricción e interpolación) requieren un número de operaciones aritméticas por punto, de la correspondiente malla, que englobaremos en una constante  $C$  independiente de la malla, es decir, independiente de  $k$ . Más concretamente, esto se traduce en la desigualdad siguiente

$$W_k^{k-1} \leq CN_k (k = 1, \dots, l-1), \quad (2.15)$$

donde la desigualdad con el punto significa desigualdad salvo términos de orden inferior y

$$C = (\nu\omega_0 + \omega_1 + \omega_2),$$

donde  $\nu = \nu_1 + \nu_2$  es el número de iteraciones de suavizado,  $\omega_0$  es la medida de coste computacional por punto en la red de una iteración del suavizado,  $\omega_1$  es la medida de coste computacional por punto en la red del cálculo del defecto y su transferencia a la red  $\Omega_{k-1}$  y  $\omega_2$  es la medida de coste computacional por punto de red de la interpolación de la corrección a  $\Omega_k$  y su suma a la aproximación previa.

Veamos a continuación el coste computacional total del método multigrid, estudiado en la sección anterior, para distintos valores de  $\gamma$ , supondremos que  $W_0$  es despreciable,

- Para  $\gamma = 1$ .

$$W_l = \sum_{k=1}^l W_k^{k-1} + W_0 \leq C(N_1 + N_2 + \dots + N_l),$$

donde hemos aplicado (2.15). Puesto que estamos trabajando con el engrosamiento estándar,  $N_{l-k} = N_l/4^k$ . Aplicando esto a la desigualdad de arriba tendremos

$$W_l \leq C(N_1 + N_2 + \dots + N_l) \leq C\left(\sum_{k=0}^{l-1} \frac{N_l}{4^k}\right) \leq CN_l \frac{4}{3} \left(1 - \frac{1}{4^l}\right) \leq CN_l \frac{4}{3}.$$

- Para  $\gamma = 2$ .

De nuevo aplicando (2.14), (2.15) tendremos

$$W_l = \sum_{k=1}^l 2^{l-k} W_k^{k-1} \leq CN_l \left(\sum_{k=0}^{l-1} \frac{1}{2^k}\right) \leq 2CN_l,$$

- Para  $\gamma = 3$ .

Razonando como en los casos anteriores tenemos

$$W_l = \sum_{k=1}^l 3^{l-k} W_k^{k-1} \leq CN_l \left(\sum_{k=0}^{l-1} \frac{1}{4^k} 3^k\right) = 4CN_l \left(1 - \left(\frac{3}{4}\right)^l\right) + W_0 \leq 4CN_l,$$

- Para  $\gamma = 4$ .

Como anteriormente tendremos

$$W_l = \sum_{k=1}^l 4^{l-k} W_k^{k-1} + 4^{l-1} W_0 \leq CN_l \left(\sum_{k=0}^{l-1} \frac{1}{4^k} 4^k\right) = CN_l l = O(N_l \log(N_l)).$$

A modo de resumen, de estos casos que acabamos de desarrollar, tenemos que cuando  $W_0$  se considera despreciable la estimación de  $W_l$  muestra que el número de operaciones aritméticas necesarias para un ciclo multigrid en 2D es proporcional al número de puntos de la red más fina para  $\gamma \leq 3$  y el engrosamiento estándar. Dicha constante de proporcionalidad dependerá del tipo de ciclo, es decir de  $\gamma$ .

### 2.7.3. Convergencia y eficiencia del método multigríd

En esta subsección, mostraremos los resultados de aplicar al problema modelo con  $f^\Omega \equiv 0$  y  $f^\Gamma \equiv 0$ , con el engrosamiento estándar, la implementación descrita en el teorema 2.7.1 del método multigríd. Las componentes del método que consideraremos serán las explicadas anteriormente en el trabajo, es decir, el método de suavizado será GSLEX, el operador restricción será el FW y el operador interpolador será el interpolador bilineal. El tamaño de la red más fina considerado será  $h_7 = 1/2^8$  y el tamaño de la red más gruesa será  $h_0 = 1/2$ . Tomaremos como iterante inicial  $u^0 = [1, 1, \dots, 1]^T$ .

Representaremos en escala semilogarítmica el error del defecto (2.7) cometido en norma  $L^2$ -discreta frente al número de iteraciones realizadas. Siguiendo con la notación vista en la sección 2.7.1, denotaremos por  $V(\nu_1, \nu_2)$  al método multigríd con parámetro  $\gamma = 1$  y por  $W(\nu_1, \nu_2)$  al método multigríd con parámetro  $\gamma = 2$ . ( $\nu_1, \nu_2$  denotan el número de pasos de pre y post suavizado respectivamente). Como se observa en la figura 2.6, los métodos  $V(1,1)$  y  $W(1,1)$  son los que

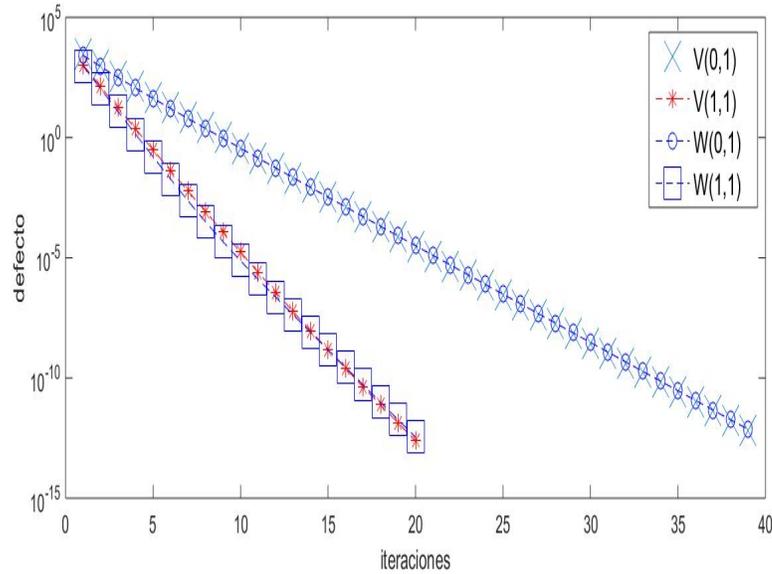


Figura 2.6: Representación de la norma del defecto frente las iteraciones para los ciclos multigríd  $V(0,1)$ ,  $V(1,1)$ ,  $W(0,1)$ ,  $W(1,1)$ .

presentan mayor velocidad de convergencia, ya que a partir de 20 iteraciones el defecto ya es menor que  $10^{-12}$  mientras que los otros métodos necesitan 40 iteraciones.

### Medición del factor de convergencia del método multigrid en la práctica

Para evaluar y analizar una iteración del método multigrid, buscaremos determinar su factor de convergencia  $\rho$  empíricamente. En general, nosotros solo dispondremos de los defectos  $d_h^m$  ( $m = 1, 2, \dots$ ) y, por lo tanto, podremos medir las siguientes cantidades

$$q^{(m)} := \frac{\|d_h^m\|}{\|d_h^{m-1}\|}, \quad \hat{q}^{(m)} := \sqrt[m]{\frac{\|d_h^m\|}{\|d_h^0\|}},$$

para la norma  $\|\cdot\|_2$ . Interpretaremos  $\hat{q}^{(m)}$  como un promedio del factor de reducción del defecto tras  $m$  iteraciones. Por supuesto, debemos suponer, para que estas cantidades tengan sentido, que  $d_h^0 \neq 0$  y parece que  $\hat{q}^{(m)}$  es un buen estimador del factor de convergencia  $\rho$  si  $m$  está suficientemente avanzado. En la tabla 2.2, mostramos los resultados obtenidos para nuestro problema modelo en las siguientes tablas.

Método	$m$	$q^m$	$\hat{q}^m$
V(0,1)	20	0.393	0.378
V(1,1)	20	0.179	0.149
W(0,1)	20	0.395	0.379
W(1,1)	20	0.188	0.152
GSLEX	1000	0.999	0.994
JACOBI	1000	0.999	0.995
SOR	1000	0.978	0.981

Tabla 2.2: Comparación de los valores de  $q^m$  y  $\hat{q}^m$  para los distintos métodos.

Método	Malla 8	Malla 7	Malla 5	Malla 3	Malla 1
V(0,1)	0.379	0.378	0.377	0.355	0.164
V(1,1)	0.147	0.149	0.154	0.142	0.055
W(0,1)	0.379	0.379	0.377	0.359	0.164
W(1,1)	0.149	0.151	0.157	0.138	0.055
GSLEX	0.994	0.994	0.993	0.959	0.500
JACOBI	0.995	0.995	0.995	0.979	0.707
SOR	0.993	0.981	0.910	0.678	0.475

Tabla 2.3: Comparación de los valores  $\hat{q}^m$  ( $m = 20$  para multigrid y  $m = 1000$  para los métodos multigrid) para los distintos métodos en distintos tamaños de mallas donde el tamaño de la malla  $k$  es  $h = 1/2^{k+1}$ .

Observamos cómo los métodos multigrid tienen un factor de convergencia mucho más pequeño que los métodos iterativos clásicos. También vemos que los

métodos multigrid que menor factor de convergencia tienen son los ciclos  $V(1,1)$  y  $W(1,1)$ , lo cual concuerda con lo visto en la gráfica de la figura 2.6 pues estos eran los métodos que antes reducían el defecto más rápidamente.

Por otro lado, en la tabla 2.3 vemos cómo el aumento en el tamaño de la malla más fina no hace que el factor de convergencia de los métodos multigrid se acerque a 1. Sin embargo, sí lo hace en los métodos clásicos. Esto se debe a la convergencia  $h$ -independiente de los métodos multigrid que veremos un poco más adelante.

### Eficiencia numérica

Por último, consideramos la eficiencia numérica del método multigrid, es decir, no solo tendremos en cuenta su velocidad de convergencia sino también su coste computacional. En este sentido, la cantidad que más nos interesa para el estudio de la eficiencia numérica será el tiempo necesario para resolver nuestro sistema lineal con una precisión dada. Mostramos en la figura 2.7 los tiempos invertidos para alcanzar una norma  $L^2$  discreta del defecto menor que  $10^{-12}$  en cada método.

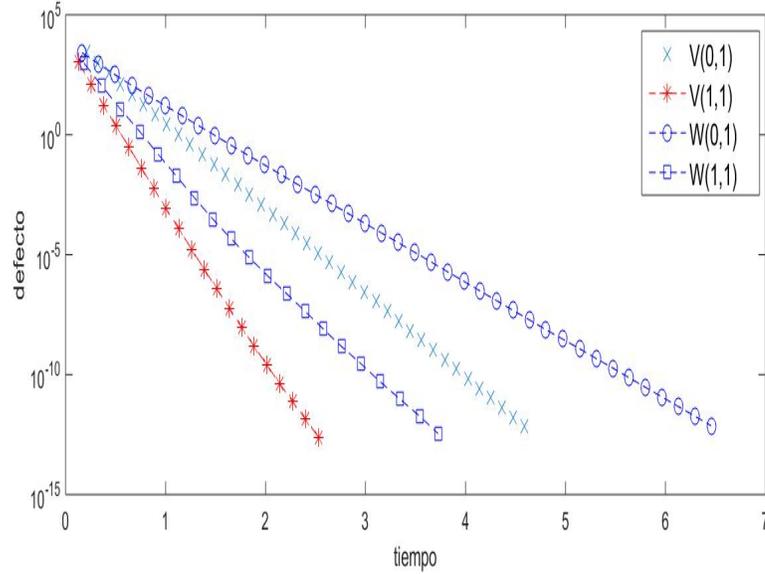


Figura 2.7: Representación de la norma del defecto tras varias iteraciones para los ciclos multigrid frente al tiempo

Podemos observar cómo el aumento del coeficiente  $\gamma$  implica un aumento en el tiempo de resolución del problema. Por otro lado, la figura 2.7 nos indica que entre  $V(1,1)$  y  $W(1,1)$ , que eran los que mejor factor de convergencia tenían y además mostraban un comportamiento muy similar en la figura 2.6, el más

eficiente numéricamente es  $V(1,1)$  puesto que es el que menos tarda en resolver.

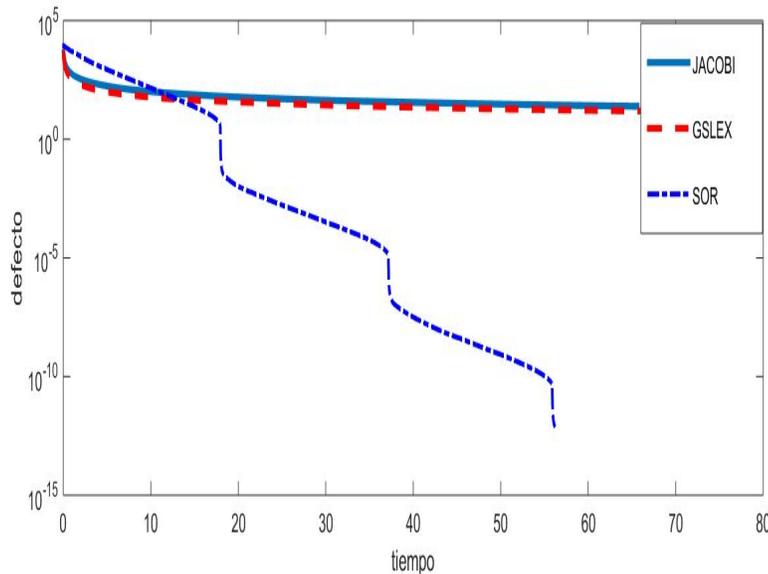


Figura 2.8: Representación de la norma del defecto tras varias iteraciones para los métodos clásicos frente al tiempo

Por lo tanto, en la figura 2.8 vemos cómo los métodos clásicos tardan mucho más tiempo en lograr reducir la norma del defecto un factor de  $10^{-12}$ . Solo el método S.O.R. lo logra antes de los 70 segundos frente a los 6.5, aproximadamente, que tardaba  $W(1,1)$  (el más lento de los métodos multigrad estudiados en la figura 2.7).

## 2.8. $h$ - Independencia de la convergencia

En esta sección veremos cómo podemos obtener cotas para  $\|M_k\|$  independientes del tamaño de la malla  $h = 1/2^{k+1}$  para  $W$ -ciclos. Para el caso de los  $V$ -ciclos es necesario un estudio más sofisticado y no será realizado en este trabajo. Veremos cómo una simple y general estimación de cotas de  $\|M_k\|$  será generada a partir de cotas de la norma del operador en dos mallas  $\|M_{k+1}^k\|$ . Dicha estimación recursiva estará basada en el teorema 2.7.1.

De forma más concreta, trataremos de demostrar que si partimos de un método de dos mallas tal que

$$\|M_{k+1}^k\| \leq \sigma^*,$$

donde  $\sigma^*$  es un valor suficientemente pequeño independiente de  $h$ , entonces el correspondiente método multigrad con  $\gamma \geq 2$  tendrá propiedades similares al

operador de dos mallas. Como vimos a lo largo de la prueba del teorema 2.7.1, el método multigrad se puede interpretar como una modificación del operador en dos mallas. Esto se resume en el siguiente corolario que es consecuencia inmediata del teorema 2.7.1:

**Corolario 2.8.1** *Para  $k = 1, \dots, l - 1$ , se tiene la ecuación*

$$M_{k+1} = M_{k+1}^k + A_k^{k+1} M_k^\gamma A_{k+1}^k, \quad (2.16)$$

donde

$$\begin{aligned} A_k^{k+1} &:= S_{k+1}^{\nu_2} I_k^{k+1} : \mathcal{G}(\Omega_k) \longrightarrow \mathcal{G}(\Omega_{k+1}), \\ A_{k+1}^k &:= L_k^{-1} I_{k+1}^k L_{k+1} S_{k+1}^{\nu_1} : \mathcal{G}(\Omega_{k+1}) \longrightarrow \mathcal{G}(\Omega_k), \end{aligned}$$

y  $M_{k+1}^k : \mathcal{G}(\Omega_{k+1}) \longrightarrow \mathcal{G}(\Omega_{k+1})$  es

$$M_{k+1}^k = S_{k+1}^{\nu_2} (I_{k+1} - I_k^{k+1} L_k^{-1} I_{k+1}^k L_{k+1}) S_{k+1}^{\nu_1}.$$

Y con esta formulación tenemos el siguiente teorema

**Teorema 2.8.1** *Si suponemos las siguientes estimaciones uniformes respecto de  $k(\leq l - 1)$ :*

$$\|M_{k+1}^k\| \leq \sigma^*, \quad \|A_k^{k+1}\| \|A_{k+1}^k\| \leq C, \quad (2.17)$$

entonces  $\|M_l\| \leq \eta_l$ , donde  $\eta_l$  viene definida recursivamente por

$$\eta_1 := \sigma^*, \quad \eta_{k+1} := \sigma^* + C\eta_k^\gamma \quad (k = 1, \dots, l - 1). \quad (2.18)$$

Si además suponemos que

$$4C\sigma^* \leq 1 \quad \text{y} \quad \gamma = 2, \quad (2.19)$$

obtenemos la estimación uniforme:

$$\|M_l\| \leq \eta := (1 - \sqrt{1 - 4C\sigma^*})/2C \leq 2\sigma^* \quad (l \geq 1). \quad (2.20)$$

**Demostración.**

La ecuación (2.18) es inmediata de aplicar (2.16), (2.17). Encuanto a la ecuación (2.20), vemos que la sucesión  $\{\eta_l\}$  es convergente probando que es creciente y acotada. Notemos para ello que

$$\begin{aligned} \eta_2 &= \sigma^* (1 + C\sigma^*) \\ \eta_3 &= \sigma^* (1 + C\sigma^* (1 + C\sigma^*)^2) \\ \eta_4 &= \sigma^* \left[ 1 + C\sigma^* \left( 1 + C\sigma^* (1 + C\sigma^*)^2 \right)^2 \right]^2 \\ &\vdots \end{aligned} \quad (2.21)$$

De aquí, como  $\sigma^* \leq 1/(4C)$ , entonces  $\eta_2 \leq 5/(16C)$ ,  $\eta_3 \leq 89/(256C)$ , ... y, en general, si  $\eta_j \leq 1/(2C)$ , entonces

$$\eta_{j+1} \leq \frac{1}{4C} + \frac{C}{4C^2} = \frac{1}{2C}.$$

Con lo que la sucesión está acotada. El hecho de que sea creciente se deduce del hecho de que  $\eta_{j+1}$  se calcula a partir de  $\eta_j$  en (2.21) multiplicando el  $C\sigma^*$  más interior por  $(1 + C\sigma^*)^2$ .

Pasando entonces al límite y llegando a la ecuación  $\eta = \sigma^* + C\eta^2$ . Resolviendo esa ecuación para  $\eta$  y quedándonos con la raíz más pequeña

$$\eta = \frac{1 - \sqrt{1 - 4C\sigma^*}}{2C} \leq 2\sigma^*.$$

□

Con este último teorema tendríamos que la norma del operador del método multigrad  $\|M_l\|$ , está acotada superiormente por un cierto valor  $\eta$  independiente del tamaño de la red. La existencia de las cotas superiores independientes de  $h$  de (2.20) y el hecho de que el número de operaciones por ciclo multigrad es  $O(N)$  (como vimos en la sección del trabajo computacional 2.7.2) implican la optimización del método. Por lo tanto, concluimos que para la construcción del método multigrad normalmente es suficiente con el análisis del método en dos redes.



## Capítulo 3

# Conclusiones Finales

En este último capítulo pondremos las conclusiones finales a las que hemos llegado tras la realización del trabajo.

- La principal consecuencia que debemos extraer de este trabajo es que, para resolver el problema discreto del problema modelo, una buena opción es utilizar métodos multigrad. Más concretamente, si consideramos el problema modelo y una malla de puntos con el orden lexicográfico, entonces el ciclo  $V(1,1)$ , con las componentes indicadas en la sección 2.7.3, nos dar excelentes resultados comparados con los métodos clásicos. Esto no quiere decir que hayamos obtenido el mejor método multigrad para resolver el problema. Es decir, como hemos indicado a lo largo del trabajo, la elección de las componentes del ciclo multigrad puede darnos mejores o peores aproximaciones. Ahora bien, suponen una mejora de los métodos iterativos clásicos, como ha quedado reflejado en el ejemplo del problema modelo.
- Los métodos multigrad no solo reducen el defecto en un menor número de iteraciones comparado con los métodos clásicos, sino que también tardan menos en reducir el defecto. Por lo tanto, son mucho más eficientes.
- Otra conclusión importante es la independencia de la convergencia de los ciclos multigrad frente al diámetro de la malla (ver teorema 2.8.1 para los  $W$ -ciclos). Esto lo que nos dice es que, si tenemos una buena convergencia del método multigrad en una cierta malla, podremos tener, sustancialmente, buena convergencia en mallas más finas. Esto no es cierto en los métodos clásicos como se veía numéricamente en la tabla 2.3.
- Hemos visto que es más eficiente suavizar con Gauss-Seidel que con  $\omega$ -Jacobi o S.O.R.
- Para el problema modelo, la convergencia es mejor con S.O.R. que con Gauss-Seidel y mejor con Gauss-Seidel que con Jacobi. Hemos obtenido también el  $\omega$  óptimo para el método S.O.R. en el capítulo 1.

- El análisis local de Fourier es una herramienta básica en este trabajo para el estudio de las propiedades de suavizado del método de Gauss-Seidel. También ayuda en el estudio del factor de convergencia del método en dos mallas. Sin embargo, hay que observar que en el análisis local del Fourier se está obviando la frontera del problema modelo, que es el que estamos tratando. De hecho, se aplica dicho análisis en una malla infinita de puntos y esto es debido a que se trata de un estudio local de las propiedades en las que estamos interesados. La filosofía de este análisis local de Fourier es el estudio cuantitativo de algunas propiedades al que después habría que aplicar un correspondiente tratamiento en la frontera.
- Una observación que se ha indicado en el trabajo es que el ciclo multigrad es una modificación del ciclo de dos mallas. En el ciclo de dos mallas aproximábamos la solución de la ecuación del defecto resolviendo en la malla más gruesa. En el ciclo multigrad tenemos la misma idea salvo que ahora la aproximación la obtenemos de aplicar sucesivamente el ciclo de dos mallas hasta llegar a la malla más gruesa, donde finalmente resolvemos.
- El coste computacional de los métodos multigrad para el modelo concreto de este trabajo es proporcional al número de puntos en la malla fina, siempre que consideremos el coste computacional de resolver en la malla más gruesa despreciable. Esto, junto a la  $h$ -independencia de la convergencia, nos proporciona un método muy adecuado para resolver sistemas lineales resultantes de la discretización de problemas elípticos donde es importante que la red sea fina para aproximar convenientemente el problema continuo. Un coste proporcional al número de puntos en la malla fina para cada iteración es comparable al coste de cada iteración de un método iterativo clásico y, debido a su mucho mejor factor de convergencia, conduce a métodos mucho más eficientes.

# Apéndice

```
%JACOBI
function U = JacobiM(b, f)
%La entrada es una matriz (2^(N)-1)x(2^(N)-1) que es
%iterante inicial de
%Jacobi y que cada entrada es el valor de la primera
%aproximacion de la funcion en la red de puntos
[m,n] = size(b);
h = 1/(m+1);
%La completamos con ceros la matriz para simular la malla
Z = zeros(1,m);
B = [Z;b;Z];
Z = zeros(m+2,1);
B = [Z,B,Z];

%Construimos el iterante que guardaremos en la matriz U
[r,s] = size(B);
U = zeros(r,s);
for i = m : -1 : 1
    for j = 1 : n

        U(i+1,j+1)= -(h^2/4)*(f(i,j)-(1/h)^2*
            (B((i+1)-1,(j+1))+B((i+1),(j+1)+1)+B((i+1)+1,(j+1))
            +B((i+1),(j+1)-1)));

    end
end

U(:,s) = [];
U(:,1) = [];
U(r,:) = [];
U(1,:) = [];

end
```

```

%Gauss Seidel
function U = GSM(b, f)
%La entrada son dos matrices (2^(N)-1)x(2^(N)-1)
%que es iterante inicial y
%el termino independiente
[m,n] = size(b);
h = 1/(m+1);
%La completamos con ceros la matriz para simular la malla
Z = zeros(1,m);
B = [Z;b;Z];
Z = zeros(m+2,1);
B = [Z,B,Z];

%Construimos el iterante que guardaremos en la matriz U
[r,s] = size(B);
U = zeros(r,s);
for i = m : -1 : 1
    for j = 1 : n

        U(i+1,j+1) = -(h^2/4)*(f(i,j)-(1/h)^2*
            (B((i+1)-1,(j+1))+B((i+1),(j+1)+1)+
            U((i+1)+1,(j+1))+U((i+1),(j+1)-1)));

    end
end

U(:,s) = [];
U(:,1) = [];
U(r,:) = [];
U(1,:) = [];

end

```

```

%SOR
function U = SORM(b, f)

%La entrada es una matriz (2^(N)-1)x(2^(N)-1) que es
%iterante inicial
[m,n] = size(b);
h = 1/(m+1);
w = 2/(1+sin(pi*h));

%La completamos con ceros la matriz para simular la malla
Z = zeros(1,m);
B = [Z;b;Z];

```

```

Z = zeros(m+2,1);
B = [Z,B,Z];

%Construimos el iterante que guardaremos en la matriz U
[r,s] = size(B);
U = zeros(r,s);
for i = m : -1 : 1
    for j = 1 : n

        U(i+1,j+1)= -w*(h^2/4)*(f(i,j)-(1/h)^2*
            (B((i+1)-1,(j+1))+B((i+1),(j+1)+1)+U((i+1)+1,(j+1))
            +U((i+1),(j+1)-1)))+(1-w)*B(i+1,j+1);

    end
end

U(:,s) = [];
U(:,1) = [];
U(r,:) = [];
U(1,:) = [];

end

```

```

%Operador FW
function U = FWM(b)

%La entrada es una matriz  $(2^{(N)}-1) \times (2^{(N)}-1)$  que es
%iterante inicial
[m,n] = size(b);

%Construimos la matriz donde guardaremos la restriccion
U = zeros(floor(m/2),floor(n/2));

for i = m-1 : -2 : 2
    for j = 2: 2 : n-1

        U(i/2,j/2) = 1/16*(4*b(i,j)+2*b(i-1,j)+2*b(i,j-1)
            +2*b(i+1,j)+2*b(i,j+1)+b(i-1,j+1)+b(i-1,j-1)
            +b(i+1,j-1)+b(i+1,j+1));

    end
end

end

```

```

%Aplicar el operador L
function U = OperadorM(b,k)

[m,n] = size(b);

%La completamos con ceros la matriz para simular la malla
Z = zeros(1,m);
B = [Z;b;Z];
Z = zeros(m+2,1);
B = [Z,B,Z];

U = zeros(m,n);
for i = 1 : 1 :m
    for j = 1 : 1 : n

        U(i,j) = 1/(1/2^(k+1))^2 *
            (-4*B(i+1,j+1) + B(i+1,j) + B(i+1,j+2) +
            B(i,j+1) + B(i+2,j+1));

    end
end
end

```

```

%Operador interpolador Bilineal
function U = IBM(b)

%La entrada es una matriz (2^(N)-1)x(2^(N)-1) que es
%iterante inicial
[m,n] = size(b);

%La completamos con ceros la matriz para simular la malla
Z = zeros(1,m);
B = [Z;b;Z];
Z = zeros(m+2,1);
B = [Z,B,Z];

%Construimos la matriz donde guardaremos la restriccion
U = zeros(2*m +1,2*n +1);

for i = 1 : 1 : 2*m +1
    for j = 1: 1 : 2*m +1

        if 0 == mod(j,2) && 0 == mod(i,2)
            U(i,j) = B (i/2 + 1 ,j/2 + 1);

```

```

end
if 1 == mod(j,2) && 0 == mod(i,2)
    U(i,j) = 1/2*(B(i/2 + 1 ,(j-1)/2 + 1) +
        B(i/2 + 1 ,(j+1)/2 + 1));
end
if 0 == mod(j,2) && 1 == mod(i,2)
    U(i,j) = 1/2*(B((i-1)/2 + 1 ,(j)/2 + 1) +
        B((i+1)/2 + 1 ,(j)/2 + 1));
end
if 1 == mod(j,2) && 1 == mod(i,2)
    U(i,j) = 1/4*(B((i-1)/2 + 1 ,(j-1)/2 + 1)
        + B((i+1)/2 + 1 ,(j-1)/2 + 1) +
        B((i-1)/2 + 1 ,(j+1)/2 + 1) +
        B((i+1)/2 + 1 ,(j+1)/2 + 1));
end

end
end
end

```

```

%Metodo Multigrid
function [x] = MGCYCM(k, g, u, f, v, w)

%k es el nivel de la malla
%g el tipo de ciclo
%u primera aproximacion a la solucion
%f termino independiente
%w numero de veces que aplicamos el pre-suavizado
%u numero de veces que aplicamos el post-suavizado

x = u;

for i = 1 : v
    x = GSM(x, f);
end

%Correccion en malla grosera

d = f - OperadorM(x, k);
d = FWM(d);

%aquí tenemos d en la red mas gruesa (k-1)

```

```
if k̃ = 1
    p = zeros((2^(k) - 1), (2^(k) - 1));
    for j = 1 : g
        z = MGCYCM(k-1, g, p, d, v, w);
        p = z;
    end
else
    z = -d/16;
end
z = IBM(z);
x = x + z;
for i = 1 : w
    x = GSM(x, f);
end
end
```

# Bibliografía

- [1] BAKHVALOV, N. S., *On the convergence of a relaxation method with natural constraints on the elliptic operator*, USSR Comp. Math. Math. Phys. 6, 101-135, 1966.
- [2] BRAMBLE, J. H., *Multigrid Methods*, Longman Scientific and Technical, 1995.
- [3] BRANDT, A., *Multi-level adaptive techniques (MLAT) . I. The multigrid method*, Research Rep. RC 6026, IBM T.J. , 1976.
- [4] BRANDT, A. , *Multi-level adaptive technique (MLAT) for fast numerical solution to boundary value problems. Proceedings of the 3rd International Conference on Numerical Methods in Fluid Mechanics*, Lecture Notes in Physics 18 (eds H. Cabannes and R. Temam). 82-89. Springer, Berlin, 1973.
- [5] BRANDT, A. , *Multi-level adaptive techniques (MLAT) for partial differential equations: ideas and software.*, Mathematical Software III (ed. J.R. Rice). 277-318. Academic Press, New York, 1977.
- [6] FEDORENKO, R. P., *A relaxation method for solving elliptic difference equations*, USSR Comp. Math. Math. Phys. 1(5), 1092-1096, 1962.
- [7] FEDORENKO, R. P., *The speed of convergence of one iterative process*, USSR Comp. Math. Math. Phys. 4(3), 227-235, 1964.
- [8] HACKBUSH, W., *Multi-Grid Methods and Applications*, Springer-Verlag, 1980.
- [9] TROTTENBERG, U., OOSTERLEE, C. y SCHÜLLER, A., *Multigrid*, Academic Press, 2000.
- [10] STRIKWERDA, J. C., *Finite difference schemes and partial differential equations*, Wadsworth & Brooks, 1989.
- [11] YOUNG, D. M. , *Iterative solution of large linear systems*, Academic Press, 1971.