

A.A. Samarski
Introducción
a los métodos numéricos

Editorial Mir Moscú

Introducción a los métodos numéricos

Introducción a los métodos numéricos

А. А. Самарский

Введение в численные методы

Издательство «Наука»

A. A. Samarski

Introducción a los métodos numéricos



Editorial Mir
Moscú

Traducido del ruso por el ingeniero K. P. Medkov

Impreso en la URSS

На испанском языке

Издательство «Наука». Главная редакция физико-математической литературы, 1982

Traducción al español, editorial «Mir», 1988

Contenido

Prólogo	7
Introducción	9

Capítulo I

Ecuaciones en diferencias

§ 1. Funciones reticulares	28
§ 2. Ecuaciones en diferencias	31
§ 3. Resolución de los problemas de contorno en diferencias para las ecuaciones de segundo orden	40
§ 4. Ecuaciones en diferencias como ecuaciones operacionales	45
§ 5. Principio del máximo para las ecuaciones en diferencias	65

Capítulo II

Interpolación e integración numérica

§ 1. Interpolación y aproximación de las funciones	72
§ 2. Integración numérica	82

Capítulo III

Resolución numérica de los sistemas de ecuaciones algebraicas lineales

§ 1. Sistemas de ecuaciones algebraicas lineales	100
§ 2. Métodos directos	106
§ 3. Métodos iterativos	113
§ 4. Esquema iterativo de dos capas con parámetros de Chébishev	129
§ 5. Método alternado triangular	140
§ 6. Métodos iterativos de tipo variacional	147
§ 7. Resolución de las ecuaciones no lineales	150

Capítulo IV

Métodos de diferencias de la resolución de los problemas de contorno para ecuaciones diferenciales ordinarias

§ 1.	Conceptos fundamentales de la teoría de esquemas de diferencias	158
§ 2.	Esquemas de diferencias homogéneos tripuntuales	172
§ 3.	Esquemas de diferencias conservativos	175
§ 4.	Esquemas homogéneos sobre las redes no uniformes	183
§ 5.	Métodos de construcción de los esquemas de diferencia	191

Capítulo V

Problema de Cauchy para las ecuaciones diferenciales ordinarias

§ 1.	Métodos de Runge-Kutta	200
§ 2.	Esquemas de varios pasos. Métodos de Adams	212
§ 3.	Aproximación del problema de Cauchy para un sistema de ecuaciones diferenciales lineales ordinarias de primer orden	224
§ 4.	Estabilidad del esquema de dos capas	230

Capítulo VI

Métodos de diferencias para las ecuaciones elípticas

§ 1.	Esquemas de diferencias para la ecuación de Poisson	241
§ 2.	Resolución de las ecuaciones en diferencias	252

Capítulo VII

Métodos de diferencias para resolver la ecuación de conductibilidad térmica

§ 1.	Ecuación de conductibilidad térmica con coeficientes constantes	364
§ 2.	Problemas multidimensionales de conductibilidad térmica	277
§ 3.	Esquemas económicos	285
	Anexo	295
	Bibliografía	302
	Lista de designaciones	304
	Índice alfabético	306

Prólogo

Este libro representa una introducción a la teoría de los métodos numéricos en la que se emplea un mínimo de información de tales apartados de las matemáticas como son el análisis, el álgebra lineal y la teoría de ecuaciones diferenciales. El libro ha surgido como resultado de elaboración de las conferencias dictadas por el autor durante varios años para los estudiantes de la facultad de matemática de cálculo y cibernética de la Universidad de Moscú Lomónosov.

El contenido del libro es tradicional: interpolación y aproximación, integración numérica, resolución de ecuaciones no lineales, métodos directos e iterativos de resolución de los sistemas de ecuaciones algebraicas lineales, métodos de diferencias destinados a resolver el problema de Cauchy y problemas de contorno para las ecuaciones diferenciales ordinarias.

La aspiración del autor fue hacer la exposición comprensible de la primera lectura, prestando una atención especial a los conceptos principales de la teoría de los métodos numéricos e ilustrándolos con los ejemplos más simples.

Para la resolución numérica de varios problemas de la física y de la técnica descritos por las ecuaciones de la física matemática se emplea actualmente el método de diferencias finitas. Los conceptos principales de la teoría de los métodos de diferencias (aproximación, estabilidad, convergencia) se ilustran con ejemplos de esquemas de diferencias para las ecuaciones diferenciales ordinarias. Al aproximar las ecuaciones diferenciales, obtenemos ecuaciones en diferencias que representan sistemas de ecuaciones lineales de orden superior con matrices del tipo especial (tienen muchos elementos nulos), por ejemplo, tridiagonales. Un papel de

importancia lo desempeña la elección de los métodos efectivos (directos e iterativos) para resolver los sistemas mencionados. Con este motivo en el libro se exponen los fundamentos de la teoría general de métodos iterativos. Una gran atención se ha dedicado a la cuestión de estabilidad de los cálculos en los ordenadores. En el capítulo V viene una exposición sencilla de la teoría de estabilidad del problema de Cauchy para el sistema de ecuaciones en diferencias de primer orden. Aquí se han obtenido las condiciones coincidentes de estabilidad *necesarias y suficientes de los esquemas de diferencias* y, además, se ha investigado la estabilidad asintótica de los esquemas de diferencias.

En los dos últimos capítulos del libro (VI y VII) se analizan métodos de diferencias para resolver las ecuaciones elípticas y la ecuación de conductibilidad térmica. Estos capítulos son complementarios y permiten realizar el paso a la teoría de esquemas de diferencias para las ecuaciones en derivadas parciales.

Una exposición más detallada de los apartados separados de los métodos numéricos se da en los libros: «Teoría de esquemas de diferencias» por Samarski A. A., «Métodos de resolución de las ecuaciones reticulares» por Samarski A. A., Nikoláev E. S., y otros que se indican en la lista al final del libro.

El libro está destinado a los estudiantes de los primeros años que eligen como su especialidad la matemática aplicada y la física matemática; este libro puede resultar útil también para postgraduados y colaboradores científicos que estudien los métodos numéricos.

A. A. Samarski

Introducción

La aparición y el perfeccionamiento incesante de los ordenadores de alta velocidad han conducido a una transformación auténticamente revolucionaria de la ciencia en general y de las matemáticas, en particular. Ha cambiado la tecnología de las investigaciones científicas, han aumentado inmensamente las posibilidades de los estudios teóricos, del pronóstico de procesos complejos, de la proyección de las construcciones de ingeniería. Únicamente gracias a la aplicación de la simulación matemática y de nuevos métodos numéricos destinados para los ordenadores se hizo posible resolver grandes problemas científico-técnicos tales como el dominio de la energía nuclear y la asimilación del cosmos.

El primer gran problema, el dominio de la energía nuclear, requiere que se resuelva un conjunto de problemas complejos de la física y mecánica (manejo del trabajo de la caldera nuclear, la utilización de la energía proveniente de la fisión de los núcleos de uranio, la protección de la irradiación penetrante, el enfriamiento de las paredes de reactor, el estudio de los campos térmicos y de tensiones elásticas en las paredes, la resolución de varios otros problemas). Todos estos problemas han de ser resueltos antes de que empiece a trabajar una caldera, usando para este fin la descripción matemática (un modelo) y realizando cálculos numéricos en el ordenador. El segundo gran problema consistente en la asimilación del cosmos está relacionado con la creación de aparatos voladores y la resolución para estos últimos de diferentes problemas aerodinámicos y balísticos (por ejemplo, el cálculo del movimiento de un cohete y la dirección de su vuelo). En este dominio también hay un conjunto de problemas complejos de la mecánica, física

y técnica los cuales pueden ser resueltos sólo aplicando los *métodos numéricos*.

Indiquemos un problema más planteado ante la humanidad, esto es, la búsqueda de nuevas fuentes de energía. Uno de los proyectos fundamentales para obtener energía consiste en emplear la reacción de fusión termonuclear dirigida de los núcleos de deuterio y de tritio. Los recursos de combustible termonuclear en la Tierra son prácticamente inagotables, mientras que los productos de reacción no ensucian el ambiente. No obstante, la reacción termonuclear comienza sólo en condiciones extremadas: a una altísima temperatura (decenas y centenas de millones de grados) y enorme compresión (miles de veces) del deuterio y tritio; además, se requiere mantener la sustancia combustible en dicho estado durante un período de tiempo que sea suficiente para que se desarrolle la reacción de combustión (del síntesis). La creación de las condiciones mencionadas es un problema científico-técnico que por ahora no está resuelto. Existen varios proyectos destinados a calentar, comprimir y mantener el combustible termonuclear (plasma). Al realizarlos surge una serie de cuestiones que *deben ser resueltas antes de proceder a la proyección de las instalaciones correspondientes, incluso experimentales*. Es menester estudiar ante todo el comportamiento del plasma a altas temperaturas y densidades, en campos magnéticos y, además, aclarar las condiciones bajo las cuales resulta posible la propia reacción de la síntesis termonuclear.

Las investigaciones de tal índole se efectúan a base de la descripción matemática (modelo matemático) de los procesos físicos y la resolución ulterior de problemas matemáticos correspondientes en el ordenador con ayuda de algoritmos de cálculo (computacionales).

Hoy día podemos decir que ha surgido un método nuevo para la investigación teórica de los procesos complejos que admiten la descripción matemática: se trata de un *experimento de cálculo, es decir, la investigación de los problemas científicos naturales por medio de la matemática de cálculo*. Expliquemos la esencia de este método de investigación con un ejemplo de resolución de un problema físico. Supongamos que se pide estudiar cierto proceso físico. A la investigación matemática le precede la elección de una

aproximación física, es decir, se debe determinar qué factores han de tomarse en consideración y cuáles pueden ser menospreciados. Resuelta la cuestión citada, se realiza la investigación del problema mediante un experimento de cálculo, en el que pueden distinguirse las siguientes etapas principales.

En la primera etapa se elige un modelo matemático, es decir, la descripción aproximada del proceso en forma de ecuaciones algebraicas, diferenciales o integrales. Estas ecuaciones expresan corrientemente las leyes de conservación de las magnitudes físicas principales (la energía, la cantidad de movimiento, la masa, etc.). El modelo matemático obtenido ha de ser investigado recurriendo a la teoría de ecuaciones diferenciales. Se debe establecer si el problema está planteado correctamente, si los datos de partida son suficientes y ellos no contradicen los unos a los otros, si existe la solución del problema planteado y si es única. En esta etapa se emplean los métodos de la matemática clásica. Hemos de señalar que muchos problemas físicos conducen a ciertos modelos matemáticos cuya elaboración teórica acaba de iniciarse. En la práctica nos vemos obligados a resolver problemas de la física matemática, para los cuales no existen teoremas de existencia y unicidad.

La segunda etapa del experimento de cálculo consiste en la construcción de un método numérico aproximado que se usa para resolver el problema, es decir, en la elección del algoritmo de cálculo. Por algoritmo de cálculo se entiende una sucesión de operaciones aritméticas y lógicas que ayudan a encontrar la solución del problema matemático formulado en la primera etapa. Más abajo se discutirán detalladamente las exigencias que se presentan a un algoritmo de cálculo destinado para el empleo en los ordenadores modernos. El presente libro está dedicado, en esencia, al estudio de los algoritmos de cálculo elementales.

En la tercera etapa se lleva a cabo la programación del algoritmo de cálculo para el ordenador y en la cuarta etapa, los cálculos en el ordenador. No nos detendremos en las cuestiones ligadas con la programación, organización y realización de los cálculos en el ordenador, puesto que todas estas cuestiones salen de los márgenes del libro. Notemos sólo que todas las operaciones referentes a la programación

deben estar en una relación estrecha con la elaboración de los algoritmos numéricos concretos.

En fin, a título de la quinta etapa del experimento de cálculo puede indicarse el análisis de los resultados numéricos obtenidos y la precisión posterior del modelo matemático. Puede suceder que el modelo es demasiado aproximado (el resultado de los cálculos no concuerda con el experimento físico) o bien el modelo es muy complejo y la solución puede obtenerse con una exactitud suficiente empleando modelos más simples. En este caso el trabajo se inicia desde la primera etapa, es decir, se precisa el modelo matemático y se repasan otra vez todas las etapas.

Hemos de notar que un experimento de cálculo no es, como regla, una operación sencilla de cálculo por fórmulas estándar, sino, ante todo, los cálculos de toda una serie de variantes para diferentes modelos matemáticos.

Ahora, fijemos nuestra atención en ciertas características y exigencias generales concernientes a los algoritmos de cálculo. La elaboración e investigación de los algoritmos de cálculo y la aplicación de éstos a la resolución de los problemas concretos constituyen el contenido de un gran apartado de la matemática moderna: matemática de cálculo.

La matemática de cálculo se determina en el amplio sentido de este término como un apartado de las matemáticas que incluye un conjunto de cuestiones relacionadas con el empleo de los ordenadores; en el sentido estrecho dicho apartado de las matemáticas se entiende como teoría de los métodos numéricos y de los algoritmos para resolver los problemas matemáticos planteados. En lo sucesivo la matemática de cálculo se considerará sólo en el sentido estrecho.

Hay un rasgo común para todos los métodos numéricos que consiste en reducir todo problema matemático a uno que sea de dimensión finita. Esto se consigue con mayor frecuencia discretizando el problema de partida, es decir, pasando de las funciones de un argumento continuo a las de argumento discreto. Discretizado el problema de partida, se debe construir un algoritmo de cálculo, es decir, indicar la sucesión de operaciones aritméticas y lógicas que se ejecutan en el ordenador y que proporcionan, tras un número finito de operaciones, la solución del problema discreto.

La solución obtenida del problema discreto se considera como solución aproximada del problema matemático de partida.

Al resolver problemas en el ordenador obtenemos siempre no la solución exacta del problema de partida, sino cierta solución aproximada. ¿A qué se debe el error que surge? Pueden ser indicadas tres razones principales a consecuencia de las cuales surgen errores en la resolución numérica del problema matemático de partida. Ante todo, los datos de entrada del problema de partida (condiciones iniciales y de frontera, coeficientes y segundos miembros de las ecuaciones) se dan siempre con cierta inexactitud. Un error del método numérico condicionado por la prefijación inexacta de los datos de entrada suele denominarse *error inevitable*. Luego, al sustituir el problema de partida por otro problema discreto aparece un error que se llama *error de discretización* o, de otra forma, *error del método*. Por ejemplo, sustituyendo la derivada $u'(x)$ por una razón de diferencias $(u(x + \Delta x) - u(x))/\Delta x$, cometemos un error de discretización que para $\Delta x \rightarrow 0$ tiene el orden Δx . Finalmente, el orden finito de los números que se suministran al ordenador lleva a *errores de redondeo* que pueden acumularse en el transcurso de los cálculos. Es natural exigir que los errores en la prefijación de la información inicial y el error que surge como resultado de discretización sean concordados con el error de la solución del problema discreto en el ordenador.

De lo dicho proviene que la exigencia principal que se levanta ante el algoritmo de cálculo es exactitud. Dicha exigencia quiere decir que el algoritmo de cálculo debe asegurar la solución del problema de partida con la *exactitud* prefijada $\varepsilon > 0$, realizadas un número finito $Q(\varepsilon)$ de operaciones. El algoritmo ha de ser realizable, es decir, debe proporcionar la solución del problema en tiempo de máquina admisible. Para la mayoría de los algoritmos el tiempo que se necesita para resolver el problema (volumen de los cálculos) $Q(\varepsilon)$ crece al aumentar la exactitud, es decir, cuando disminuye ε . Por supuesto, se puede prefijar ε tan pequeño que el tiempo de resolución del problema se hará inadmisiblemente grande. Resulta importante conocer que el algoritmo da en principio una posibilidad de obtener la solución del problema con cualquier exactitud. Sin embargo, en la

práctica la magnitud de ε se elige tomando en consideración una posibilidad de realizar el algoritmo en el ordenador dado. Para cualquier problema, algoritmo y ordenador existe un valor individual de ε .

Es natural de aspirar a que el número de operaciones (y, de este modo, el tiempo de máquina para la resolución del problema) $Q(\varepsilon)$ sea mínimo para el problema dado. Para cualquier problema se pueden ofrecer varios algoritmos que proporcionen (para $\varepsilon \rightarrow 0$) una exactitud $\varepsilon > 0$ igual en orden, pero con diferente número de operaciones $Q(\varepsilon)$. Entre estos algoritmos (de los cuales suele decirse que ellos son equivalentes según el orden de exactitud) se debe elegir uno que proporcione la solución con un gasto mínimo de tiempo de máquina (número de operaciones $Q(\varepsilon)$). Tales algoritmos se denominarán *económicos*.

He aquí una exigencia más que ha de ser satisfecha por el algoritmo de cálculo, es decir, el requisito de que no haya parada de emergencia (de indisponibilidad) del ordenador en el proceso de los cálculos.

Es necesario tener en cuenta que todo ordenador opera con números que tienen una cantidad finita de cifras significativas y que pertenecen (en módulo) no a todo el eje numérico, sino a cierto intervalo (M_0, M_∞) , $M_0 > 0$, $M_\infty < \infty$, donde M_0 es un cero de máquina y M_∞ , un infinito de máquina. Si la condición $|M| < M_\infty$ no se cumple en el proceso de los cálculos, ocurre una parada de emergencia del ordenador («*parem*»), a consecuencia de que queda rellena la red de órdenes y los cálculos se dan por terminado. La posibilidad de una *parem* depende tanto del algoritmo como del problema de partida.

Si la solución del problema de partida se expresa en términos de números muy grandes (muy pequeños) $|M| > M_\infty$ ($|M| < M_0$), entonces, como regla, variando la escala, el problema puede ser reducido a una forma que contiene sólo las magnitudes pertenecientes (en módulo) al intervalo prefijado (M_0, M_∞) . La posibilidad de la *parem* se elimina frecuentemente cambiando el orden de operaciones. Expliquémoslo con un ejemplo sencillo.

EJEMPLO. Sea $M_\infty = 10^p$, $M_0 = 10^{-p}$, $p = 2^n$, n es un número entero. Se pide calcular el producto de los números $10^{p/2}$, $10^{p/4}$, $10^{-p/2}$, $10^{3p/4}$, $10^{-3p/4}$.

1^{er} MÉTODO. Fijemos los números en el orden decreciente:

$$\begin{aligned} q_1 &= 10^{3P/4}, & q_2 &= 10^{P/2}, & q_3 &= 10^{P/4}, \\ q_4 &= 10^{-P/2} & q_5 &= 10^{-3P/4} \end{aligned}$$

y formemos los productos $S_{k+1} = S_k q_{k+1}$, $S_1 = q_1$. En este caso, ya en el primer paso tendrá lugar una parem, puesto que $S_2 = q_1 q_2 = 10^{5P/4} > M_\infty$.

2^{do} MÉTODO. Fijemos los números en el orden creciente:

$$\begin{aligned} q_1 &= 10^{-3P/4}, & q_2 &= 10^{-P/2}, & q_3 &= 10^{P/4}, \\ q_4 &= 10^{P/2}, & q_5 &= 10^{3P/4}. \end{aligned}$$

En este caso obtendremos en el primer paso

$$S_2 = q_1 q_2 = 10^{-5P/4} < M_0,$$

es decir, S_2 es un cero de máquina; todos los productos sucesivos S_3, S_4, S_5 son también nulos; de este modo aquí ocurre una pérdida total de exactitud.

3^{er} MÉTODO. Mezclemos estos números suponiendo $q_1 = 10^{-2P/4}$, $q_2 = 10^{P/2}$, $q_3 = 10^{3P/4}$, $q_4 = 10^{-P/2}$, $q_5 = 10^{P/4}$. Entonces hallaremos sucesivamente:

$$\begin{aligned} S_2 &= q_1 q_2 = 10^{-P/4}, & S_3 &= S_2 q_3 = 10^{P/2}, \\ S_4 &= S_3 q_4 = 10^0, & S_5 &= S_4 q_5 = 10^{P/4}, \end{aligned}$$

es decir, en el proceso de los cálculos no aparecen números superiores a $10^{P/2}$ e inferiores a $10^{-P/4}$. Tal algoritmo está privado de parem. En el cap. III nos encontraremos con un método iterativo de resolución de sistemas de ecuaciones algebraicas lineales que puede realizarse con una parem y sin ésta, según sea la forma de numeración de los parámetros que determina la sucesión de los cálculos.

En cada etapa de los cálculos surgen errores de redondeo. Estos errores de redondeo pueden crecer o ir disminuyendo, en dependencia del algoritmo.

Si, en el transcurso de los cálculos, la magnitud de los errores de redondeo crece indefinidamente, el algoritmo se llamará *inestable* (desde el punto de vista de cálculos). En cambio, si los errores de redondeo no se acumulan, el algoritmo será estable.

EJEMPLOS. 1. Supongamos que se pide hallar y_i ($0 < i \leq i_0$) según la fórmula $y_{i+1} = y_i + d$ ($i \geq 0$) para y_0 , d prefijados. Supongamos, además, que al calcular y_i se ha introducido un error (por ejemplo, un error de redondeo) cuya magnitud es δ_i , es decir, en lugar del valor exacto de y_i tenemos un valor aproximado $\tilde{y}_i = y_i + \delta_i$. Entonces, en vez del valor exacto de y_{i+1} obtendremos el valor aproximado $\tilde{y}_{i+1} = (\tilde{y}_i + d) = y_{i+1} + \delta_i$. De este modo, un error cometido en cualquier paso intermedio no aumenta en el proceso de los cálculos. El algoritmo es estable.

2. Examinemos la ecuación $y_{i+1} = qy_i$ ($i \geq 0$, y_0 y q están prefijados). Supongamos, al igual que en el ejemplo 1, se ha obtenido, en lugar de y_i , el valor $\tilde{y}_i = y_i + \delta_i$. Entonces, en lugar de y_{i+1} obtendremos un valor aproximado

$$\tilde{y}_{i+1} = q(y_i + \delta_i) = y_{i+1} + q\delta_i.$$

De aquí se ve que el error $\delta_{i+1} = \tilde{y}_{i+1} - y_{i+1}$, que surge al calcular y_{i+1} , está ligado con el error δ_i mediante una ecuación

$$\delta_{i+1} = q\delta_i, \quad i = 0, 1, 2, \dots$$

Por consiguiente, si $|q| > 1$, el valor absoluto del error crecerá en el proceso de los cálculos (el algoritmo es inestable). Si, en cambio, $|q| \leq 1$, entonces el error no aumenta, es decir, el algoritmo es estable. La inestabilidad se liga corrientemente con la propiedad de crecimiento exponencial del error de redondeo. Si el error de redondeo crece según la ley potencial al pasar de una operación a la otra («de paso a paso»), el algoritmo se considera *convencionalmente estable* (estable con ciertas restricciones que se imponen sobre el volumen de cálculos y la exactitud requerida). El proceso de los cálculos puede interpretarse así: al pasar de un paso a otro tiene lugar una alteración (a cuenta de los errores de redondeo) de las últimas cifras significativas («una onda del error de redondeo» se mueve de derecha a izquierda, partiendo de las últimas cifras significativas). Nuestra tarea consiste en conservar justas unas cuantas primeras cifras significativas (4—5 signos) y por esta razón los cálculos deben darse por terminado antes de que «la

onda del error de redondeo alcance dichas cifras. Si el error de redondeo ε_0 crece de un paso al otro según la ley exponencial, esto conduce, como regla, a una parem en cierta etapa intermedia de los cálculos, si (lo mismo que en el ejemplo 2) $|q| \varepsilon_0 \geq M_\infty$.

Si $M_\infty = 10^p$, $\varepsilon_0 = 10^{-k}$, la parem llega cuando $i_0 > (p + k_0)/\lg |q|$. Otra cosa ocurre cuando el error de redondeo crece según la ley potencial. Sea $|\delta y_i| \approx i^n \varepsilon_0$ ($n \geq 1$); entonces, la parem tiene lugar para $i_0^n \varepsilon_0 \geq M_\infty$, es decir, para $i_0 \geq \left(\frac{1}{\varepsilon_0} M_\infty\right)^{1/n} = 10^{(p+k_0)/n}$.

De aquí se ve que para $n = 1$ la parem no tendrá lugar en virtud de una restricción evidente $i < M_\infty = 10^p$. La desigualdad $|\delta y_i| \leq \varepsilon$, donde $\varepsilon = 10^{-k}$ es la exactitud prefijada, se verifica para $i \leq \left(\frac{\varepsilon}{\varepsilon_0}\right)^{1/n} = 10^{(k_0-k)/n} = i_0$. Si están prefijados ε y ε_0 , esta desigualdad significa una restricción para el número de ecuaciones $i \leq i_0$. Por ejemplo, para $k_0 = 12$, $k = 6$, tenemos $i \leq 10^{6/n}$, de suerte que $i \leq 10^3$ para $n = 2$. Está claro que puede elegirse tal n grande que el número admisible de ecuaciones i_0 sea muy pequeño. Sin embargo, en la práctica se encuentran corrientemente casos de n pequeño (por ejemplo, para el método de factorización (§ 3 cap. I) $n = 2$, es decir, el error se acumula según la ley cuadrática a medida que crece el número de ecuaciones).

Al resolver un problema (cualquiera que sea) es necesario conocer ciertos datos de entrada (de partida): datos iniciales, valores de frontera de la función buscada, coeficientes y el segundo miembro de la ecuación, etc.

Para todo problema se buscan respuestas a las preguntas de un mismo género: si existe la solución del problema, si será única y cómo depende la solución de los datos de entrada. Son posibles dos casos:

El problema está correctamente planteado (es correcto); esto quiere decir que: 1) el problema es resoluble para cualesquiera datos de entrada admisibles; 2) se tiene una única solución; 3) la solución del problema depende continuamente de los datos de entrada (a una variación pequeña de los datos de entrada le corresponde una variación pequeña de la solución), en otras palabras, el problema es estable.

El problema no está correctamente planteado (no es correcto), si la solución de éste es inestable respecto a los datos de entrada (a una variación pequeña de los datos de entrada le puede corresponder una variación grande de la solución).

Como ejemplo de un problema correcto puede servir el problema de integración y como ejemplo de un problema no correcto, el problema de diferenciación.

EJEMPLOS. 1. PROBLEMA DE INTEGRACION. Sea dada una función $f(x)$; hállese la integral

$$J = \int_0^1 f(x) dx.$$

Sustituyamos f por \tilde{f} y veamos $\tilde{J} = \int_0^1 \tilde{f}(x) dx$ y la diferen-

cia $\delta J = \tilde{J} - J = \int_0^1 \delta f dx$ ($\delta f = \tilde{f}(x) - f(x)$). De aquí se ve que

$|\delta J| \leq \max_{0 \leq x \leq 1} |\delta f(x)|$, $|\delta J| \leq \epsilon$, si $|\delta f| \leq \epsilon$, es decir J depende

continuamente de f . Con el fin de calcular la integral J hagamos uso de la fórmula de cuadratura:

$$J_N = \sum_{h=1}^N c_h f(x_h), \quad c_h > 0, \quad \sum_{h=1}^N c_h = 1.$$

Al repetir los razonamientos aducidos más arriba, llegamos a que

$$\delta J_N = \tilde{J}_N - J_N = \sum_{h=1}^N c_h (\tilde{f}_h - f_h) = \sum_{h=1}^N c_h \delta f_h,$$

$$|\delta J_N| \leq \sum_{h=1}^N c_h \max_{1 \leq h \leq N} |\delta f_h| = \max_{1 \leq h \leq N} |\delta f_h|.$$

De este modo, el problema de cálculo de una integral por la fórmula de cuadratura es correcto.

2. PROBLEMA DE DIFERENCIACION. El problema de diferenciación de una función $u(x)$ definida aproximadamente no es correcto.

En efecto, sea $\tilde{u}(x) = u(x) + \frac{1}{N} \sin N^2 x$, donde N es suficientemente grande. Entonces, en la métrica C (en cierto segmento $0 \leq x \leq \delta$ ($\delta > \pi/N^2$)) tenemos $\|\delta u\|_C = \|\tilde{u} - u\|_C = 1/N \leq \varepsilon$ para $N \geq 1/\varepsilon$. Para el error de las derivadas $\delta u' = \tilde{u}' - u' = N \cos N^2 x$ tenemos $\|\delta u'\|_C = N \geq 1/\varepsilon$. De este modo, a una variación pequeña $O(\varepsilon)$ en C de la función $u(x)$ le corresponde la variación grande $O(1/\varepsilon)$ en C de su derivada.

Por eso la diferenciación numérica tampoco es correcta. Para encontrar un valor aproximado de una derivada según la fórmula de la derivada de diferencias con cierta exactitud $\varepsilon > 0$ a condición de que la función viene definida con un error δ_i ($|\delta_i| \leq \delta_0$), hace falta que se cumplan las condiciones de concordancia entre ε , δ_0 y el paso h de la red, por ejemplo, del tipo $\varepsilon \geq k \sqrt{\delta_0}$ ($k = \text{const} > 0$ no depende de h , δ_0), con la particularidad de que el paso de la red está acotado tanto inferiormente como superiormente. De este modo, la exactitud alcanzable de la diferenciación numérica está limitada por la exactitud con la que viene definida la propia función.

En este libro se estudian sólo problemas correctos y métodos numéricos correctos destinados para aplicarlos en el ordenador.

Los métodos numéricos dan la solución aproximada del problema. Esto significa que en lugar de la solución exacta u (de una función o de una funcional) de algún problema encontramos una solución y de otro problema, próxima en cierto sentido (según la norma, por ejemplo) a la buscada. De acuerdo con lo dicho, la idea principal de todos los métodos consiste en discretizar o aproximar (sustitución, aproximación) el problema de partida a algún otro que sea más cómodo para resolverlo en un ordenador, con la particularidad de que la resolución del problema que aproxima depende de ciertos parámetros que, siendo manejados de una manera adecuada, permiten determinar la solución con una exactitud requerida. Por ejemplo, en el problema de integración numérica los nodos y los pesos de la fórmula de cuadratura representan precisamente los parámetros de este tipo. Luego, la solución de un problema discreto es elemento de

un espacio de dimensión finita. Expliquemos esto más detalladamente.

Veamos, por ejemplo, la discretización de un espacio $H = \{f(x)\}$ de funciones $f(x)$ de argumento continuo $x \in [a, b]$. Introduzcamos en un segmento $a \leq x \leq b$ un conjunto finito de puntos $\omega = \{x_i, i = 0, 1, \dots, N, x_0 = a, x_N = b, x_i < x_{i+1}\}$ el cual se llamará *red*. Los puntos x_i se denominarán *nodos* de la red ω . Si la distancia $h_i = x_i - x_{i-1}$ entre los nodos vecinos es constante (no depende de i), $h_i = h$ para todos los $i = 1, 2, \dots, N$, entonces la red se denomina *uniforme* (de paso h); en caso contrario se denomina *no uniforme*. En lugar de la función $f(x)$ definida para cualesquiera $x \in [a, b]$ consideraremos una *función reticular* $y_i = f(x_i)$ de argumento i ($i = 0, 1, \dots, N$), que es un número entero, o bien de nodo x_i de la red ω , y sustituiremos $H = \{f(x), x \in [a, b]\}$ por un espacio de dimensión finita (de dimensión $N + 1$) $H_{N+1} = \{y_i, 0 \leq i \leq N\}$ de funciones reticulares. Es evidente, que la función reticular $y_i = f(x_i)$ puede considerarse como un vector $y = (y_0, y_1, \dots, y_N)$.

Podemos discretizar también el espacio de funciones $f(x)$ de varias variables, si $x = (x_1, x_2, \dots, x_p)$ es un punto del espacio euclídeo p -dimensional ($p > 1$). Por ejemplo, en un plano (x_1, x_2) se puede introducir una red $\omega = \{x_i = (i_1 h_1, i_2 h_2), i_1, i_2 = 0, \pm 1, \pm 2, \dots\}$ como conjunto de puntos (nodos) de intersección de las rectas perpendiculares $x_1^{(i_1)} = i_1 h_1, x_2^{(i_2)} = i_2 h_2, h_1 > 0, h_2 > 0, i_1, i_2 = 0, \pm 1, \pm 2, \dots$, donde h_1 y h_2 son los pasos de la red según las direcciones de x_1 y x_2 , respectivamente. La red ω es, evidentemente, uniforme según cada una de las variables por separado. En lugar de la función $f(x) = f(x_1, x_2)$ analizaremos una función reticular

$$y_{i_1, i_2} = f(i_1 h_1, i_2 h_2).$$

Si la red ω contiene sólo los nodos pertenecientes al rectángulo $(0 \leq x_1 \leq l_1, 0 \leq x_2 \leq l_2)$ de modo que $h_1 = l_1/N_1, h_2 = l_2/N_2$, entonces la red cuenta con un número finito $N = (N_1 + 1)(N_2 + 1)$ de nodos, mientras que el espacio H_N de funciones reticulares $y_i = y_{i_1, i_2}$ es de dimensión finita.

En todos los casos consideramos sólo un espacio de funciones reticulares cuya dimensión es finita. Al sustituir el espacio $H = \{f(x)\}$ de funciones de argumento continuo por el espacio H_N de funciones reticulares y el problema de partida, por su aproximación discreta, debemos estar seguros de que nos aproximamos mejor a la solución del problema de partida aumentando el número de nodos. La estimación de la calidad de aproximación y la elección del método de aproximación constituyen la tarea principal de la teoría de los métodos numéricos.

El contenido fundamental del libro está relacionado de una u otra manera con la aplicación de los métodos de diferencias para solucionar ecuaciones diferenciales. Destaquemos dos momentos de importancia:

- obtención de la aproximación discreta (de diferencias) de las ecuaciones diferenciales e investigación de las ecuaciones en diferencias que aparecen en este caso;
- resolución de las ecuaciones en diferencias.

Al obtener una aproximación discreta (esquema de diferencias) un papel importante lo desempeña el requisito general referente a que el esquema de diferencias aproxime al máximo (simule) las propiedades principales de la ecuación diferencial inicial. Tales esquemas de diferencias pueden obtenerse, por ejemplo, con ayuda de los principios variacionales y las relaciones integrales (véase el cap. IV). La estimación de la exactitud de un esquema de diferencias se reduce al estudio del error de aproximación y de la estabilidad del esquema. El estudio de la estabilidad es una cuestión central de la teoría de los métodos numéricos y a ella se le presta gran atención en el presente libro. Los algoritmos de los problemas complejos se pueden representar como una sucesión (cadena) de algoritmos simples (módulos). Por eso muchos problemas de principio de la teoría de los métodos numéricos pueden aclararse con algoritmos simples.

En el primer capítulo se examinan las ecuaciones en diferencias unidimensionales (que dependen de un argumento entero). Nos limitamos al estudio de las ecuaciones en diferencias de primero y segundo órdenes. Las ecuaciones en diferencias de segundo orden representan un sistema de ecuaciones algebraicas lineales con matriz tridiagonal. Con

el objeto de resolver los problemas de contorno para estas ecuaciones se emplea el así llamado método de factorización. En el primer capítulo se dan, a título de un material de información, algunos conocimientos sobre los operadores lineales en un espacio de dimensión finita. Posteriormente se investigan las propiedades de los operadores de diferencias en su calidad de operadores lineales en un espacio de dimensión finita provisto de un producto escalar. En este caso se emplea un aparato matemático más sencillo, es decir, las fórmulas para la diferenciación de diferencias de un producto y para la adición por partes.

En el segundo capítulo se expone el material tradicional del análisis numérico: la interpolación, la aproximación media cuadrática y la integración numérica.

Al aproximar las ecuaciones diferenciales en una red se obtienen ecuaciones en diferencias que representan un sistema de ecuaciones algebraicas lineales de orden superior (igual al número de nodos de la red) con una matriz especial (enrarecida, es decir, una matriz que tiene muchos elementos nulos). Un ejemplo más simple de tal matriz (una matriz tridiagonal) fue indicado anteriormente.

En el tercer capítulo se exponen los métodos numéricos de resolución de las ecuaciones algebraicas lineales

$$\sum_{j=1}^N a_{ij}u^j = f^i \quad i = 1, 2, \dots, N, \quad (1)$$

las cuales pueden ser escritas en forma matricial

$$Au = f, \quad (2)$$

donde $A = (a_{ij})$ es una matriz cuadrada de dimensión $N \times N$, $u = (u^1, u^2, \dots, u^N)$ es el vector buscado y $f = (f^1, f^2, \dots, f^N)$, el vector prefijado (el segundo miembro).

Para resolver los sistemas de ecuaciones se usan métodos directos e iterativos.

En el § 2 del cap. III se analizan el método de eliminación de Gauss y el de raíz cuadrada que representan métodos directos los cuales requieren $O(N^3)$ operaciones aritméticas para resolver el sistema.

Al estudiar los métodos iterativos, resulta cómodo interpretar el sistema de ecuaciones algebraicas lineales (2)

como ecuación operacional de primera especie con un operador que actúa en el espacio N -dimensional H_N ($A: H_N \rightarrow H_N$), $u, f \in H_N$. Para subrayar la equivalencia existente entre las formas de escritura matricial y operacional, la matriz y el operador correspondiente se designarán con una misma letra A .

En la teoría de los métodos iterativos (de un paso o de dos capas) es de mucha importancia la forma canónica del esquema iterativo

$$B \frac{y_{k+1} - y_k}{\tau_{k+1}} + Ay_k = f, \quad k=0, 1, \dots \text{ para todo } (3)$$

$$y_0 \in H_N,$$

donde $A, B: H_N \rightarrow H_N$, $\{\tau_k\}$ son parámetros iterativos.

En cualquier caso se supone que el operador A es autoconjugado y definido positivo ($A = A^* > 0$). Está demostrado el teorema general sobre la convergencia del método estacionario con $\tau_k = \tau = \text{const}$. Como condición suficiente de la convergencia interviene una desigualdad

$$(By, y) > \frac{\tau}{2} (Ay, y) \text{ para todo } y \in H, \quad (4)$$

donde $B \neq B^*$ es, en el caso general, un operador no autoconjugado. De aquí se desprende la convergencia del método de iteración simple, del método de Seidel y del de relajación superior.

Si se conocen unas constantes $\gamma_1 > 0$, $\gamma_2 > \gamma_1$, tales que

$$\gamma_1 (Bx, x) \leq (Ax, x) \leq \gamma_2 (Bx, x) \text{ para cualquier } x \in H_N, \quad (5)$$

donde $B = B^* > 0$, entonces podemos encontrar una totalidad optimal de parámetros de Chébishev $\{\tau_k^*\}$, con los cuales el proceso de cálculo es estable y se realiza sin parem.

Se examina el método universal alternado triangular con la totalidad $\{\tau_k^\alpha\}$ y un operador

$$B = (D + \omega A_1) D^{-1} (D + \omega A_2), \quad (6)$$

donde $D = D^* > 0$, $A_1^* = A_2$, $A_1 + A_2 = A$, las matrices A_1 y A_2 son triangulares. Hemos obtenido la fórmula para

el parámetro ω . El algoritmo para este método es muy simple. En todo caso se dan a conocer las fórmulas para el número de iteraciones con las cuales se alcanza la exactitud requerida. Los diferentes métodos fueron comparados a base de un problema modelo para la ecuación en diferencias de segundo orden $y_{i-1} - 2y_i + y_{i+1} = -h^2 f_i$, $i = 1, 2, \dots, N-1$, $y_0 = y_N = 0$, $h = 1/N$, que corresponde al problema de contorno $u''(x) = -f(x)$ ($0 < x < 1$), $u(0) = u(1) = 0$. Esta ecuación es un análogo unidimensional de la ecuación de Laplace. Por cuanto el número de iteraciones no depende prácticamente del número de mediciones, entonces en el proceso de comparación podemos limitarnos a este problema unidimensional. El método alternativo triangular exige $O\left(\frac{1}{\sqrt{h}} \ln \frac{1}{\varepsilon}\right)$ iteraciones, donde $\varepsilon > 0$ es la exactitud prefijada.

Ha de ser notado que en el cap. III, en forma lo suficientemente completa, está expuesta de hecho, con ayuda de los medios matemáticos más sencillos, la teoría general de los métodos iterativos para resolver la ecuación $Au = f$ ($A = A^* > 0$).

Los conceptos fundamentales de la teoría de esquemas de diferencias —error de aproximación, estabilidad, convergencia y exactitud— se exponen a base de ejemplos de los problemas de contorno y del problema de Cauchy para ecuaciones diferenciales ordinarias (cap. IV y cap. V). En el cap. IV se analizan esquemas de diferencias tripuntuales para una ecuación diferencial ordinaria de segundo orden

$$\frac{d}{dx} \left(k(x) \frac{du}{dx} \right) - q(x) u = -f(x), \quad 0 < x < 1,$$

$$u(0) = u_1, \quad u(1) = u_2, \quad k(x) > 0, \quad q(x) \geq 0. \quad (7)$$

Se han investigado las cuestiones referentes a la velocidad de convergencia (orden de exactitud) de los esquemas homogéneos de diferencias sobre las redes no uniformes y para el caso de coeficientes discontinuos. Esto ha exigido que se obtengan estimaciones apriorísticas bastante finas que expresen la estabilidad del esquema de diferencias respecto al segundo miembro.

Para obtener los esquemas de diferencias pueden utilizarse algunos métodos más diferentes: integral de interpolación, de aproximación de la funcional cuadrática, los de Ritz y Galerkin (§ 5, cap. IV).

Con el fin de resolver el problema de Cauchy para la ecuación de primer orden

$$\frac{du}{dt} = f(t, u), \quad t > 0, \quad u(0) = u_0 \quad (8)$$

se emplean los métodos de Runge—Kutta y de Adams expuestos en el cap. V. Estos métodos son también aplicables para un sistema de ecuaciones en que f y u son vectores.

Una atención especial en el cap. V se presta al problema de Cauchy para el sistema de ecuaciones lineales

$$\frac{du}{dt} + Au = f(t), \quad t > 0, \quad u(0) = u_0, \quad (9)$$

donde $A = (a_{ij})$ es una matriz cuadrada $N \times N$, $u(t) = (u^1, u^2, \dots, u^N)$, $f(t) = (f^1, f^2, \dots, f^N)$ es una función vectorial de N -ésima dimensión.

Tal problema surge, en particular, si en la ecuación de conductibilidad térmica

$$\frac{\partial u}{\partial t} = \Delta u + f(x, t), \quad \Delta u = \frac{\partial^2 u}{\partial x_1^2} + \frac{\partial^2 u}{\partial x_2^2}, \quad x = (x_1, x_2) \quad (10)$$

sustituimos el operador de Laplace Δu por el operador de diferencias correspondiente. Entonces, (9) puede interpretarse como un método de las rectas para la ecuación de conductibilidad térmica (10). Empleando para resolver este problema algún esquema de un paso, llegamos a un esquema operacional de diferencias de dos capas de forma general, el cual se escribe en la forma canónica

$$B \frac{y_{k+1} - y_k}{\tau} + Ay_k = \varphi_k, \quad k = 0, 1, \dots, \\ \text{para todo } y_0 \in H_{Nz} \quad (11)$$

donde $A, B: H_N \rightarrow H_N$ son los operadores lineales, τ es el paso de la red según t .

Se ha demostrado que la condición necesaria y suficiente de estabilidad del esquema tiene por expresión

$$B \geq \frac{\tau}{2} A, \text{ o bien } (Bx, x) \geq \frac{\tau}{2} (Ax, x)$$

para todo $x \in H_N$. (12)

Este es el teorema fundamental de la teoría general de estabilidad de esquemas operacionales de diferencias (véase «Teoría de los esquemas de diferencias» por Samarski A.A.) aplicable en la investigación de la estabilidad de esquemas de diferencias para las ecuaciones con derivadas parciales de la física matemática (véase el cap. VII). En realidad, en el § 4 están expuestos los fundamentos de la teoría general de estabilidad de los esquemas de diferencias, incluida la estabilidad asintótica.

Los conocimientos dados a conocer en los capítulos III, IV y V permiten pasar sin dificultad alguna al estudio de la teoría de los métodos de diferencias para resolver ecuaciones en derivadas parciales. En el cap. VI este estudio se ha realizado para esquemas de diferencias que aproximan la ecuación de Poisson y las ecuaciones elípticas en un rectángulo con condiciones de contorno de primera especie. Aquí están analizadas tanto las cuestiones de convergencia como los métodos de resolución de las ecuaciones en diferencias.

La teoría general de estabilidad de los esquemas de diferencias de dos capas (cap. V) simplifica la exposición de los métodos de diferencias para la ecuación de conductibilidad térmica con coeficientes constantes y variables realizada en el cap. VII. En el mismo capítulo se analizan también los esquemas económicos (de direcciones variables, de fisión, etc.) para los problemas multidimensionales, como también el principio general de la aproximación sumaria el que permite efectuar la partición de los problemas complejos en una sucesión de problemas más sencillos y debido a ello simplificar considerablemente la resolución de los problemas multidimensionales de la física matemática.

Se debe observar que el contenido principal de este libro se expone desde un punto de vista único. El carácter único se logra debido a que los esquemas de diferencias se

tratan como ecuaciones operacionales u operacionales de diferencias con operadores que actúan en un espacio de dimensión finita dotado de un producto escalar. Al construir la teoría de los métodos iterativos y la de estabilidad de los esquemas de diferencias se emplean las propiedades más simples de los operadores (de las matrices): el carácter constante de los signos, la autoconjugación, ciertas propiedades de los valores propios y de los vectores propios; no se hacen ningunas suposiciones referentes a la estructura de los operadores. Todas las condiciones de la teoría resultaron ser muy cómodas para la comprobación en el caso de esquemas concretos de diferencias. El material expuesto en los cap. VI y VII puede servir para un estudio más completo de la teoría la cual se da en los libros [6, 9].

Capítulo I

Ecuaciones en diferencias

En el presente capítulo se estudian funciones reticulares, cuyo argumento es un número entero, y además ecuaciones en diferencias de segundo orden. Se da a conocer un aparato matemático más simple para el estudio de las funciones reticulares y de los operadores de diferencias. Para resolver las ecuaciones en diferencias de segundo orden se emplea el método de eliminación llamado método de factorización.

§ 1. Funciones reticulares

1. **Funciones reticulares y operaciones sobre ellas.** Ya se ha mencionado que en los métodos aproximados las funciones de un argumento continuo se sustituyen habitualmente por las de argumento discreto, esto es, por las funciones reticulares. La *función reticular* puede, pues, considerarse como una función cuyo argumento es un número entero:

$$y(i) = y_i, \quad i = 0, \pm 1, \pm 2, \dots$$

Podemos introducir para $y(i)$ las operaciones que representan un análogo discreto (de diferencias) de las operaciones de diferenciación e integración.

El análogo de la primera derivada lo constituyen las diferencias de *primer orden*:

$$\Delta y_i = y_{i+1} - y_i, \text{ la diferencia derecha;}$$

$$\nabla y_i = y_i - y_{i-1}, \text{ la diferencia izquierda;}$$

$$\delta y_i = \frac{1}{2} (\Delta y_i + \nabla y_i) = \frac{1}{2} (y_{i+1} - y_{i-1}), \text{ la diferencia central;}$$

resulta fácil notar en este caso que $\Delta y_i = \nabla y_{i+1}$.

Ahora podemos escribir las diferencias de *segundo orden*:

$$\begin{aligned}\Delta^2 y_i &= \Delta (\Delta y_i) = \Delta (y_{i+1} - y_i) = y_{i+2} - 2y_{i+1} + y_i, \\ \Delta \nabla y_i &= \Delta (y_i - y_{i-1}) = (y_{i+1} - y_i) - (y_i - y_{i-1}) = \\ &= y_{i+1} - 2y_i + y_{i-1},\end{aligned}$$

de modo que

$$\Delta^2 y_i = \Delta \nabla y_{i+1}.$$

Análogamente se define la diferencia de *m-ésimo orden*:

$$\Delta^m y_i = \Delta (\Delta^{m-1} y_i),$$

que contiene los valores de $y_i, y_{i+1}, \dots, y_{i+m}$. Es evidente que

$$\sum_{j=h}^i \Delta y_j = y_{i+1} - y_h, \quad \sum_{j=h}^i \nabla y_j = y_i - y_{h-1}.$$

2. Análogos en diferencias de las fórmulas de diferenciación de un producto y de integración por partes. Sean y_i, v_i las funciones arbitrarias cuyo argumento es un número entero. En este caso serán válidas las fórmulas

$$\Delta (y_i v_i) = y_i \Delta v_i + v_{i+1} \Delta y_i = y_{i+1} \Delta v_i + v_i \Delta y_i, \quad (1)$$

$$\nabla (y_i v_i) = y_{i-1} \nabla v_i + v_i \nabla y_i = y_i \nabla v_i + v_{i-1} \nabla y_i, \quad (2)$$

que se comprueban inmediatamente. Por ejemplo,

$$\Delta (y_i v_i) = y_{i+1} v_{i+1} - y_i v_i;$$

$$\begin{aligned}y_i \Delta v_i + v_{i+1} \Delta y_i &= y_i (v_{i+1} - v_i) + v_{i+1} (y_{i+1} - y_i) = \\ &= y_{i+1} v_{i+1} - y_i v_i = \Delta (y_i v_i).\end{aligned}$$

Al deducir la fórmula para $\nabla (y_i v_i)$ es suficiente tomar en consideración que $\nabla (y_i v_i) = \Delta (y_{i-1} v_{i-1})$.

Las fórmulas (1), (2) representan los análogos de la fórmula de diferenciación del producto $(y(x)v(x))' = yv' + vy'$.

Como análogo de la fórmula de integración por partes interviene la fórmula de sumación por partes:

$$\sum_{i=0}^{N-1} y_i \Delta v_i = - \sum_{i=1}^N v_i \nabla y_i + (yv)_N - (yv)_0, \quad (3)$$

la cual se anota también en la forma

$$\sum_{i=1}^{N-1} y_i \Delta v_i = - \sum_{i=1}^{N-1} v_i \nabla y_i + y_{N-1} v_N - y_0 v_1. \quad (4)$$

Para deducir la fórmula (3) hagamos uso de la fórmula (1); tenemos

$$y_i \Delta v_i = \Delta (y_i v_i) - v_{i+1} \Delta y_i = \Delta (y_i v_i) - v_{i+1} \nabla y_{i+1},$$

puesto que $\Delta y_i = \nabla y_{i+1}$; de aquí obtenemos

$$\begin{aligned} \sum_{i=0}^{N-1} y_i \Delta v_i + \sum_{i=1}^N v_i \nabla y_i &= \\ &= \sum_{i=0}^{N-1} \Delta (y_i v_i) - \sum_{i=0}^{N-1} v_{i+1} \nabla y_{i+1} + \sum_{i=1}^N v_i \nabla y_i = \\ &= y_N v_N - y_0 v_0 - \sum_{i=1}^N v_i \nabla y_i + \sum_{i=1}^N v_i \nabla y_i = (y v)_N - (y v)_0. \end{aligned}$$

Si $y_0 = 0$, $y_N = 0$, entonces $\sum_{i=0}^{N-1} y_i \Delta v_i = - \sum_{i=1}^N v_i \nabla y_i$.

La fórmula de sumación por partes puede emplearse para calcular sumas.

EJEMPLOS. 1. Calcúlese la suma $S_N = \sum_{i=1}^N i 2^i$. Ponemos $v_i = i$, $\nabla y_i = 2^i$, de suerte que

$$y_i = y_{i-1} + 2^i = y_0 + \sum_{j=1}^i 2^j = y_0 + 2^{i+1} - 2.$$

Elijamos $y_0 = 2 - 2^{N+1}$; entonces $y_N = 0$. Como $v_0 = 0$ y $\Delta v_i = 1$, de (3) se infiere

$$\begin{aligned} S_N = \sum_{i=1}^N v_i \nabla y_i &= - \sum_{i=0}^{N-1} v_i \Delta v_i = - \sum_{i=0}^{N-1} y_i = \\ &= -N(y_0 - 2) - \sum_{i=0}^{N-1} 2^{i+1} = N 2^{N+1} - (2^{N+1} - 2), \end{aligned}$$

de manera que $S_N = (N-1) 2^{N+1} + 2$.

2. Calcúlese $S_N = \sum_{i=1}^N i(i-1) = \sum_{i=1}^{N-1} i(i+1)$. Ponemos

$y_i = i$, $\nabla v_i = i + 1$. En este caso $v_{i+1} = v_i + (i + 1) = v_1 + (2 + 3 + \dots + (i + 1)) = (v_1 - 1) + (i + 1) \times (i + 2)/2$, $v_i = v_1 - 1 + i(i + 1)/2$. Elijamos v_1 de la condición $v_N = 0$, es decir, $v_1 = 1 - N(N + 1)/2$. Aplicando la fórmula (3) y teniendo en cuenta que $y_0 = 0$, $v_N = 0$, $\nabla y_i = 1$, encontramos

$$\begin{aligned} S_N &= \sum_{i=1}^{N-1} i(i+1) = \sum_{i=0}^{N-1} y_i \Delta v_i = - \sum_{i=1}^N v_i \nabla y_i = \\ &= - \sum_{i=1}^{N-1} v_i = -(N-1)(v_1 - 1) - \frac{1}{2} \sum_{i=1}^{N-1} i(i+1) = \\ &= -\frac{1}{2} S_N + \frac{(N-1)N(N+1)}{2}, \end{aligned}$$

de modo que $S_N = \frac{1}{3} (N-1)N(N+1)$. De aquí se deduce que

$$\sum_{i=1}^N i^2 = 1^2 + 2^2 + \dots + N^2 = S_N + \sum_{i=1}^N i = \frac{N(N+1)(2N+1)}{6}.$$

§ 2. Ecuaciones en diferencias

1. **Ecuaciones en diferencias.** Una ecuación lineal respecto de la función reticular $y_i = y(i)$ ($i = 0, \pm 1, \pm 2, \dots$) $a_0(i)y(i) + a_1(i)y(i+1) + \dots + a_m(i)y(i+m) = f(i)$, (1)

donde $a_k(i)$ ($k = 0, 1, \dots, m$), $f(i)$ son las funciones reticulares prefijadas, $a_0(i) \neq 0$, $a_m(i) \neq 0$, lleva el nombre de *ecuación en diferencias lineal* de m -ésimo orden. Dicha ecuación contiene $m + 1$ valores de la función $y(i)$.

Haciendo uso de las fórmulas para las diferencias Δy_i , $\Delta^2 y_i, \dots, \Delta^{m-1} y_i$, podemos expresar los valores y_{i+1} , y_{i+2}, \dots, y_{m+1} en términos de y_i y las diferencias citadas: $y_{i+1} = y_i + \Delta y_i$, $y_{i+2} = \Delta^2 y_i + 2y_{i+1} - y_i = \Delta^2 y_i +$

+ $2\Delta y_i + y_i$, etc. Como resultado, obtendremos de (1) una nueva notación de la *ecuación en diferencias* de m -ésimo orden:

$$\alpha_0(i) y_i + \alpha_1(i) \Delta y_i + \dots + \alpha_m(i) \Delta^m y_i = f(i),$$

$$i = 0, \pm 1, \pm 2, \dots, \quad (2)$$

(por lo que se determina precisamente el término «ecuación en o de diferencias»). Si los coeficientes $\alpha_0, \alpha_1, \dots, \alpha_m$ no dependen de i , $\alpha_0 \neq 0$ y $\alpha_m \neq 0$, entonces (1) se denomina *ecuación en diferencias lineal de m -ésimo orden con coeficientes constantes*.

Para $m = 1$ de (1) se obtiene una ecuación en diferencias de *primer orden*

$$a_0(i) y_i + a_1(i) y_{i+1} = f(i), \quad a_0(i) \neq 0, \quad a_1(i) \neq 0; \quad (3)$$

cuando $m = 2$, obtenemos una ecuación en diferencias de *segundo orden*

$$a_0(i) y_i + a_1(i) y_{i+1} + a_2(i) y_{i+2} = f(i), \quad a_0(i) \neq 0, \\ a_2(i) \neq 0.$$

Nos limitaremos al estudio de las ecuaciones en diferencias de *primero y segundo órdenes*.

2. Ecuaciones de primer orden. Examinemos la ecuación en diferencias de primer orden (3). Al sustituir $y_{i+1} = y_i + \Delta y_i$, obtendremos

$$\bar{a}_0(i) y_i + a_1(i) \Delta y_i = f(i), \quad \bar{a}_0 = a_0 + a_1.$$

Como ejemplos más simples de ecuaciones en diferencias de primer orden pueden servir las ecuaciones para los términos de una progresión aritmética $y_{i+1} = y_i + d$ y de una progresión geométrica $y_{i+1} = qy_i$.

Escribamos la ecuación (3) en la forma

$$y_{i+1} = q_i y_i + \varphi_i, \quad (4)$$

donde $q_i = -a_0(i)/a_1(i)$, $\varphi_i = f(i)/a_1(i)$. De aquí se ve que la solución $y(i)$ está definida unívocamente para $i > i_0$, si está prefijado el valor $y(i_0)$. Supongamos que para $i = 0$ viene prefijado $y_0 = y(0)$. En tal caso podemos determinar

$y_1, y_2, \dots, y_i, \dots$ Eliminando sucesivamente según la fórmula (4) y_i, y_{i-1}, \dots, y_1 , obtendremos

$$y_{i+1} = q_i q_{i-1} \dots q_0 y_0 + \varphi_i + q_i \varphi_{i-1} + \dots + q_i q_{i-1} \dots q_1 \varphi_0,$$

o bien

$$y_{i+1} = \left(\prod_{h=0}^i q_h \right) y_0 + \sum_{h=0}^{i-1} \left(\prod_{s=h+1}^i q_s \right) \varphi_h + \varphi_i. \quad (5)$$

Para la ecuación con coeficientes constantes $q_i = q$, de lo que se tiene

$$y_{i+1} = q^{i+1} y_0 + \sum_{h=0}^i q^{i-h} \varphi_h, \quad i = 0, 1, 2, \dots, \quad (6)$$

que es una solución de la ecuación en diferencias (4) con coeficientes constantes.

3. Desigualdades de primer orden. Si el signo de igualdad en las expresiones de tipo (1) ó (2) lo sustituimos por los signos de desigualdad $<$, $>$, \leq , \geq , obtendremos *desigualdades en diferencias* de m -ésimo orden. Sea dada una desigualdad en diferencias de primer orden

$$y_{i+1} \leq q y_i + f_i, \quad i = 0, 1, 2, \dots, \quad q \geq 0; \quad (7)$$

sin restringir la generalidad de nuestros razonamientos, en adelante consideramos siempre que $q > 0$ (y_0, q, f_i son conocidos). Hallemos la solución de la desigualdad citada. Sea v_i una solución de la ecuación en diferencias

$$v_{i+1} = q v_i + f_i, \quad i = 0, 1, \dots, \quad v_0 = y_0. \quad (8)$$

En este caso queda lícita la estimación

$$y_i \leq v_i. \quad (9)$$

En efecto, al sustraer (8) de (7), encontramos

$$y_{i+1} - v_{i+1} \leq q (y_i - v_i) \leq q^2 (y_{i-1} - v_{i-1}) \leq \dots \leq q^{i+1} (y_0 - v_0) = 0.$$

Al poner en (9) la expresión explícita para v_i , tenemos

$$y_i \leq q^i y_0 + \sum_{h=0}^{i-1} q^{i-1-h} f_h, \quad i = 0, 1, 2, \dots, \quad (10)$$

lo que es la solución de la desigualdad (7).

4. Ecuaciones de segundo orden con coeficientes constantes. Analicemos una ecuación en diferencias de segundo orden

$$by_{i+1} - cy_i + ay_{i-1} = f_i, \quad i = 0, 1, \dots, \\ a \neq 0, \quad b \neq 0, \quad (11)$$

cuyos coeficientes no dependen de i . Si $f_i = 0$, la ecuación

$$by_{i+1} - cy_i + ay_{i-1} = 0, \quad i = 0, 1, \dots, \quad (12)$$

se llamará *homogénea*. Su solución se halla en la forma explícita.

Sea \bar{y}_i una solución de la ecuación homogénea (12), y sea y_i^* una solución cualquiera de la ecuación no homogénea (11). Entonces, la suma $y_i = \bar{y}_i + y_i^*$ será también una solución de la ecuación no homogénea:

$$b(\bar{y}_{i+1} + y_i^*) - c(\bar{y}_i + y_i^*) + a(\bar{y}_{i-1} + y_{i-1}^*) = \\ = [b\bar{y}_{i+1} - c\bar{y}_i + a\bar{y}_{i-1}] + [by_i^* - cy_i^* + ay_{i-1}^*] = f_i.$$

Esta propiedad se debe a la linealidad de la ecuación (11); ella queda en vigor para la ecuación en diferencias (1) de cualquier orden. Es evidente que si \bar{y}_i es una solución de la ecuación homogénea (12), entonces también $c\bar{y}_i$ (donde c es una constante arbitraria) satisface la citada ecuación.

Sean $y_i^{(1)}$ e $y_i^{(2)}$ dos soluciones de la ecuación (12). Se denominarán *linealmente independientes*, si la igualdad

$$c_1 y_i^{(1)} + c_2 y_i^{(2)} = 0, \quad i = 0, 1, 2, \dots,$$

se verifica sólo cuando $c_1 = c_2 = 0$. Esta afirmación es equivalente a la exigencia de que el determinante del sistema

$$c_1 y_i^{(1)} + c_2 y_i^{(2)} = 0, \\ c_1 y_{i+m}^{(1)} + c_2 y_{i+m}^{(2)} = 0, \quad m = \pm 1, \pm 2, \dots,$$

sea distinto de cero para cualesquiera i, m . En particular,

$$\Delta_{i, i+1} = \begin{vmatrix} y_i^{(1)} & y_i^{(2)} \\ y_{i+1}^{(1)} & y_{i+1}^{(2)} \end{vmatrix} \neq 0.$$

Al igual que en la teoría de las ecuaciones diferenciales, se puede introducir la noción de *solución general* de la

ecuación en diferencias (12) y mostrar que si las soluciones $y_i^{(1)}$, $y_i^{(2)}$ son linealmente independientes, la solución general de la ecuación (12) tendrá la forma

$$y_i = c_1 y_i^{(1)} + c_2 y_i^{(2)},$$

donde c_1 y c_2 son unas constantes arbitrarias. La solución general de la ecuación no homogénea (11) puede representarse en la forma

$$y_i = c_1 y_i^{(1)} + c_2 y_i^{(2)} + y_i^*, \quad (13)$$

donde y_i^* es una solución (particular) cualquiera de la ecuación (11). Lo mismo que en el caso de las ecuaciones diferenciales, para determinar c_1 y c_2 se deben definir las condiciones complementarias iniciales o las de contorno.

La solución particular de la ecuación (12) puede hallarse en la forma explícita. Buscaremos dicha solución en la forma $y_i = q^i$, donde $q \neq 0$ es un número hasta ahora desconocido. Al realizar la sustitución $y_k = q^k$ en (12), obtendremos una ecuación cuadrática $bq^2 - cq + a = 0$, cuyas raíces son

$$q_1 = \frac{c + \sqrt{c^2 - 4ab}}{2b}, \quad q_2 = \frac{c - \sqrt{c^2 - 4ab}}{2b}. \quad (14)$$

Según sean los valores del discriminante $D = c^2 - 4ab$, son posibles tres casos:

1) $D = c^2 - 4ab > 0$. Las raíces q_1 y q_2 son reales y distintas. Les corresponden las soluciones particulares

$$y_k^{(1)} = q_1^k, \quad y_k^{(2)} = q_2^k;$$

estas soluciones son linealmente independientes, puesto que es distinto de cero el determinante:

$$\Delta_{k, k+1} = \begin{vmatrix} q_1^k & q_1^{k+1} \\ q_2^k & q_2^{k+1} \end{vmatrix} = q_1^k q_2^k (q_2 - q_1) \neq 0$$

Ha de notarse que $q_1 \neq 0$ y $q_2 \neq 0$, pues en el caso contrario $a = 0$ y la ecuación (12) no sería ecuación en diferencias de segundo orden. La solución general de la ecuación (12) tiene por expresión

$$y_k = c_1 q_1^k + c_2 q_2^k. \quad (15)$$

2) $D = c^2 - 4ab < 0$. La ecuación cuadrática cuenta con las raíces complejas conjugadas

$$q_1 = \frac{c + i\sqrt{|D|}}{2b}; \quad q_2 = \frac{c - i\sqrt{|D|}}{2b},$$

donde i es la unidad imaginaria. Resulta cómodo representar estas raíces en la forma

$$q_1 = \rho e^{i\varphi}, \quad q_2 = \rho e^{-i\varphi}, \quad \rho = \sqrt{\frac{a}{b}},$$

$$\varphi = \arctg \frac{\sqrt{|D|}}{c}.$$

No sólo las funciones

$$q_1^k = \rho^k e^{ik\varphi} = \rho^k (\cos k\varphi - i \operatorname{sen} k\varphi),$$

$$q_2^k = \rho^k e^{-ik\varphi} = \rho^k (\cos k\varphi + i \operatorname{sen} k\varphi)$$

representan soluciones particulares sino también las funciones siguientes:

$$y_k^{(1)} = \rho^k \cos k\varphi, \quad y_k^{(2)} = \rho^k \operatorname{sen} k\varphi,$$

las cuales son linealmente independientes en virtud de la independencia lineal de las funciones $\operatorname{sen} k\varphi$ y $\cos k\varphi$. La solución general tiene la forma

$$y_k = \rho^k (c_1 \cos k\varphi + c_2 \operatorname{sen} k\varphi). \quad (16)$$

3) $D = c^2 - 4ab = 0$. Las raíces son reales e iguales: $q_1 = q_2 = c/(2b) = q_0$. Las soluciones

$$y_k^{(1)} = q_0^k, \quad y_k^{(2)} = kq_0^k \quad (17)$$

son linealmente independientes. Mostraremos que $y_k^{(2)}$ es una solución de la ecuación (12):

$$by_{k+1}^{(2)} - cy_k^{(2)} + ay_{k-1}^{(2)} = b(k+1)q_0^{k+1} - ckq_0^k + a(k-1)q_0^{k-1} =$$

$$= k(bq_0^{k+1} - cq_0^k + aq_0^{k-1}) + (bq_0^2 - a)q_0^{k-1} = 0,$$

puesto que $bq_0^2 - a = b \frac{c^2}{4b^2} - a = \frac{D}{4b} = 0$. Como $\Delta_{k, k+1} =$

$$= \begin{vmatrix} q_0^k & kq_0^k \\ q_0^{k+1} & (k+1)q_0^{k+1} \end{vmatrix} = q_0^{2k+1} \neq 0, \quad \text{entonces las solu-}$$

ciones (17) son linealmente independientes y la solución general tiene por expresión

$$y_h = c_1 q_0^h + c_2 k q_0^h.$$

5. Ejemplos. Veamos algunos ejemplos de resolución de las ecuaciones en diferencias de segundo orden (11).

1. Hállese la solución general de la ecuación

$$y_{h+1} - 2p y_h + y_{h-1} = 0, \quad a = b = 1, \quad c = 2p > 0.$$

Son posibles tres casos. 1) $p < 1$. Ponemos $p = \cos \alpha$; entonces $D = 4(\cos^2 \alpha - 1) = -4\text{sen}^2 \alpha < 0$. Las soluciones particulares tienen la forma

$$y_h^{(1)} = \cos k\alpha, \quad y_h^{(2)} = \text{sen } k\alpha.$$

2) $p > 1$. Suponiendo $p = \text{ch } \alpha$, obtendremos para q una ecuación cuadrática $q^2 - 2\text{ch } \alpha q + 1 = 0$; su discriminante es $D = 4(\text{ch}^2 \alpha - 1) = 4\text{sh}^2 \alpha$, mientras que las raíces tienen por expresión $q_{1,2} = \text{ch } \alpha \pm \text{sh } \alpha = e^{\pm \alpha}$. El papel de soluciones particulares desempeñan las funciones

$$y_h^{(1)} = \text{ch } k\alpha, \quad y_h^{(2)} = \text{sh } k\alpha.$$

3) $p = 1$. En este caso $q^2 - 2q + 1 = 0$, $q_{1,2} = 1$, las soluciones particulares tienen la forma $y_h^{(1)} = 1$, $y_h^{(2)} = k$, y la solución general es

$$y_h = c_1 + c_2 k.$$

2. Hállese la solución de la ecuación

$$y_{h+2} - y_{h+1} = 2y_h = 0.$$

El discriminante es igual a $D = 1 + 8 = 9$, las raíces serán $q_{1,2} = (1 \pm 3)/2$, $q_1 = 2$, $q_2 = -1$. La solución general es de la forma

$$y_h = c_1 2^h + c_2 (-1)^h.$$

3. Hállese la solución general de la ecuación

$$y_{h+1} - y_h - 6y_{h-1} = 2^{h+1}. \quad (18)$$

La solución general de una ecuación no homogénea es la suma $y_h = \bar{y}_h + y_h^*$ de la solución general \bar{y}_h de la ecuación homogénea y de la solución particular y_h^* de la ecuación

no homogénea. Hallemos primero la solución general de la ecuación homogénea. El discriminante es $D = 1 + 24 = 25 > 0$, y las raíces de la ecuación cuadrática $q^2 - q - 6 = 0$ son $q_1 = 3$, $q_2 = -2$, de suerte que $y_k^{(1)} = 3^k$, $y_k^{(2)} = (-2)^k$. La solución particular y_k^* buscaremos en la forma $y_k^* = c2^k$, donde $c = \text{const.}$ Sustituyendo $y_k^* = c2^k$ en (18), obtendremos $c(2^{k+1} - 2^k - 6 \cdot 2^{k-1}) = c \cdot 2^{k-1}(-4) = 2^{k+1}$, $c = -1$.

La solución general de la ecuación (18) tiene por expresión

$$y_k = c_1 \cdot 3^k + c_2 (-2)^k - 2^k.$$

6. Ecuación en diferencias de segundo orden con coeficientes variables. Problema de Cauchy y problema de contorno. Examinemos ahora una ecuación en diferencias con coeficientes variables

$$\begin{aligned} b_i y_{i+1} - c_i y_i + a_i y_{i-1} &= f_i, \\ a_i \neq 0, \quad b_i \neq 0, \quad i &= 0, 1, 2, \dots \end{aligned} \quad (19)$$

Dado que $b_i \neq 0$, de (19) obtenemos la siguiente relación recurrente:

$$y_{i+1} = \frac{c_i y_i - a_i y_{i-1} + f_i}{b_i}, \quad b_i \neq 0. \quad (20)$$

Expresemos y_{i+1} e y_{i-1} en términos de y_i y las diferencias de primero y segundo órdenes. La ecuación (19) se anotará en este caso en la forma

$$\begin{aligned} \Delta \nabla y_i + (b_i - a_i) \Delta y_i - (c_i - a_i - b_i) y_i &= f_i, \\ a_i \neq 0, \quad b_i \neq 0. \end{aligned}$$

La solución de una ecuación en diferencias de primer orden depende de una constante arbitraria y se determina unívocamente, si está prefijada una condición complementaria, por ejemplo, $y_0 = c_0$. La solución de la ecuación de segundo orden se determina por dos constantes arbitrarias y se puede hallarla, siempre que vienen dadas dos condiciones complementarias. Si ambas condiciones están dadas en dos puntos vecinos, se trata de un *problema de Cauchy*. Si las dos condiciones están dadas en dos puntos diferentes (pero no vecinos), obtenemos un *problema de contorno*. De mayor interés para nosotros serán precisamente los

problemas de contorno. Introduzcamos las designaciones

$$Ly_i = b_i y_{i+1} - c_i y_i + a_i y_{i-1}$$

y enunciemos los problemas mencionados más detalladamente.

PROBLEMA DE CAUCHY: hállese la solución de la ecuación

$$Ly_i = f_i, \quad i = 1, 2, \dots, \quad (21)$$

con las condiciones complementarias

$$y_0 = \mu_1, \quad y_1 = \mu_2. \quad (22)$$

La segunda condición de (22) puede notarse de otro modo: $\Delta y_0 = y_1 - y_0 = \mu_2 - \mu_1 = \bar{\mu}_1$; podemos, entonces, decir que en el caso del problema de Cauchy vienen dadas en un punto $i = 0$ las magnitudes

$$y_0 = \mu_1, \quad \Delta y_0 = \bar{\mu}_1. \quad (22')$$

PROBLEMA DE CONTORNO: hállese la solución de la ecuación

$$Ly_i = f_i, \quad i = 1, 2, \dots, N - 1$$

para las condiciones complementarias

$$y_0 = \mu_1, \quad y_N = \mu_2, \quad N \geq 2. \quad (23)$$

En los nodos de frontera $i = 0$ e $i = N$ pueden definirse no sólo los valores de las funciones, sino también sus diferencias y combinaciones, es decir, las expresiones $\alpha_1 \Delta y_0 + \beta_1 y_0$ para $i = 0$, y $\alpha_2 \nabla y_N + \beta_2 y_N$ para $i = N$. Tales condiciones se pueden anotar en la forma

$$y_0 = \kappa_1 y_1 + \mu_1, \quad y_N = \kappa_2 y_{N-1} + \mu_2. \quad (24)$$

Si $\kappa_1 = \kappa_2 = 0$, obtenemos de aquí las *condiciones de primer género*; cuando $\kappa_1 = 1$, $\kappa_2 = 1$, tenemos *condiciones de segundo género*

$$\Delta y_0 = -\mu_1, \quad \nabla y_N = \mu_2. \quad (25)$$

Si $\kappa_{1,2} \neq 0; 1$, (24) llevan el nombre de *condiciones de tercer género*:

$$\begin{aligned} -\kappa_1 \Delta y_0 + (1 - \kappa_1) y_0 &= \mu_1, \\ \kappa_2 \nabla y_N + (1 - \kappa_2) y_N &= \mu_2. \end{aligned} \quad (26)$$

Además, pueden encontrarse problemas de contorno con ciertas combinaciones de las citadas condiciones de con-

torno: las condiciones de un tipo para $i = 0$, y condiciones de otro tipo, para $i = N$.

La solución del problema de Cauchy se halla directamente de la ecuación (21) según la fórmula recurrente (20), tomando en consideración los datos iniciales $y_0 = \mu_1$, $y_1 = \mu_2$. Para la resolución de los problemas de contorno se emplea un método más complejo (método de eliminaciones) el cual se expondrá en adelante.

Para una ecuación con coeficientes constantes la solución del problema de contorno puede expresarse en la forma explícita.

EJEMPLO. Hállese la solución del problema de contorno $\Delta^2 y_{i-1} = 1$, $i = 1, 2, \dots, N-1$, $y_0 = 0$, $y_N = 0$. (27)

La ecuación homogénea $\Delta^2 y_{i-1} = y_{i+1} - 2y_i + y_{i-1} = 0$ tiene su solución general $\bar{y}_i = c_1 + c_2 i$. La solución particular y_i^* de la ecuación no homogénea $\Delta^2 y_{i-1} = y_{i+1} - 2y_i + y_{i-1} = 1$ se buscará en la forma $y_i^* = ci^2$. Al sustituir esta expresión en la ecuación (27), encontramos $\Delta^2 y_{i-1}^* = c((i+1)^2 - 2i^2 + (i-1)^2) = 1$, es decir, $c = 1/2$, de suerte que $y_i = \bar{y}_i + y_i^* = c_1 + c_2 i + i^2/2$. Con el fin de determinar c_1 y c_2 se usan las condiciones de contorno para $i = 0$, $i = N$: $y_0 = c_1 = 0$, $y_N = c_2 N + N^2/2 = 0$, $c_2 = -N/2$. De este modo,

$$y_i = -\frac{1}{2} iN + \frac{1}{2} i^2 = -\frac{1}{2} i(N-i)$$

es la solución del problema (27).

§ 3. Resolución de los problemas de contorno en diferencias para las ecuaciones de segundo orden

1. Resolución de los problemas de contorno en diferencias por el método de factorización. Un problema de contorno

$$a_i y_{i-1} - c_i y_i + b_i y_{i+1} = -f_i, \quad a_i \neq 0, \quad b_i \neq 0, \\ i = 1, 2, \dots, N-1, \quad (1)$$

$$y_0 = \alpha_1 y_1 + \mu_1, \quad y_N = \alpha_2 y_{N-1} + \mu_2$$

representa el sistema de ecuaciones algebraicas lineales con matriz tridiagonal de dimensión $(N + 1) \times (N + 1)$:

$$A = \begin{bmatrix} 1 & -\kappa_1 & 0 & \dots & 0 & 0 & 0 & \dots & 0 & 0 & 0 & 0 \\ a_1 & -c_1 & b_1 & \dots & 0 & 0 & 0 & \dots & 0 & 0 & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & a_i & -c_i & b_i & \dots & 0 & 0 & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & 0 & 0 & \dots & a_{N-1} & -c_{N-1} & b_{N-1} & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 & \dots & 0 & -\kappa_2 & 1 & 0 & 0 \end{bmatrix}.$$

En lugar de (1) podemos escribir

$$Ay = f, \quad y = (y_0, y_1, \dots, y_N),$$

$$f = (\mu_1, -f_1, \dots, -f_{N-1}, \mu_2). \quad (2)$$

En el caso del primer problema del contorno la matriz correspondiente tiene la dimensión $(N - 1) \times (N - 1)$.

Para resolver el problema de contorno (1) puede emplearse el siguiente método de eliminación llamado *método de factorización*. Supongamos que tiene lugar la relación

$$y_i = \alpha_{i+1}y_{i+1} + \beta_{i+1} \quad (3)$$

con coeficientes indeterminados α_{i+1} y β_{i+1} , y sustituyamos $y_{i-1} = \alpha_i y_i + \beta_i$ en (1):

$$(a_i \alpha_i - c_i) y_i + b_i y_{i+1} = -(f_i + a_i \beta_i);$$

al comparar esta identidad con (3), encontramos

$$\alpha_{i+1} = \frac{b_i}{c_i - a_i \alpha_i}, \quad i = 1, 2, \dots, N-1, \quad (4)$$

$$\beta_{i+1} = \frac{a_i \beta_i + f_i}{c_i - a_i \alpha_i}, \quad i = 1, 2, \dots, N-1. \quad (5)$$

Con el fin de hallar α_1, β_1 hagamos uso de la condición de contorno para $i = 0$. De las fórmulas (3) y (1) encontramos para $i = 0$:

$$\alpha_1 = \kappa_1, \quad \beta_1 = \mu_1. \quad (6)$$

Conociendo α_1, β_1 y pasando de i a $i + 1$ en las fórmulas (4) y (5), determinamos α_i y β_i para cualquier $i = 2, 3, \dots, N$. Los cálculos según la fórmula (3) se llevan a cabo

pasando de $i + 1$ a i (es decir, conociendo y_{i+1} , hallamos y_i), y para iniciar los cálculos mencionados se debe fijar y_N . Determinemos y_N partiendo de la condición de contorno $y_N = \kappa_2 y_{N-1} + \mu_2$ y de la condición (3) para $i = N - 1$: $y_{N-1} = \alpha_N y_N + \beta_N$. De aquí encontramos

$$y_N = \frac{\mu_2 + \kappa_2 \beta_N}{1 - \alpha_N \kappa_2}. \quad (7)$$

Reunamos ahora todas las fórmulas de factorización y escribámoslas en el orden de su aplicación:

$$\stackrel{(-)}{\alpha}_{i+1} = \frac{b_i}{c_i - a_i \alpha_i}, \quad i = 1, 2, \dots, N-1, \quad \alpha_1 = \kappa_1; \quad (8)$$

$$\stackrel{(-)}{\beta}_{i+1} = \frac{a_i \beta_i + f_i}{c_i - a_i \alpha_i}, \quad i = 1, 2, \dots, N-1, \quad \beta_1 = \mu_1; \quad (9)$$

$$\stackrel{(-)}{y}_i = \alpha_{i+1} y_{i+1} + \beta_{i+1}, \quad i = N-1, N-2, \dots, 2, 1, 0,$$

$$y_N = \frac{\mu_2 + \kappa_2 \beta_N}{1 - \alpha_N \kappa_2}. \quad (10)$$

Las flechas indican la dirección de cálculo: (\rightarrow) de i a $i + 1$, (\leftarrow) de $i + 1$ a i .

De este modo, el problema de contorno para la ecuación de segundo orden se ha reducido a tres problemas de Cauchy para las ecuaciones de primer orden.

2. Estabilidad del método de factorización. Las fórmulas de factorización pueden emplearse, si los denominadores de las fracciones (8) y (10) no se reducen a cero. Las condiciones suficientes para ello están representadas por las desigualdades

$$\begin{aligned} |c_i| &\geq |a_i| + |b_i|, & i = 1, 2, \dots, N-1, \\ | \kappa_1 | &\leq 1, & | \kappa_2 | \leq 1, & | \kappa_1 | + | \kappa_2 | < 2. \end{aligned} \quad (11)$$

Problemos que siendo cumplidas las condiciones (11), los denominadores $c_i - a_i \alpha_i$ y $1 - \alpha_N \kappa_2$ no se reducen a cero y que

$$| \alpha_i | \leq 1, \quad i = 1, 2, \dots, N. \quad (12)$$

Supongamos que $| \alpha_i | \leq 1$, y mostremos que $| \alpha_{i+1} | \leq 1$; entonces de aquí y de la condición $| \alpha_i | = | \kappa_1 | \leq 1$ se deducirá (12). Examinaremos la diferencia $|c_i - a_i \alpha_i| -$

$-|b_i| \geq |c_i| - |a_i| |\alpha_i|$, $-|b_i| \geq |a_i| (1 - |\alpha_i|) \geq 0$, de modo que $|c_i - a_i \alpha_i| \geq |b_i| > 0$, y $|\alpha_{i+1}| = |b_i| / |c_i - a_i \alpha_i| \leq 1$.

Observemos que si $|c_{i_0}| > |a_{i_0}| + |b_{i_0}|$ siquiera en un solo punto $i = i_0$, entonces $|\alpha_i| < 1$ para todo $i > i_0$, incluso para $i = N$: $|\alpha_N| < 1$. En este caso $|1 - \alpha_N \kappa_2| \geq 1 - |\alpha_N| |\kappa_2| \geq 1 - |\alpha_N| > 0$, y la condición $|\kappa_1| + |\kappa_2| < 2$ será superflua. Si $|\kappa_1| < 1$, entonces $|\alpha_N| < 1$. En cambio, si $|\kappa_1| = 1$, entonces $|\kappa_2| < 1$ y $|\alpha_N| \leq 1$, y tenemos $|1 - \alpha_N \kappa_2| \geq 1 - |\alpha_N| |\kappa_2| \geq 1 - |\kappa_2| > 0$. De este modo, si se cumplen las condiciones (11), el problema (1) tiene la única solución la cual se halla según las fórmulas de factorización (8)–(10).

Los cálculos según las fórmulas (8)–(10) se realizan en un ordenador aproximadamente, con un número finito de cifras significativas. A consecuencia de los errores de redondeo se halla, de hecho, no la función y_i (la cual representa solución del problema (1)), sino \tilde{y}_i , esto es, la solución del mismo problema con coeficientes perturbados $\tilde{a}_i, \tilde{b}_i, \tilde{c}_i, \tilde{\kappa}_1, \tilde{\kappa}_2$ y segundos miembros $\tilde{f}_i, \tilde{\mu}_1, \tilde{\mu}_2$. Surge naturalmente la cuestión de si ocurre o no, en transcurso de los cálculos, el aumento del error de redondeo, lo que puede conducir tanto a la pérdida de precisión, como a la imposibilidad de continuar los cálculos a causa del crecimiento de las magnitudes que se determinan. A título de ejemplo puede servir la búsqueda de y_i según la fórmula $y_{i+1} = qy_i$ para $q > 1$. Puesto que $y_n = q^n y_0$, para cualquier y_0 puede indicarse tal n_0 que y_{n_0} será el infinito de ordenador. Realmente, en virtud de los errores de redondeo, se determina no el valor exacto de y_i , sino el valor de \tilde{y}_i a partir de la ecuación $\tilde{y}_{i+1} = q\tilde{y}_i + \eta$, donde η es el error de redondeo. Para el error $\delta y_i = \tilde{y}_i - y_i$ obtendremos una ecuación $\delta y_{i+1} = q\delta y_i + \eta$ ($i = 0, 1, \dots, \delta y_0 = \eta$). De la fórmula $\delta y_i = q^i \eta + \eta (q^i - 1)/(q - 1)$ se ve que el error δy_i va creciendo, para $q > 1$, de manera exponencial a medida que crece i .

Volvamos al método de factorización y probemos que el error δy_i no aumenta cuando $|\alpha_i| \leq 1$. Efectivamente, de

las igualdades $\tilde{y}_i = \alpha_{i+1}\tilde{y}_{i+1} + \beta_{i+1}$, $y_i = \alpha_{i+1}y_{i+1} + \beta_{i+1}$ proviene $\delta y_i = \alpha_{i+1}\delta y_{i+1}$, $|\delta y_i| \leq |\alpha_{i+1}| |\delta y_{i+1}| \leq |\delta y_{i+1}|$, porque $|\alpha_{i+1}| \leq 1$.

Tomando en consideración que en el transcurso de los cálculos se perturban también los coeficientes α_{i+1} , β_{i+1} , se puede señalar que el error δy_i es proporcional al cuadrado del número de nodos N :

$$\max_{1 \leq i \leq N} |\delta y_i| \leq \varepsilon_0 N^2,$$

donde ε_0 es el error de redondeo. De aquí se ve la relación que existe entre la precisión requerida ε de la solución del problema, el número N de ecuaciones y el número de cifras significativas del ordenador, puesto que $\varepsilon_0 N^2 \approx \varepsilon$.

3. Otras variantes del método de factorización. El método de factorización (8)–(10) analizado más arriba, en el cual la determinación de y_i se realiza sucesivamente de derecha a izquierda, se denomina *factorización derecha*. Análogamente se anotan las fórmulas de la *factorización izquierda*:

$$\xi_i^{(-)} = \frac{a_i}{c_i - b_i \xi_{i+1}}, \quad i = N-1, N-2, \dots, 2, 1, \quad \xi_N = \kappa_2, \quad (13)$$

$$\eta_i^{(-)} = \frac{b_i \eta_{i+1} + f_i}{c_i - b_i \xi_{i+1}}, \quad i = N-1, N-2, \dots, 2, 1, \quad \eta_N = \mu_2, \quad (14)$$

$$y_{i+1} = \xi_{i+1} y_i + \eta_{i+1}, \quad i = 0, 1, \dots, N-1,$$

$$y_0 = \frac{\mu_1 + \kappa_1 \eta_1}{1 - \xi_1 \kappa_1}. \quad (15)$$

En efecto, suponiendo que $y_{i+1} = \xi_{i+1} y_i + \eta_{i+1}$, eliminemos de (1) y_{i+1} ; obtendremos

$$-f_i = a_i y_{i-1} + (b_i \xi_{i+1} - c_i) y_i + b_i \eta_{i+1},$$

o bien

1258

$$y = \frac{a_i}{c_i - b_i \xi_{i+1}} y_{i-1} + \frac{f_i + b_i \eta_{i+1}}{c_i - b_i \xi_{i+1}}.$$

Al cotejar con la fórmula $y_i = \xi_i y_{i-1} + \beta_i$, obtendremos (13) y (14). El valor de y_0 hallamos de la condición $y_0 = \kappa_1 y_1 + \mu_1$ y de la fórmula $y_0 = \xi_1 y_1 + \eta_1$. De la desi-

gualdad $|c_i - b_i \xi_{i+1}| \geq |c_i| - |b_i| |\xi_{i+1}| \geq |a_i| + |b_i| (1 - |\xi_{i+1}|)$, $|1 - \xi_1 \kappa_1| \geq 1 - |\xi_1| |\kappa_1|$ se ve que las condiciones (11) garantizan que las fórmulas de factorización izquierda sean aplicables y su cálculo sea estable, puesto que $|\xi_i| \leq 1$ ($i = 1, 2, \dots, N$).

La combinación de las factorizaciones izquierda y derecha da el *método de factorizaciones opuestas*. Empleándose este método, en la región $0 \leq i \leq i_0 + 1$ se calculan, según las fórmulas (8), (9), los coeficientes de factorización α_i , β_i , y en la región $i_0 \leq i \leq N$ se hallan, por las fórmulas (13), (14), ξ_i y η_i . Cuando $i = i_0$, se realiza la conjugación de soluciones en la forma (10) y (15).

De las fórmulas $y_{i_0} = \alpha_{i_0+1} y_{i_0+1} + \beta_{i_0+1} y_{i_0+1} = \xi_{i_0+1} y_{i_0} + \eta_{i_0+1}$ hallamos

$$y_{i_0} = \frac{\beta_{i_0+1} + \alpha_{i_0+1} \eta_{i_0}}{1 - \alpha_{i_0+1} \xi_{i_0+1}}.$$

La citada fórmula tiene sentido, puesto que por lo menos una de las magnitudes $|\xi_{i_0+1}|$ ó $|\alpha_{i_0+1}|$ es, en virtud de (11), inferior a la unidad y, por lo tanto, $1 - \alpha_{i_0+1} \xi_{i_0+1} > 0$. Al conocer y_{i_0} , podemos hallar, sirviéndonos de la fórmula (10), todos los valores de y_i para $i < i_0$, y, por la fórmula (15), todos los valores de y_i para $i > i_0$. Cuando $i > i_0$ e $i < i_0$, los cálculos son autónomos (se llevan a cabo paralelamente). El método de factorizaciones opuestas es particularmente cómodo, si, por ejemplo, se pide hallar y_i sólo en un nodo $i = i_0$.

§ 4. Ecuaciones en diferencias como ecuaciones operacionales

1. Espacio lineal*). Veamos un conjunto H de elementos x, y, z, \dots , respecto de los cuales se sabe que a cada par de elementos x e y , pertenecientes a H , se le pone en correspondencia de tal o cual modo un elemento tercero $z \in H$, llamado suma de los dos elementos primeros y designado $z = x + y$: a todo elemento $x \in H$ y a cada número λ

*) Véase, por ejemplo, V. Ilyin, E. Poznyak, "Linear algebra". Editorial Mir, Moscú, 1985.

se les pone en correspondencia un elemento $u \in H$, denominado producto de x por el número λ y designado $u = \lambda x$.

El conjunto H se llamará *espacio lineal*, si las operaciones de sumación y multiplicación por un número, determinadas para sus elementos x, y, z, \dots , satisfacen los siguientes axiomas:

1) $x + y = y + x$ para cualesquiera $x, y \in H$ (conmutatividad de la sumación);

2) $(x + y) + z = x + (y + z)$ para cualesquiera $x, y, z \in H$ (asociatividad de la sumación);

3) existe un elemento «cero», designado 0 , tal que $x + 0 = x$ para cualquier $x \in H$;

4) para todo elemento $x \in H$ existe un elemento opuesto $(-x)$ tal que $x + (-x) = 0$;

5) $1 \cdot x = x$;

6) $(\lambda\mu)x = \lambda(\mu x)$ (asociatividad de la multiplicación);

7) $\lambda(x + y) = \lambda x + \lambda y$; $(\lambda + \mu)x = \lambda x + \mu x$ (distributividad de la multiplicación respecto a sumación), donde λ y μ son unos números cualesquiera.

Un espacio lineal se denomina *complejo*, si para sus elementos está definida la multiplicación por números complejos y se llama *real*, si viene definida solamente la multiplicación por números reales.

Los elementos x, y, z, \dots del espacio lineal H llevan el nombre de *vectores*.

Los vectores x_1, x_2, \dots, x_N se denominan *linealmente independientes*, siempre que la igualdad

$$c_1 x_1 + c_2 x_2 + \dots + c_N x_N = 0 \quad (1)$$

se verifica sólo cuando $c_1 = c_2 = \dots = c_N = 0$. Si, en cambio, existen tales c_1, c_2, \dots, c_N (no todos iguales a cero) que tiene lugar la igualdad (1), entonces los vectores x_1, \dots, x_N se llamarán *linealmente dependientes*. El número máximo (si existe) de vectores linealmente independientes del espacio lineal H se denomina *dimensión* del espacio citado. Un espacio que posee una infinidad de vectores linealmente independientes, se denomina de *dimensión infinita*.

El espacio H se llama *normado*, si para cada $x \in H$ viene definido un número real $\|x\|$, denominado *norma*, que satisface las siguientes condiciones:

- 1) $\|x\| > 0$ para $x \neq 0$; $\|x\| = 0$, si $x = 0$;
- 2) $\|x + y\| \leq \|x\| + \|y\|$ (desigualdad triangular);
- 3) $\|cx\| = |c| \cdot \|x\|$, donde c es un número.

Se denomina espacio *euclídeo* (*unitario*, respectivamente) el espacio lineal real de dimensión finita H (espacio lineal complejo de dimensión finita H , respectivamente), en el cual a todo par de vectores x, y se les ha puesto en correspondencia un número real (complejo) (x, y) , denominado *producto escalar* de dichos vectores, con la particularidad de que se consideran cumplidas las siguientes condiciones:

Para el caso de un espacio euclídeo:

- 1) $(x, y) = (y, x)$ (simetría);
- 2) $(x_1 + x_2, y) = (x_1, y) + (x_2, y)$ (distributividad);
- 3) $(\lambda x, y) = \lambda (x, y)$ (homogeneidad), donde λ es un número real cualquiera;
- 4) si $x \neq 0$, entonces $(x, x) > 0$.

Para el caso de un espacio unitario;

- 1) $(x, y) = \overline{(y, x)}$;
- 2) $(x_1 + x_2, y) = (x_1, y) + (x_2, y)$;
- 3) $(\lambda x, y) = \lambda (x, y)$ para cualquier número complejo λ ;
- 4) si $x \neq 0$, entonces $(x, x) > 0$.

Hemos de observar que el producto escalar introducido (x, y) engendra en H la norma

$$\|x\| = \sqrt{(x, x)}. \quad (2)$$

Resulta válida aquí la desigualdad de Cauchy—Buniakovski

$$|(x, y)|^2 \leq (x, x) \cdot (y, y), \quad (3)$$

la cual puede escribirse, tomando en consideración (2), en la forma

$$|(x, y)| \leq \|x\| \cdot \|y\|.$$

2. Operadores lineales en un espacio de dimensión finita.

Sea H un espacio lineal de dimensión finita provisto de producto escalar (x, y) . Designemos con D cierto subespacio de H . Si a todo vector $x \in D$ se le ha puesto en correspondencia, de acuerdo con una regla determinada, el vector $y = Ax$ de H , suele decirse que en H está dado el *operador* A . El conjunto $D \subset U$ se llama *dominio de definición* del opera-

por A y se designa $D(A)$. Un conjunto de todos los vectores del tipo $y = Ax$, $x \in D(A)$ se denomina *campo de valores* del operador A y se denota $R(A)$. Si $D(A) = H$, dicen que el operador A está prefijado sobre H .

El operador A se llama *lineal*, si a) es aditivo, es decir, $A(x_1 + x_2) = Ax_1 + Ax_2$ para cualesquiera $x_1, x_2 \in H$; b) es homogéneo, es decir, $A(cx) = cAx$ para todo $x \in H$ y cualesquiera números c . Los requisitos a) y b) son equivalentes a la condición $A(c_1x_1 + c_2x_2) = c_1Ax_1 + c_2Ax_2$, cualesquiera que sean $x_1, x_2 \in H$ y los números c_1 y c_2 .

Un operador lineal se denomina *acotado*, si existe tal constante $M > 0$ que

$$\|Ax\| \leq M \|x\| \quad \text{para todo } x \in H. \quad (4)$$

La cota inferior exacta del conjunto de números M que satisfacen la condición (4) lleva el nombre de *norma* del operador A y se denota $\|A\|$. Está claro que

$$\|Ax\| \leq \|A\| \cdot \|x\|. \quad (5)$$

En adelante se considerarán siempre operadores lineales acotados A prefijados sobre H con el campo de valores $R(A) \subseteq H$. Tal operador A aplica H en H , lo que se escribe en la forma: $A: H \rightarrow H$.

En el espacio de dimensión finita cualquier operador lineal es acotado.

Si a cada $y \in H$ le corresponde sólo un vector $x \in H$, para el cual $Ax = y$, entonces mediante esta correspondencia queda definido un operador A^{-1} , denominado *inverso*: $A^{-1}: H \rightarrow H$. De la definición de operador inverso A^{-1} proviene que

$$A^{-1}(Ax) = x, \quad A(A^{-1}y) = y \quad \text{para cualesquiera } x, y \in H.$$

Un operador D que actúa según la regla $Dx = A(Bx)$ recibe el nombre de *producto* de los operadores A y B y se designa $D = AB$. Un operador E se denomina *operador unidad (idéntico)*, si $Ex = x$ para todos los $x \in H$. Si existe A^{-1} , entonces $A^{-1}A = AA^{-1} = E$. Los operadores A y B se llaman *permutables* o conmutativos, si $AB = BA$.

Es evidente que A^{-1} es un operador lineal, si lo es el operador A . Resulta válida la siguiente afirmación:

Para que un operador lineal $A: H \rightarrow H$ cuente con su inverso, es necesario y suficiente que la ecuación $Ax = 0$ tenga la única solución $x = 0$.

Un operador $A^*: H \rightarrow H$ se denomina *conjugado* del operador $A: H \rightarrow H$, si

$$(Ax, y) = (x, A^*y) \text{ para cualesquiera } x, y \in H.$$

El operador A es *autoconjugado* (*simétrico*), siempre que $A = A^*$ (o bien $(Ax, y) = (x, Ay)$ para cualesquiera $x, y \in H$). El operador lineal A se llamará: *positivo*, si $(Ax, x) > 0$ ($x \in H; x \neq 0$); *definido positivo*, si $(Ax, x) \geq \delta \|x\|^2$ ($x \in H$), donde $\delta > 0$ es un número; *no negativo*, si $(Ax, x) \geq 0$ ($x \in H$). Cualquier operador A puede ser representado como una suma:

$$A = A_0 + A_1, \quad A_0 = \frac{1}{2}(A + A^*), \quad A_1 = \frac{1}{2}(A - A^*),$$

donde $A_0 = A_0^*$ es un operador autoconjugado y $A_1 = -A_1^*$, operador antisimétrico, para el cual en un espacio real se verifica $(A_1x, x) = -(x, A_1x) = -(A_1x, x)$, y, por consiguiente, $(A_1x, x) = 0$. Por eso, para cualquier operador A en el espacio real H se verifica la igualdad

$$(Ax, x) = (A_0x, x) \text{ para todos los } x \in H. \quad (6)$$

Hagamos uso de las siguientes desigualdades operacionales:

$$\begin{aligned} A \geq 0, & \text{ si } (Ax, x) \geq 0, & \text{ para todos los } x \in H; \\ A > 0, & \text{ si } (Ax, x) > 0, & \text{ para todos los } x \in H, x \neq 0; \\ A \geq \delta E, & \text{ si } (Ax, x) \geq \delta \|x\|^2, & \text{ para todos los } x \in H, \end{aligned} \quad (7)$$

donde E es un operador *unidad*.

La desigualdad

$$B \geq \alpha A$$

significa que queda cumplida la condición $B - \alpha A \geq 0$, es decir, $((B - \alpha A)x, x) \geq 0$ (para todos los $x \in H$).

Si en un espacio real $A \neq A^*$, entonces la desigualdad $A \geq 0$ ($A > 0$) será equivalente a la desigualdad $A_0 \geq 0$ ($A_0 > 0$), lo que se deduce de (6)

Sea A un operador positivo. Entonces, existe un operador inverso $A^{-1}: H \rightarrow H$, siendo $A^{-1} > 0$ para $A > 0$, $(A^{-1})^* = A^{-1}$ cuando $A^* = A$. En efecto, el operador A^{-1} existe, siempre que la ecuación $Ax = 0$ tiene solamente la ecuación trivial. Admitamos que $Ax = 0$ cuando $x \neq 0$; entonces $0 = (Ax, x)$ cuando $x \neq 0$, lo que contradice la condición $A > 0$, o bien $(Ax, x) > 0$ cuando $x \neq 0$. De este modo, si $A > 0$, entonces la ecuación $Ax = y$ tiene la solución única.

3. Valores propios del operador lineal. Sea A un operador autoconjugado en el espacio N -dimensional H provisto de producto escalar (\cdot, \cdot) . Analicemos un problema sobre los valores propios del operador A : se pide hallar los valores del parámetro λ (valores propios), para los cuales la ecuación homogénea

$$A\xi = \lambda\xi \quad (8)$$

tenga soluciones no triviales (vectores propios). He aquí algunas afirmaciones fundamentales del álgebra lineal sobre el problema de valores propios.

1) El operador autoconjugado A tiene N vectores propios ortonormalizados $\xi_1, \xi_2, \dots, \xi_N$:

$$(\xi_s, \xi_m) = \delta_{sm}. \quad \delta_{sm} = \begin{cases} 1, & s = m. \\ 0, & s \neq m. \end{cases} \quad (9)$$

2) Los valores propios correspondientes son reales y pueden disponerse en el orden de crecimiento de sus magnitudes absolutas:

$$0 \leq |\lambda_1| \leq |\lambda_2| \leq \dots \leq |\lambda_N|. \quad (10)$$

3) Si A es un operador positivo, entonces todos los valores propios $\{\lambda_k\}$ son positivos:

$$0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_N. \quad (11)$$

Efectivamente, $\lambda_s = (A\xi_s, \xi_s) / \|\xi_s\|^2 = (A\xi_s, \xi_s) > 0$, puesto que $\xi_s \neq 0$.

4) Un vector arbitrario $x \in H$ puede ser descompuesto según los vectores propios del operador $A = A^*$:

$$x = \sum_{k=1}^N c_k \xi_k, \quad c_k = (x, \xi_k), \quad (12)$$

quedando en este caso válida la igualdad

$$\|x\|^2 = \sum_{k=1}^N c_k^2. \quad (13)$$

En efecto, debido a la condición (9) de ortonormalidad del sistema $\{\xi_k\}$ tenemos

$$\begin{aligned} \|x\|^2 = (x, x) &= \left(\sum_{k=1}^N c_k \xi_k, \sum_{k'=1}^N c_{k'} \xi_{k'} \right) = \\ &= \sum_{k=1}^N \sum_{k'=1}^N c_k c_{k'} (\xi_k, \xi_{k'}) = \sum_{k=1}^N \sum_{k'=1}^N c_k c_{k'} \delta_{kk'} = \sum_{k=1}^N c_k^2. \end{aligned}$$

5) Si $A = A^* > 0$, entonces la solución de la ecuación $Ax = f$ puede ser representada en la forma

$$x = \sum_{k=1}^N \frac{f_k}{\lambda_k} \xi_k, \quad (14)$$

donde $f_k = (f, \xi_k)$ es el coeficiente de Fourier de la función f . Hagamos uso de las representaciones

$$x = \sum_{k=1}^N c_k \xi_k, \quad f = \sum_{k=1}^N f_k \xi_k$$

y escribamos

$$0 = Ax - f = \sum_{k=1}^N (\lambda_k c_k - f_k) \xi_k.$$

Al multiplicar esta igualdad escalarmente por ξ_k y teniendo en cuenta que $(\xi_k, \xi_{k'}) = \delta_{kk'}$, hallemos $0 = \lambda_k c_k - f_k$, es decir, $c_k = f_k / \lambda_k$.

6) La norma de un operador autoconjugado A es igual al módulo de su valor propio máximo:

$$\|A\| = \max_{1 \leq k \leq N} |\lambda_k| = |\lambda_N|. \quad (15)$$

Efectivamente, aprovechando (12), obtenemos

$$Ax = \sum_{k=1}^N c_k A \xi_k = \sum_{k=1}^N \lambda_k c_k \xi_k,$$

y, en virtud de (10) y (13), tenemos

$$\|Ax\|^2 = \sum_{k=1}^N \lambda_k^2 c_k^2 \leq \lambda_N^2 \sum_{k=1}^N c_k^2 = \lambda_N^2 \|x\|^2,$$

es decir, $\|A\| \leq |\lambda_N|$. Esta estimación se consigue. En efecto, para $x = \xi_N$ tenemos $\|Ax\|^2 = \|A\xi_N\|^2 = \|\lambda_N \xi_N\|^2 = |\lambda_N|^2$, puesto que $\|\xi_N\|^2 = 1$. De aquí precisamente se desprende que $\|A\| = |\lambda_N|$.

7) Si $A = A^*$, entonces

$$\|A\| = \sup_{\|x\|=1} |(Ax, x)|. \quad (16)$$

8) Si $A = A^* > 0$, entonces $\lambda_1 E \leq A \leq \lambda_N E$, o bien $\lambda_1 \|x\|^2 \leq (Ax, x) \leq \lambda_N \|x\|^2$, $\lambda_1 > 0$, $x \in H$. (17)

9) Si el operador A es positivo, será definido positivo, es decir, existe una constante $\delta > 0$ tal que de la condición $A > 0$ proviene la desigualdad $A \geq \delta E$. Para un operador autoconjugado esta propiedad se deduce de la propiedad 8). En el caso general representemos A en forma de una suma $A = A_0 + A_1$, donde $A_0 = A_0^* > 0$, $A_1 = -A_1^*$ es un operador antisimétrico. Puesto que $(A_1 x, x) = 0$, se tiene $(Ax, x) = (A_0 x, x) > 0$. Para A_0 es cierta la propiedad 8). Suponiendo $\lambda_1 = \lambda_1(A_0) = \delta > 0$, obtenemos $(A_0 x, x) = (Ax, x) \geq \delta \|x\|^2$ para todos los $x \in H$.

10) Si existe Q^{-1} , las desigualdades operacionales

$$C \geq 0, \quad Q^* C Q \geq 0 \quad (18)$$

serán equivalentes. Esto se deduce de la identidad

$$(Q^* C Q x, x) = (C Q x, Q x) = (C y, y),$$

donde $y = Qx$, $x = Q^{-1}y$.

11) Sean A_1 y A_2 los operadores en H autoconjugados, positivos y permutables:

$$A_1 = A_1^* > 0, \quad A_2 = A_2^* > 0, \quad A_1 A_2 = A_2 A_1. \quad (19)$$

En este caso los operadores A_1 y A_2 , la suma de ellos $A_1 + A_2$ y el producto $A_1 A_2$ tienen un sistema común de fun-

ciones propias $\{\xi_k\}$:

$$\begin{aligned} A_1 \xi_k &= \lambda_k^{(1)} \xi_k, & A_2 \xi_k &= \lambda_k^{(2)} \xi_k, \\ \lambda(A_1 + A_2) &= \lambda(A_1) + \lambda(A_2), \\ \lambda(A_1 A_2) &= \lambda(A_1) \lambda(A_2). \end{aligned}$$

12) Si $A = A^* > 0$, entonces el operador $A^{-1} = (A^{-1})^* > 0$ es también autoconjugado, tiene los mismos vectores propios que el operador A , y los valores propios $\lambda(A^{-1}) = 1/\lambda(A)$.

Efectivamente, de $A \xi_k = \lambda_k \xi_k$ proviene $\xi_k = \lambda_k A^{-1} \xi_k$, es decir, $(A^{-1}) \xi_k = (1/\lambda_k) \xi_k$. De aquí concluimos que las desigualdades $\lambda_1 E \leq A \leq \lambda_N E$ y $(1/\lambda_N) E \leq A^{-1} \leq (1/\lambda_1) E$ son equivalentes.

4. Problema generalizado sobre valores propios. Sea dado un operador autoconjugado positivo B . Introduzcamos un producto escalar nuevo $(x, y)_B = (Bx, y)$ y una norma $\|y\|_B = \sqrt{(By, y)}$. Un espacio H provisto de producto escalar $(x, y)_B$ recibe el nombre de *espacio energético* y se designa H_B .

Examinemos un problema generalizado sobre valores propios que consiste en buscar las soluciones no triviales v de la ecuación

$$Av = \mu Bv, \quad v \neq 0, \quad (20)$$

donde A es un operador autoconjugado positivo.

Supongamos que los operadores A y B están representados por las matrices respectivas $A = (a_{ij})$, $B = (b_{ij})$ ($i, j = 1, 2, \dots, N$). La ecuación operacional (20) puede escribirse en forma de un sistema de ecuaciones algebraicas lineales

$$\sum_{j=1}^N a_{ij} v^{(j)} = \mu \sum_{j=1}^N b_{ij} v^{(j)}, \quad i = 1, 2, \dots, N,$$

donde $v^{(1)}, \dots, v^{(N)}$ son componentes del vector v . Para determinar los valores propios se obtiene una ecuación algebraica de N -ésimo grado

$$\det(a_{ij} - \mu b_{ij}) = 0. \quad (21)$$

Para el problema (20) son justas las propiedades análogas a las del problema corriente sobre valores propios, a saber:

existen N vectores propios ortonormalizados en el sentido del producto escalar $(x, y)_B$

$$(v_k, v_m)_B = \delta_{km}, \quad k, m = 1, 2, \dots, N, \quad (22)$$

a los cuales corresponden los valores propios

$$0 < \mu_1 \leq \dots \leq \mu_N. \quad (23)$$

Por analogía con el p. 3 tenemos

$$x = \sum_{k=1}^N c_k v_k, \quad c_k = (x, v_k)_B, \\ \|x\|_B^2 = \sum_{k=1}^N c_k^2. \quad (24)$$

Se verifican las desigualdades operacionales

$$\mu_1 B \leq A \leq \mu_N B, \quad (25)$$

con la particularidad de que μ_N es la norma del operador A en H_B . Esto significa que

$$\|Ax\|_B \leq \|A\|_B \|x\|_B.$$

OBSERVACION. Las desigualdades

$$\gamma_1 B \leq A \leq \gamma_2 B, \quad \gamma_1 > 0, \quad (26)$$

$$\gamma_1 \leq \mu_h \leq \gamma_2, \quad h = 1, 2, \dots, N, \quad (27)$$

son equivalentes. En efecto, descompongamos un vector

arbitrario $x = \sum_{k=1}^N c_k v_k$, hallemos $(A - \gamma B)x = \sum_{k=1}^N c_k (\mu_k - \gamma) B v_k$ y el producto escalar

$$((A - \gamma B)x, x) = \sum_{k=1}^N c_k^2 (\mu_k - \gamma) (B v_k, v_k) = \sum_{k=1}^N (\mu_k - \gamma) c_k^2,$$

donde γ es uno de los números γ_1 ó γ_2 . Suponiendo $x = v_h$, determinemos $((A - \gamma B)v_h, v_h) = \mu_h - \gamma$. Sea $\gamma = \gamma_2$ y supongamos cumplida la condición $A \leq \gamma_2 B$; entonces $\mu_h \leq \gamma_2$. La afirmación recíproca es también cierta. Análogamente se realizan los razonamientos para $\gamma = \gamma_1$.

5. Espacios lineales de las funciones reticulares. Operadores de diferencias. En lo que sigue se examinarán sólo las funciones definidas sobre la red con nodos de números enteros:

$$\omega_N = \{i: i = 0, 1, \dots, N\}.$$

Al introducir en el segmento $0 \leq x \leq 1$ los nodos $x_i = ih$, $h = 1/N$ ($i = 0, 1, \dots, N$), obtendremos una red uniforme de paso h como una variedad de nodos $x_i = ih$ con índices de números enteros:

$$\omega_h = \{x_i = ih: i = 0, 1, \dots, N; h = 1/N\}.$$

El paso de una red a la otra es evidente y en algunos casos (bastante frecuentes) no las distinguiremos.

Denotemos con $\Omega_{N+1} = \{y_i, i = 0, 1, \dots, N\}$ el espacio de funciones reticulares definidas sobre la red ω_N , con $\dot{\Omega}_{N+1} = \{y_i, i = 0, 1, \dots, N; y_0 = 0, y_N = 0\}$ el subespacio de funciones reticulares que están definidas sobre la red ω_N y se reducen a cero en los nodos de frontera de la red ω_N : $y_0 = y_N = 0$. Las funciones de $\dot{\Omega}_{N+1}$ se designarán \dot{y} (i) = \dot{y}_i .

Veamos unos ejemplos de operadores de diferencias más simples. Para el operador de la diferencia derecha Δ tenemos

$$\Delta y_i = y_{i+1} - y_i, \quad i = 0, 1, \dots, N-1;$$

aquí el dominio de definición es Ω_{N+1} , el campo de valores está representado por el espacio $\Omega_N^* = \{y_i, i = 0, 1, \dots, N-1\}$ de N -ésima dimensión.

Para el operador de la diferencia izquierda ∇ tenemos

$$\nabla y_i = y_i - y_{i-1}, \quad i = 1, 2, \dots, N;$$

el dominio de definición es Ω_{N+1} , el campo de valores está representado por el espacio $\Omega_N^* = \{y_i, i = 1, 2, \dots, N\}$.

De la fórmula

$$\Delta^2 y_{i-1} = \Delta (\Delta y_{i-1}) = \Delta (\nabla y_i) = y_{i+1} - 2y_i + y_{i-1}$$

se ve que el operador de la segunda diferencia está definido para las funciones reticulares y_i con $i = 1, 2, \dots, N-1$, es decir, aplica Ω_{N+1} en el espacio $\Omega_{N-1} = \{y_i, i = 1, 2, \dots, N-1\}$. La misma propiedad posee el operador de

diferencias Λ :

$$\begin{aligned}\Lambda y_i &= b_i y_{i+1} - c_i y_i + a_i y_{i-1} = \\ &= b_i \Delta (\nabla y_i) - (b_i - a_i) (\nabla y_i) - (c_i - a_i - b_i) y_i, \\ & \qquad \qquad \qquad i = 1, 2, \dots, N-1,\end{aligned}$$

es decir, $\Lambda y_i \in \Omega_{N-1}$, si $y_i \in \Omega_{N+1}$, o bien, en la notación reducida, $\Lambda: \Omega_{N+1} \rightarrow \Omega_{N-1}$.

Analícemos un problema de contorno en diferencias

$$\begin{aligned}\Lambda y_i &= -f_i, \quad i = 1, 2, \dots, N-1, \\ y_0 &= \mu_1, \quad y_N = \mu_2\end{aligned}\quad (28)$$

y escribámosla en la forma matricial:

$$AY = \Phi, \quad (29)$$

donde $\Phi = (f_1 + a_1 \mu_1, f_2, \dots, f_{N-2}, f_{N-1} + b_{N-1} \mu_2)$ es el vector conocido e $Y = (y_1, y_2, \dots, y_{N-2}, y_{N-1})$ es un vector desconocido, ambos de dimensión $N-1$; A es una matriz tridiagonal cuadrada de dimensión $(N-1) \times (N-1)$:

$$A = - \begin{bmatrix} -c_1 & b_1 & \dots & 0 \\ a_2 & -c_2 & b_2 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & 0 & a_{N-1} & -c_{N-1} \end{bmatrix}.$$

Al comparar (28) y (29), vemos que se puede escribir

$$\begin{aligned}\tilde{\Lambda} y_i &= -\varphi_i, \quad i = 1, 2, \dots, N-1, \\ \tilde{\Lambda} y_1 &= -c_1 y_1 + b_1 y_2, \quad \varphi_1 = f_1 + a_1 \mu_1, \\ \tilde{\Lambda} y_i &= \Lambda y_i, \quad \varphi_i = f_i, \quad i = 2, 3, \dots, N-2, \\ \tilde{\Lambda} y_{N-1} &= a_{N-1} y_{N-2} - c_{N-1} y_{N-1}, \quad \varphi_{N-1} = f_{N-1} + b_{N-1} \mu_2.\end{aligned}\quad (28')$$

El operador de diferencias $\tilde{\Lambda}$ aplica Ω_{N-1} en Ω_{N-1} . No es difícil observar que $\Lambda \tilde{y}_i = \tilde{\Lambda} y_i$. En lugar de (28') obtendremos

$$\Lambda \tilde{y}_i = -\varphi_i, \quad i = 1, 2, \dots, N-1.$$

Introduzcamos ahora el operador A correspondiente a la matriz (29), suponiendo

$$Ay_i = -\tilde{\Lambda}y_i = -\dot{\Lambda}y_i, \quad i = 1, 2, \dots, N-1.$$

En tal caso, en vez del problema de contorno en diferencias (28) obtendremos una ecuación operacional

$$Ay = \varphi,$$

donde $A: \Omega_{N-1} \rightarrow \Omega_{N-1}$, $\varphi \in \Omega_{N-1}$, es decir, el operador A actúa de Ω_{N-1} en Ω_{N-1} . Es evidente que A será un operador lineal. Ha de notarse que también puede considerarse (teniendo en cuenta que $Ay = -\Lambda\dot{y}$), que A aplica $\dot{\Omega}_{N-1}$ en Ω_{N-1} .

En el espacio $H = \Omega_{N-1}$ se puede introducir un producto escalar

$$(y, v) = \frac{1}{N} \sum_{i=1}^{N-1} y_i v_i$$

y una norma

$$\|y\| = \sqrt{(y, y)}.$$

Si se estudian el segundo ($\kappa_1 = \kappa_2 = 1$) o el tercero ($\kappa_1 \neq 0$, $\kappa_2 \neq 0$) de los problemas de contorno (véase (1) del § 3), la matriz A será cuadrada de dimensión $(N+1) \times (N+1)$ y el operador A se definirá de la manera siguiente:

$$Ay_i = -\Lambda y_i = -(b_i y_{i+1} - c_i y_i + a_i y_{i-1}), \\ i = 1, 2, \dots, N-1,$$

$$Ay_0 = -(\kappa_1 y_1 - y_0), \quad Ay_N = -(y_N - \kappa_2 y_{N-1}).$$

En este caso el operador A aplica el espacio de funciones reticulares $H = \Omega_{N+1}$ en sí mismo: $A: H \rightarrow H$.

En adelante se analizará el primer problema de contorno para una ecuación en diferencias de segundo orden; en este caso, según lo mostrado más arriba, $H = \Omega_{N-1}$.

6. Fórmulas de Green de diferencias. Examinemos un operador de diferencias L :

$$Ly_i = b_i y_{i+1} - c_i y_i + a_i y_{i-1}, \quad i = 1, \dots, N-1. \quad (30)$$

Si $b_i \neq a_{i+1}$, la matriz correspondiente no será simétrica. Es simétrica sólo en el caso

$$b_i = a_{i+1}, \quad i = 1, 2, \dots, N-1. \quad (31)$$

Al tomar en consideración esta condición, escribimos Ly_i en la forma siguiente:

$$\begin{aligned} Ly_i &= a_{i+1}y_{i+1} - c_i y_i + a_i y_{i-1} = \\ &= a_{i+1} (y_{i+1} - y_i) - a_i (y_i - y_{i-1}) - (c_i - a_i - \\ &- a_{i+1}) y_i = a_{i+1} \nabla y_{i+1} - a_i \nabla y_i - (c_i - a_i - \\ &- a_{i+1}) y_i = \Delta (a_i \nabla y_i) - (c_i - a_i - a_{i+1}) y_i. \end{aligned} \quad (32)$$

Dividamos el segmento $[0, 1]$ con los puntos x_i en N partes iguales, hagamos $y(x_i) = y_i = y(i)$ e introduzcamos las designaciones que siempre se usarán en lo sucesivo:

$$\begin{aligned} h &= \frac{1}{N}, \quad x_i = ih, \quad i = 0, 1, \dots, N, \quad x_0 = 0, \quad x_N = 1, \\ y_{x,i} &= \frac{\Delta y_i}{h} = \frac{y_{i+1} - y_i}{h}, \quad y_{\bar{x},i} = \frac{\nabla y_i}{h} = \frac{y_i - y_{i-1}}{h}, \quad (33) \\ y_{\bar{x}\bar{x},i} &= y_{\bar{x}\bar{x}}(i) = \frac{y_{i-1} - 2y_i + y_{i+1}}{h^2} = \frac{\Delta(\nabla y_i)}{h^2}. \end{aligned}$$

Dividamos la expresión (32) por h^2 y obtendremos un operador de diferencias

$$\begin{aligned} \Lambda y_i &= (ay_{\bar{x}})_{x,i} - d_i y_i, \\ d_i &= \frac{1}{h^2} (c_i - a_i - a_{i+1}), \quad i = 1, \dots, N-1. \end{aligned} \quad (34)$$

En el § 1 se ha obtenido la fórmula de sumación por partes

$$\sum_{i=0}^{N-1} y_i \Delta v_i = - \sum_{i=1}^N v_i \nabla y_i + (yv)_N - (yv)_0. \quad (35)$$

Haciendo uso de las designaciones (33), escribamos esta fórmula en la forma

$$\sum_{i=0}^{N-1} y_i v_{x,i} h = - \sum_{i=1}^N v_i y_{\bar{x},i} h + (yv)_N - (yv)_0, \quad (36)$$

puesto que $\sum_{i=0}^{N-1} y_i \Delta v_i = \sum_{i=0}^{N-1} y_i \left(\frac{\Delta v_i}{h} \right) h = \sum_{i=0}^{N-1} y_i v_{x,i} h.$

Para que la exposición ulterior sea más cómoda, introduzcamos en el primer miembro de (36) la sumación entre $i = 1$ e $i = N - 1$; esto nos conduce a la fórmula

$$\sum_{i=1}^{N-1} y_i v_{x,i} h = - \sum_{i=1}^N v_i y_{\bar{x},i} h + (yv)_N - y_0 v_1. \quad (37)$$

Sustituyamos aquí $v_i = a_i z_{\bar{x},i}$; se obtendrá

$$\sum_{i=1}^{N-1} y_i (az_{\bar{x}})_{x,i} h = - \sum_{i=1}^N a_i y_{\bar{x},i} z_{\bar{x},i} h + (ayz_{\bar{x}})_N - y_0 (az_{\bar{x}})_1. \quad (38)$$

Esta es la *primera fórmula de Green de diferencias*. Cambiemos de lugar en ella y_i y z_i :

$$\sum_{i=1}^{N-1} z_i (ay_{\bar{x}})_{x,i} h = - \sum_{i=1}^N a_i z_{\bar{x},i} y_{\bar{x},i} h + (ay_{\bar{x}}z)_N - z_0 (ay_{\bar{x}})_1. \quad (38')$$

Al sustraer (38') de (38) obtenemos la *segunda fórmula de Green de diferencias*

$$\sum_{i=1}^{N-1} y_i (az_{\bar{x}})_{x,i} h = \sum_{i=1}^{N-1} z_i (ay_{\bar{x}})_{x,i} h + a_N (yz_{\bar{x}} - zy_{\bar{x}})_N - (y_0 (az_{\bar{x}})_1 - z_0 (ay_{\bar{x}})_1). \quad (39)$$

Si están cumplidas las condiciones

$$y_0 = z_0 = 0, \quad y_N = z_N = 0, \quad (40)$$

es decir, si $y = \dot{y}$, $z = \dot{z} \in \dot{\Omega}_{N+1}$, entonces en el segundo miembro de la igualdad (39) dos últimos sumandos se anulan y

$$\sum_{i=1}^{N-1} \dot{y}_i (a\dot{z}_{\bar{x}})_{x,i} h = \sum_{i=1}^{N-1} \dot{z}_i (a\dot{y}_{\bar{x}})_{x,i} h. \quad (41)$$

Al sustraer de ambos miembros de la identidad (41) la suma

$\sum_{i=1}^{N-1} d_i \dot{y}_i \dot{z}_i h$, obtenemos la *segunda fórmula de Green para*

$y, z \in \dot{\Omega}_{N+1}$:

$$\sum_{i=1}^{N-1} \dot{y}_i \Lambda \dot{z}_i h = \sum_{i=1}^{N-1} \dot{z}_i \Lambda \dot{y}_i h \quad (42)$$

para el operador de diferencias

$$\Lambda \dot{y}_i = (a \dot{y}_{\bar{x}})_{x,i} - d_i \dot{y}_i, \text{ cualquiera que sea } \dot{y} \in \dot{\Omega}_{N+1}. \quad (43)$$

Sea $H = \Omega_{N-1}$ un espacio de funciones reticulares y_i , prefijadas para $i = 1, 2, \dots, N-1$, con el producto escalar

$$(y, v) = \sum_{i=1}^{N-1} y_i v_i h$$

y la norma

$$\|y\| = \sqrt{(y, y)}.$$

Introduzcamos el operador A :

$$Ay = -\Lambda \dot{y}, \quad y \in H. \quad (44)$$

Entonces la segunda fórmula de Green puede anotarse en la forma

$$(y, Az) = (Ay, z). \quad (45)$$

Esta fórmula expresa la propiedad de autoconjugación del operador A : $A^* = A$ y, por lo tanto, $\Lambda^* = \Lambda$. Cuando $\dot{z} = \dot{y} \in \dot{\Omega}_{N+1}$, la primera fórmula de Green (38) nos da:

$$-\sum_{i=1}^{N-1} \dot{y}_i (a \dot{y}_{\bar{x}})_{x,i} h = \sum_{i=1}^N a_i (\dot{y}_{\bar{x},i})^2 h > 0$$

para $\dot{y}_i \neq 0, a_i > 0, \quad (46)$

(ya que $\dot{y}_0 = \dot{y}_N = 0$, (46) puede ser igual a cero sólo en el caso en que $\dot{y}_i = 0$ ($i = 1, \dots, N-1$)). Teniendo presente la definición del operador A , hallamos

$$(Ay, y) = \sum_{i=1}^N a_i (y_{\bar{x},i})^2 h + \sum_{i=1}^N d_i y_i^2 h > 0, \quad a_i > 0,$$

$d_i \geq 0. \quad (47)$

De este modo, el operador de diferencias A , definido por las fórmulas (43), (44), es autoconjugado y positivo: $A = A^* > 0$, siempre que

$$a_i > 0, \quad d_i \geq 0, \quad i = 1, 2, \dots, N-1, \quad a_N > 0. \quad (48)$$

7. Condición de autoconjugación del operador de diferencias de segundo orden. Nos hemos convencido de que la condición (31) es suficiente para que el operador de diferencias (30) sea autoconjugado en el espacio $H = \dot{\Omega}_{N+1}$. Mostremos que la condición (31) es necesaria para que sea autoconjugado L . Representemos L en forma de una suma:

$$\begin{aligned} Ly_i &= L_1 y_i + L_2 y_i, \\ L_1 y_i &= a_{i+1} (y_{i+1} - y_i) - a_i (y_i - y_{i-1}) - (c_i - a_i - b_i) y_i, \\ L_2 y_i &= (b_i - a_{i+1}) y_{i+1}. \end{aligned}$$

Como ya se ha mostrado en el punto antecedente, el operador $L_1 y_i = h^2 \Lambda y_i$, $\Lambda y_i = (ay_x)_{x,i} - d_i y_i$, es autoconjugado en el espacio $H = \dot{\Omega}_{N+1}$ ó en $H = \Omega_{N-1}$ con el producto escalar $(y, v) = \sum_{i=1}^{N-1} y_i v_i h$. Por eso podemos escribir:

$$\begin{aligned} \left(\frac{1}{h^2} L \dot{y}, \dot{v} \right) - \left(\dot{y}, \frac{1}{h^2} L \dot{v} \right) &= \\ &= (\Lambda \dot{y}, \dot{v}) - (\dot{y}, \Lambda \dot{v}) + \left(\frac{1}{h^2} L_2 \dot{y}, \dot{v} \right) - \left(\dot{y}, \frac{1}{h^2} L_2 \dot{v} \right) = \\ &= \sum_{i=1}^N \frac{1}{h^2} (b_i - a_{i+1}) (y_{i+1} v_i - y_i v_{i+1}) h. \end{aligned}$$

De aquí se ve que $(L \dot{y}, \dot{v}) = (\dot{y}, L \dot{v})$, es decir, $L = L^*$ sólo a la condición de que

$$\sum_{i=1}^{N-1} (b_i - a_{i+1}) (y_{i+1} v_i - y_i v_{i+1}) h = 0. \quad (49)$$

Por ser arbitrarias y_i y v_i , podemos tomar $y_i = \delta_{i, i_0+1}$, $v_i = \delta_{i, i_0}$, donde i_0 es un nodo fijo cualquiera ($i_0 = 1, 2, \dots, N-1$), mientras que δ_{i, i_0} es el símbolo de Kronecker.

Obtenemos, pues, $y_{i+1}v_i - y_iv_{i+1} = \delta_{i,i_0}$, y la condición (49) nos da $b_{i_0} = a_{i_0+1}$. Con esto queda demostrada la necesidad de la condición (31).

Se debe notar que la ecuación

$$Ly_i = -f_i \quad (50)$$

puede ser reducida a la forma

$$\tilde{L}y_i = \Delta (A_i \nabla y_i) - D_i y_i = -F_i, \quad (51)$$

donde \tilde{L} es un operador autoconjugado. En efecto, multipliquemos ambos miembros de la ecuación (50) por $\mu_i \neq 0$:

$$\tilde{L}y_i = \mu_i a_i y_{i-1} - \mu_i c_i y_i + b_i \mu_i y_{i+1} = -\mu_i f_i$$

y exijamos que para la ecuación obtenida se cumpla la condición (31), es decir,

$$b_i \mu_i = (\mu a)_{i+1} = a_{i+1} \mu_{i+1} = A_{i+1}.$$

De aquí obtenemos $\mu_{i+1} = \left(\frac{b_i}{a_{i+1}} \right) \mu_i = \mu_i \prod_{k=1}^i b_k/a_{k+1}$ y la ecuación (51), donde $A_i = a_i \mu_i$, $D_i = \mu_i (c_i - a_i - b_i)$, $F_i = -\mu_i f_i$.

8. Valores propios del operador de diferencias de segundo orden. Examinemos un problema de diferencias sobre valores propios:

$$(ay_{\bar{x}})_{x,i} - d_i y_i + \lambda y_i = 0, \quad i = 1, 2, \dots, N-1$$

$$y_0 = y_N = 0, \quad (52)$$

o bien $Ay = \lambda y$, $y \in \Omega_{N-1}$, donde A se determina por la igualdad (44). El operador A es autoconjugado y positivo, razón por la cual a él se refiere todo lo dicho en el p. 4.

En el caso más simple, $a_i = 1$, $d_i = 0$, los valores propios y los vectores propios pueden hallarse en la forma explícita. Así pues, se requiere encontrar soluciones no triviales de la ecuación homogénea con condiciones de contorno homogéneas

$$y_{\bar{x}x,i} + \lambda y_i = 0, \quad i = 1, 2, \dots, N-1, \quad hN = 1,$$

$$y_0 = 0, \quad y_N = 0, \quad y_i \neq 0. \quad (53)$$

Escribamos la ecuación (53) en la forma siguiente

$$y_{i-1} - 2 \cos \alpha y_i + y_{i+1} = 0, \quad 2 \cos \alpha = 2 - \lambda h^2. \quad (54)$$

La solución general de esta ecuación tiene por expresión

$$y_i = c_1 \cos i\alpha + c_2 \operatorname{sen} i\alpha. \quad (55)$$

Exigimos que se cumplan las condiciones de contorno: $y_0 = c_1 = 0$, $y_N = c_2 \operatorname{sen} N\alpha = 0$. Como se busca una solución no trivial, entonces $c_2 \neq 0$ y $\operatorname{sen} N\alpha = 0$, es decir, $N\alpha = m\pi$ ($m = 0, 1, 2, \dots$), $\alpha = \alpha_m = m\pi/N = = m\pi h$. De la relación $2 \cos \alpha = 2 - \lambda h^2$ encontramos

$$\lambda h^2 = 2(1 - \cos \alpha) = 4 \operatorname{sen}^2 \frac{\alpha}{2},$$

$$\lambda = \lambda_m = \frac{4}{h^2} \operatorname{sen}^2 \frac{m\pi h}{2}. \quad (56)$$

A este valor de λ_m le corresponde una función propia

$$y_m(i) = c \operatorname{sen} \pi m x_i, \quad c \neq 0, \quad x_i = ih, \quad i = 0, 1, 2, \dots, N \quad (57)$$

definida con una exactitud de hasta un factor constante arbitrario. No es difícil notar que

$$\begin{aligned} y_N(i) &= c \operatorname{sen} \pi N x_i = c \operatorname{sen} \pi i = 0, \quad i = 0, 1, 2, \dots, \\ y_{N+1}(i) &= c \operatorname{sen} \pi (N+1) x_i = c \operatorname{sen} [\pi N x_i + \pi x_i] = \\ &= c \operatorname{sen} \pi x_i \cos \pi i = (-1)^i y_1(i), \\ y_{N+m+1}(i) &= (-1)^i y_m(i), \quad m = 1, 2, \dots, N-1. \end{aligned}$$

Por consiguiente, para $m < N$, sólo las funciones $y_m(i)$ son linealmente independientes. Así pues, se ha encontrado la solución no trivial (funciones propias $y_m(i)$ que corresponden a los valores propios λ_m).

Elijamos el factor c de un modo tal que la norma de las funciones $y_m(i)$ sea igual a la unidad: $\|y_m(i)\| = = c \|\operatorname{sen} \pi m x_i\| = 1$, $c > 0$. Con este fin se debe calcular

$$\|\operatorname{sen} \pi m x_h\|^2 = \sum_{h=1}^{N-1} h \operatorname{sen}^2 \pi m x_h = \frac{1}{2} \sum_{h=1}^{N-1} h (1 - \cos 2\pi m x_h).$$

Al denotar $\alpha = 2\pi mh$ y sustituir $\cos 2\pi mx_h = \cos \alpha k = \operatorname{Re} e^{i\alpha k}$, llegamos a que

$$\sum_{k=1}^{N-1} h \cos 2\pi mx_h = \operatorname{Re} \sum_{k=1}^{N-1} h e^{i\alpha k} = h \operatorname{Re} \frac{e^{i\alpha} - e^{i\alpha N}}{1 - e^{i\alpha}} = -h,$$

$$\| \operatorname{sen} \pi mx_h \|^2 = \frac{(N-1)h}{2} - \frac{1}{2} \sum_{k=1}^{N-1} h \cos 2\pi mx_h = \frac{Nh}{2} = \frac{1}{2},$$

$$\| \operatorname{sen} \pi mx \| = 1/\sqrt{2};$$

por consiguiente, $c = \sqrt{2}$. De este modo, la función

$$y_m(i) = \sqrt{2} \operatorname{sen} \pi mx_i \quad (58)$$

está normalizada hacia la unidad.

Las funciones propias $y_s(i)$ e $y_m(i)$, correspondientes a los diferentes valores propios λ_s y λ_m , son ortogonales en el sentido del producto escalar

$$(y, v) = \sum_{i=1}^{N-1} y_i v_i h.$$

El problema (53) constituye un caso particular del problema (8) con el operador $Ay(i) = -\ddot{y}_{xx}(i)$. Dicho operador es, evidentemente, autoconjugado y positivo, puesto que

$$(Ay, y) = \sum_{i=1}^{N-1} (y_{x,i})^2 h > 0.$$

Por esta razón todo lo dicho en el p. 3 queda vigente también en el caso dado.

Los valores propios λ_s crecen a medida que crece s , puesto que $\operatorname{sen} \frac{\pi h}{2} s < \operatorname{sen} \frac{\pi h}{2} (s+1) < 1$ para $s \leq N$. El valor propio mínimo es $\lambda_1 = \frac{4}{h^2} \operatorname{sen}^2 \frac{\pi h}{2}$. El valor propio máximo es igual a $\lambda_{N-1} = \frac{4}{h^2} \cos^2 \frac{\pi h}{2}$, ya que $\operatorname{sen} \frac{\pi h}{2} (N-1) = \operatorname{sen} \left(\frac{\pi}{2} - \frac{\pi h}{2} \right) = \cos \frac{\pi h}{2}$.

Escribiendo λ_1 en la forma $\lambda_1 = \pi^2 \left(\frac{\operatorname{sen} \xi}{\xi} \right)^2$, $\xi = \pi h/2 = \pi h/2 \leq \pi/4$, y teniendo presente que $\operatorname{sen} \xi/\xi$ decrece y tiene mínimo para $\xi = \pi/4$, obtenemos $\lambda_1 \geq 8$ para $h \leq 1/2$.

Para λ_{N-1} tenemos una estimación $\lambda_{N-1} < 4/h^2$, y, por consiguiente,

$$8 < \lambda_k < 4/h^2, \quad k = 1, 2, \dots, N - 1.$$

§ 5. Principio del máximo para las ecuaciones en diferencias

1. Principio del máximo y sus corolarios. Para las ecuaciones en diferencias de segundo orden con coeficientes positivos

$$Ly_i = a_i y_{i-1} - c_i y_i + b_i y_{i+1} = -\varphi_i, \\ i = 1, 2, \dots, N - 1, \quad y_0 = \mu_1, \quad y_N = \mu_2, \quad (1) \\ a_i > 0, \quad b_i > 0, \quad c_i \geq a_i + b_i, \quad i = 1, 2, \dots, N - 1 \quad (2)$$

tiene lugar el siguiente principio del máximo.

TEOREMA 1 (principio del máximo). *Supongamos que un operador de diferencias L está definido por las fórmulas (1), (2). Si para una función y_i , prefijada sobre la red $\bar{\omega}$ y diferente de una constante ($1 \leq i \leq N - 1$), se cumple la condición $Ly_i \geq 0$ ($Ly_i \leq 0$) para todo $i = 1, 2, \dots, N - 1$, entonces dicha función no puede tomar el valor positivo máximo (negativo mínimo) en los nodos interiores de la red.*

DEMOSTRACION. Sea $Ly_i \geq 0$ ($i = 1, 2, \dots, N - 1$). Supongamos que el teorema no es cierto e y_i alcanza su valor positivo máximo en un nodo interior $i = i_*$, $1 \leq i_* \leq N - 1$: $y_{i_*} = \max_{0 \leq i \leq N} y_i = M_0 > 0$. Como $y_i \neq \text{const}$, se encontrará un nodo interior i_0 (i_0 puede coincidir con i_*) en el cual $y_{i_0} = y_{i_*} = M_0 > 0$, y en uno de los nodos vecinos, por ejemplo, en el nodo $i = i_0 - 1$, se verifica la desigualdad rigurosa $y_{i_0-1} < y_{i_0}$. Escribamos la expresión para Ly_i en la forma $Ly_i = b_i (y_{i+1} - y_i) - a_i (y_i - y_{i-1}) - (c_i - a_i - b_i) y_i$. En el nodo $i = i_0$

tenemos

$$Ly_{i_0} = b_{i_0}(y_{i_0+1} - y_{i_0}) - a_{i_0}(y_{i_0} - y_{i_0-1}) - (c_{i_0} - a_{i_0} - b_{i_0})y_{i_0} < 0,$$

lo que contradice la suposición $Ly_i \geq 0$ para cualquiera de los $i = 1, 2, \dots, N-1$, incluso para $i = i_0$. La primera afirmación del teorema queda demostrada. La segunda afirmación se demuestra análogamente (basta sustituir y_i por $-y_i$ y aprovechar la afirmación que acabamos de demostrar).

COROLARIO 1. Si se cumplen las condiciones (2), es decir, si $Ly_i \leq 0$ ($i = 1, 2, \dots, N-1$), $y_0 \geq 0$, $y_N \geq 0$, entonces $y_i \geq 0$ ($i = 0, 1, \dots, N$).

$$\text{Si } Ly_i \geq 0, \quad y_0 \leq 0, \quad y_N \leq 0, \quad \text{entonces } y_i \leq 0 \\ (i = 0, 1, \dots, N).$$

DEMOSTRACION. Supongamos que $Ly_i \leq 0$ e $y_i < 0$ por lo menos en uno de los nodos interiores $i = i_*$; entonces y_i alcanza el valor negativo mínimo en el nodo interior, lo que es imposible en virtud del principio del máximo.

COROLARIO 2. Si $\varphi_i \geq 0$, $\mu_1 \geq 0$, $\mu_2 \geq 0$, entonces la solución del problema (1)-(2) es no negativa: $y_i \geq 0$ ($i = 0, 1, \dots, N$).

COROLARIO 3. Si quedan cumplidas las condiciones (2), el problema

$$Ly_i = 0, \quad i = 1, 2, \dots, N-1, \quad y_0 = 0, \quad y_N = 0 \quad (3)$$

tiene sólo una solución trivial y el problema (1), (2) es resoluble unívocamente, cualesquiera que sean φ_i , μ_1 , μ_2 .

DEMOSTRACION. Suponiendo que la solución y_i del problema (3) es diferente de cero por lo menos en un solo punto $i = i_*$, llegamos a una contradicción con el principio del máximo: si $y_{i_*} > 0$ ($y_{i_*} < 0$), entonces y_i alcanza el valor máximo positivo (mínimo negativo) en cierto punto interior $i = i_0$, lo que es imposible. Por consiguiente, $y_i \equiv 0$.

TEOREMA 2. (teorema de comparación). Supongamos que y_i es la solución del problema (1), (2) e \bar{y}_i , la solución del problema

$$L\bar{y}_i = -\bar{\varphi}_i, \quad i = 1, 2, \dots, N-1, \quad \bar{y}_0 = \bar{\mu}_1, \\ \bar{y}_N = \bar{\mu}_2$$

y, además, admitamos cumplidas las condiciones

$$|\varphi_i| \leq \bar{\varphi}_i, \quad |\mu_1| \leq \bar{\mu}_1, \quad |\mu_2| \leq \bar{\mu}_2.$$

En este caso resulta válida la estimación

$$|y_i| \leq \bar{y}_i \text{ para todo } i = 0, 1, \dots, N.$$

DEMOSTRACION. En virtud del corolario 2, tenemos $\bar{y}_i \geq 0$. Para la diferencia $\bar{y}_i - y_i$ y para la suma $\bar{y}_i + y_i$ obtenemos una ecuación del tipo (1) con los segundos miembros $\bar{\varphi}_i - \varphi_i \geq 0$, $\bar{\mu}_1 - \mu_1 \geq 0$, $\bar{\mu}_2 - \mu_2 \geq 0$ y $\bar{\varphi}_i + \varphi_i \geq 0$, $\bar{\mu}_1 + \mu_1 \geq 0$, $\bar{\mu}_2 + \mu_2 \geq 0$, respectivamente. Puesto que $\bar{\varphi}_i \pm \varphi_i \geq 0$ y $\bar{\mu}_\alpha \pm \mu_\alpha \geq 0$ ($\alpha = 1, 2$), entonces, debido al corolario 2, $\bar{y}_i - y_i \geq 0$, $\bar{y}_i + y_i \geq 0$, de lo que se deduce que $-\bar{y}_i \leq y_i \leq \bar{y}_i$, $|y_i| \leq \bar{y}_i$, lo que se trataba de demostrar.

La función \bar{y}_i se denomina *mayorante* para la solución del problema (1), (2).

2. Estimación de la solución del problema de contorno. Representemos la solución del problema de contorno (1), (2) en forma de una suma $y_i = y_i^{(1)} + y_i^{(2)}$, donde $y_i^{(1)}$ es la solución de la ecuación no homogénea con condiciones de contorno homogéneas:

$$Ly_i = -\varphi_i, \quad i = 1, 2, \dots, N-1, \quad y_0 = y_N = 0, \quad (4)$$

mientras que $y_i^{(2)}$ representa la solución de la ecuación homogénea con condiciones de contorno no homogéneas

$$Ly_i = 0, \quad i = 1, 2, \dots, N-1, \quad y_0 = \mu_1, \quad y_N = \mu_2. \quad (5)$$

Demostremos que para $y_i^{(2)}$ es justa la estimación

$$\max_{0 \leq i \leq N} |y_i^{(2)}| \leq \max(|\mu_1|, |\mu_2|). \quad (6)$$

Sea \bar{y}_i la solución del problema

$$L\bar{y}_i = 0, \quad i = 1, 2, \dots, N-1,$$

$$\bar{y}_0 = \bar{y}_N = \bar{\mu}, \quad \bar{\mu} = \max(|\mu_1|, |\mu_2|).$$

Entonces, de acuerdo con el teorema de comparación, $|y_i^{(2)}| \leq |\bar{y}_i|$, mientras que, en virtud del principio del máximo $\max_{0 \leq i \leq N} |\bar{y}_i| \leq \mu$, puesto que $\bar{y}_i \geq 0$ puede alcanzar el valor positivo máximo sólo en la frontera, es decir, cuando $i = 0$ o $i = N$.

No es difícil demostrar que la magnitud $\max_{0 \leq i \leq N} |y_i|$ es una norma. La norma suele designarse con el símbolo $\|y\|_c$. De este modo, hemos obtenido la estimación $\|y^{(2)}\|_c \leq \leq \max(|\mu_1|, |\mu_2|)$.

TEOREMA 3. *Supongamos que están cumplidas las condiciones*

$$|a_i| > 0, \quad |b_i| > 0, \quad \bar{d}_i = |c_i| - |a_i| - |b_i| > 0 \\ i = 1, 2, \dots, N-1. \quad (7)$$

Entonces, para la solución del problema (4) es justa la estimación

$$\|y\|_c \leq \|\varphi/\bar{d}\|_c. \quad (8)$$

DEMOSTRACION. Con el fin de demostrar la citada afirmación escribamos (4) en la forma

$$c_i y_i = a_i y_{i-1} + b_i y_{i+1} + \varphi_i. \quad (4')$$

Supongamos que $|y_i|$ alcanza su valor máximo $|y_{i_0}| > 0$ cuando $i = i_0$ ($0 < i_0 < N$), de suerte que $|y_i| \leq \leq |y_{i_0}|$ para cualquier $i = 0, 1, \dots, N$. Entonces, de (4') se deduce para $i = i_0$:

$$|c_{i_0}| |y_{i_0}| = |a_{i_0} y_{i_0-1} + b_{i_0} y_{i_0+1} + \varphi_{i_0}| \leq |a_{i_0}| |y_{i_0-1}| + \\ + |b_{i_0}| |y_{i_0+1}| + |\varphi_{i_0}| \leq \\ \leq (|a_{i_0}| + |b_{i_0}|) |y_{i_0}| + |\varphi_{i_0}|, \\ (|c_{i_0}| - |a_{i_0}| - |b_{i_0}|) |y_{i_0}| \leq |\varphi_{i_0}|, \quad |y_{i_0}| \leq \frac{|\varphi_{i_0}|}{\bar{d}_{i_0}} \leq \|\varphi/\bar{d}\|_c.$$

Con esto queda demostrada la estimación (8).

OBSERVACION. Si la condición $\bar{d}_i = c_i - a_i - b_i > 0$ o $\bar{d}_i = |c_i| - |a_i| - |b_i| > 0$ no se cumple, por ejemplo, $\bar{d}_i = c_i - a_i - b_i \geq 0$, $a_i > 0$, $b_i > 0$, $i = 1, 2, \dots, N-1$,

(9)

es decir, d_i puede reducirse a cero en ciertos nodos, entonces el teorema 3 no puede ser aplicado. En este caso, con el fin de estimar la solución y_i del problema (4), se puede proceder de la manera siguiente. Representemos y_i en forma de una suma $y_i = v_i + w_i$, donde w_i es la solución del problema

$$\begin{aligned} Lw_i &= b_i (w_{i+1} - w_i) - a_i (w_i - w_{i-1}) = -\varphi_i, \\ i &= 1, 2, \dots, N-1, \quad w_0 = w_N = 0. \end{aligned} \quad (10)$$

Entonces, v_i se determina partiendo de las condiciones

$$\begin{aligned} Lv_i &= b_i (v_{i+1} - v_i) - a_i (v_i - v_{i-1}) - d_i v_i = -d_i w_i, \\ i &= 1, 2, \dots, N-1, \quad v_0 = v_N = 0. \end{aligned} \quad (11)$$

Se puede convencerse de esto sumando término a término las ecuaciones (10) y (11). La función w_i puede estimarse inmediatamente (véase el cap. IV, § 3), al escribirla en la forma explícita, mientras que para la estimación de v_i necesitaremos el

TEOREMA 4. *Para resolver el problema (11) bajo las condiciones (9) es válida la estimación*

$$\|v\|_c \leq \|w\|_c. \quad (12)$$

DEMOSTRACION. Si $d_i \equiv 0$, entonces, en virtud del corolario 3, $v_i \equiv 0$, y la estimación (12) queda cumplida. Sea $d_i \neq 0$ siquiera en un solo punto. Construyamos la mayorante \bar{v}_i como una solución del problema

$$L\bar{v}_i = -d_i |w_i|, \quad i = 1, 2, \dots, N-1, \quad v_0 = \bar{v}_N = 0.$$

Supongamos que $\bar{v}_i \geq 0$ alcanza su valor máximo para $i = i_0$; entonces $\bar{v}_{i_0+1} - \bar{v}_{i_0} \leq 0$, $\bar{v}_{i_0} - \bar{v}_{i_0-1} \geq 0$, y de (4) proviene

$$\begin{aligned} d_{i_0} \bar{v}_{i_0} &\leq -b_{i_0} (\bar{v}_{i_0+1} - \bar{v}_{i_0}) + a_{i_0} (\bar{v}_{i_0} - \bar{v}_{i_0-1}) + d_{i_0} v_{i_0} = \\ &= d_{i_0} |w_{i_0}|. \end{aligned}$$

Si $d_{i_0} > 0$, entonces $\bar{v}_{i_0} < |w_{i_0}|$ y obtenemos en seguida la estimación (12), puesto que $|v_i| \leq \bar{v}_i$. Si $d_{i_0} = 0$, la ecuación (11) tomará la forma $-b_{i_0} (\bar{v}_{i_0+1} - \bar{v}_{i_0}) + a_{i_0} (\bar{v}_{i_0} - \bar{v}_{i_0-1}) = 0$, y de esta última se deduce que $\bar{v}_{i_0+1} = \bar{v}_{i_0-1} = \bar{v}_{i_0}$. Por

cuanto $\bar{v}_i \neq \text{const}$, existe tal punto $i = i_1$, en el cual $\bar{v}_{i_1} = \bar{v}_{i_1}$, y en el punto vecino, por ejemplo, $i = i_1 + 1$, $\bar{v}_{i_1+1} < \bar{v}_{i_1}$; entonces aquí $d_{i_1} \neq 0$ y obtenemos, pues, el caso analizado más arriba: $\bar{v}_{i_1} = \|\bar{v}\|_c \leq |w_{i_1}| \leq \|w\|_c$.

3. **Estimación de la solución de una ecuación en diferencias con ayuda de las fórmulas de factorización.** Para el caso en que $b_i = a_{i+1}$, es decir, cuando el operador Ly_i sea autoconjugado, la solución del problema (4) puede ser estimada con ayuda de las fórmulas de factorización derecha. Es más conveniente escribir la ecuación (4) en la forma

$$\begin{aligned} \Delta y_i &= (ay_{\bar{x}})_{x,i} - d_i y_i = -\varphi_i, \\ i &= 1, \dots, N-1, \quad y_0 = 0, \quad y_N = 0, \quad (13) \\ a_i &> 0, \quad d_i > 0. \end{aligned}$$

Escribámosla en la forma habitual:

$$\begin{aligned} a_i y_{i-1} - c_i y_i + a_{i+1} y_{i+1} &= -h^2 \varphi_i, \quad y_0 = y_N = 0, \\ c_i &= a_i + a_{i+1} + h^2 d_i, \quad a_i > 0, \quad i = 1, 2, \dots, N-1. \end{aligned}$$

Veamos las fórmulas de factorización

$$\begin{aligned} y_i &= \alpha_{i+1} y_{i+1} + \beta_{i+1}, \quad y_N = 0, \quad i = 1, 2, \dots, N-1, \\ \alpha_{i+1} &= \frac{a_{i+1}}{c_i - a_i \alpha_i}, \quad \alpha_1 = 0, \quad i = 1, 2, \dots, N-1, \\ \beta_{i+1} &= \frac{a_i \beta_i + \varphi_i h^2}{c_i - a_i \alpha_i}, \quad \beta_i = 0 \quad i = 1, 2, \dots, N-1 \end{aligned}$$

Bajo las condiciones (7) tenemos $|\alpha_{i+1}| \leq 1$, y

$$|y_i| \leq |y_{i+1}| + |\beta_{i+1}| \leq |y_N| + \sum_{s=i+1}^N |\beta_s| = \sum_{s=i+1}^N |\beta_s|.$$

Al introducir la función $a_i \beta_i = \eta_i$, obtenemos

$$\begin{aligned} \eta_{i+1} &= (\eta_i + h^2 \varphi_i) \alpha_{i+1}, \\ |\eta_{i+1}| &\leq |\eta_i| + h^2 |\varphi_i| \leq |\eta_i| + \sum_{h=1}^i h^2 |\varphi_h|, \end{aligned}$$

de modo que

$$|\beta_{s+1}| \leq \frac{1}{a_{s+1}} h \sum_{k=1}^s h |\varphi_k|.$$

De resultados se obtiene para la solución del problema una estimación apriorística

$$\|y\|_c \leq \sum_{s=1}^N h \frac{1}{a_s} \sum_{k=1}^s h |\varphi_k| \leq \frac{1}{c_1} \sum_{s=1}^N h \sum_{k=1}^s h |\varphi_k|$$

para $a_s \geq c_1 > 0$.

Esta estimación nos será útil al estudiar la convergencia de los esquemas de diferencias.

Interpolación e integración numérica

§ 1. Interpolación y aproximación de las funciones

1. Planteamiento del problema. Uno de los problemas fundamentales del análisis numérico es la interpolación de las funciones. Se requiere a menudo restablecer la función $f(x)$ para todos los valores de x en el segmento $a \leq x \leq b$, si están conocidos sus valores en cierto número finito de puntos del segmento mencionado. Dichos valores pueden ser determinados como resultado de las mediciones (observaciones) en un experimento natural, o bien como resultado de los cálculos. Además, puede ocurrir que la función $f(x)$ viene definida por cierta fórmula y el cálculo de sus valores, rigiéndose por dicha fórmula, es muy engorroso, razón por la cual resulta deseable tener para la función otra fórmula, más simple, (menos engorrosa para los cálculos) que permitiría hallar valores aproximados de la función en consideración con una exactitud necesaria en cualquier punto del segmento. De resultas, surge el siguiente problema matemático.

Supongamos que en el segmento $a \leq x \leq b$ viene pre-fijada una red $\omega = \{x_0 = a < x_1 < \dots < x_n = b\}$ y en los nodos de la red están definidos los valores de la función $y(x)$ iguales a $y(x_0) = y_0, \dots, y(x_i) = y_i, \dots, y(x_n) = y_n$. Se pide construir una *interpolante*, esto es, una función $f(x)$ que coincida con la función $y(x)$ en los nodos de la red:

$$f(x_i) = y_i, \quad i = 0, 1, \dots, n. \quad (1)$$

El objetivo principal de la interpolación es obtener un algoritmo rápido (económico) para calcular los valores de $f(x)$ en aquellos puntos x que no están contenidos en la tabla de los datos.

La cuestión principal es: cómo elegir la interpolante $f(x)$ y cómo estimar el error $y(x) - f(x)$. Las funciones interpoladoras $f(x)$ se construyen, como regla, en forma de las combinaciones lineales de ciertas funciones elementales:

$$f(x) = \sum_{h=0}^n c_h \Phi_h(x),$$

donde $\{\Phi_h(x)\}$ son funciones linealmente independientes fijas; c_0, c_1, \dots, c_n , unos coeficientes hasta ahora desconocidos.

De las condiciones (1) obtenemos un sistema de $n + 1$ ecuaciones respecto de los coeficientes $\{c_h\}$:

$$\sum_{h=0}^n c_h \Phi_h(x_i) = y_i, \quad i = 0, 1, \dots, n.$$

Supongamos que el sistema de funciones $\Phi_h(x)$ es de tal índole que, cualquiera que sea la elección de los nodos $a = x_0 < x_1 < \dots < x_n = b$, queda distinto de cero el determinante del sistema

$$\Delta(\Phi) = \begin{vmatrix} \Phi_0(x_0) & \Phi_1(x_0) & \dots & \Phi_n(x_0) \\ \Phi_0(x_1) & \Phi_1(x_1) & \dots & \Phi_n(x_1) \\ \dots & \dots & \dots & \dots \\ \Phi_0(x_n) & \Phi_1(x_n) & \dots & \Phi_n(x_n) \end{vmatrix}.$$

En este caso los coeficientes c_k ($k = 0, 1, \dots, n$) se determinan unívocamente según las y_i ($i = 0, 1, \dots, n$) prefijadas.

A título de sistema de las funciones linealmente independientes $\{\Phi_h(x)\}$ se eligen más a menudo: funciones potenciales $\Phi_h(x) = x^h$ (en este caso $f = P_n(x)$ es un polinomio de grado n); funciones trigonométricas $\{\Phi_h(x) = \cos kx, \sin kx\}$ (f es un polinomio trigonométrico). Se emplean también funciones racionales

$$\frac{\alpha_0 + \alpha_1 x + \dots + \alpha_m x^m}{\beta_0 + \beta_1 x + \dots + \beta_p x^p}$$

y otros tipos de funciones interpoladoras. Examinaremos aquí los polinomios de interpolación y la spline-interpolación: un caso de interpolación polinomial a trozos.

Lagrange:

$$l_k(x_i) = \begin{cases} 1, & \text{si } i = k, \\ 0, & \text{si } i \neq k, \end{cases} \quad i, k = 0, 1, \dots, n.$$

No es difícil ver que el polinomio de grado n

$$l_k(x) = l_k^{(n)}(x) = \frac{(x-x_0)(x-x_1)\dots(x-x_{k-1})(x-x_{k+1})\dots(x-x_n)}{(x_k-x_0)(x_k-x_1)\dots(x_k-x_{k-1})(x_k-x_{k+1})\dots(x_k-x_n)}$$

satisface estas condiciones. El polinomio $l_k(x)$ se define, evidentemente, del modo unívoco. Efectivamente, supongamos que existe un polinomio más $\bar{l}_k(x)$; entonces la diferencia entre ellos $\bar{l}_k(x) - l_k(x) = q_n(x)$ es un polinomio de grado n que se reduce a cero en $n + 1$ puntos x_i ($i = 0, 1, \dots, n$). Esto será posible sólo cuando $\bar{l}_k(x) - l_k(x) \equiv 0$.

El polinomio $l_k(x) y_k$ toma el valor y_k en el punto x_k y es nulo en todos los demás nodos x_j para $j \neq k$. De aquí se desprende que el polinomio de interpolación

$$P_n(x) = \sum_{k=0}^n l_k(x) y_k = \sum_{k=0}^n y_k \prod_{i \neq k} \frac{x-x_i}{x_k-x_i} \quad (3)$$

tiene el grado no superior a n y $P_n(x_i) = y_i$. La fórmula (3) lleva el nombre de *Lagrange*. El número de operaciones aritméticas para el cálculo según (3) es proporcional a n^2 . Para estimar la proximidad del polinomio $P_n(x)$ a la función $f(x)$ se supone que existe la $n + 1$ -ésima derivada continua $f^{(n+1)}(x)$. En este caso resulta verídica la fórmula siguiente para el error:

$$f(x) - P_n(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \prod_{j=1}^{n+1} (x-x_j), \quad \xi \in [a, b].$$

3. Fórmula de interpolación de Newton. Al emplear en los cálculos los ordenadores, es cómoda la fórmula de interpolación de Newton. Con el fin de escribirla se debe introducir las así llamadas diferencias divididas:

— la diferencia dividida de primer orden: $y(x_i, x_j) = [y(x_i) - y(x_j)] / (x_i - x_j)$;

— la diferencia dividida de segundo orden: $y(x_i, x_j, x_k) = [y(x_i, x_j) - y(x_j, x_k)] / (x_i - x_k)$, etc. Si $y(x) = P_n(x)$ es un polinomio de grado n , entonces para él la primera diferencia dividida $P(x, x_0) = [P(x) - P(x_0)] / (x - x_0)$ será un polinomio de grado $n - 1$; la segunda diferencia $P(x, x_0, x_1)$, polinomio de grado $n - 2$, etc., de suerte que la $(n + 1)$ -ésima diferencia dividida es igual a cero.

De la definición de diferencias divididas se deduce:

$$P(x) = P(x_0) + (x - x_0) P(x, x_0),$$

$$P(x, x_0) = P(x_0, x_1) + (x - x_1) P(x, x_0, x_1),$$

$$P(x, x_0, x_1) = P(x_0, x_1, x_2) + (x - x_2) P(x, x_0, x_1, x_2),$$

etc. De aquí obtenemos para $P(x)$ la fórmula

$$P(x) = P(x_0) + (x - x_0) P(x_0, x_1) + (x - x_0)(x - x_1) \times \\ \times P(x_0, x_1, x_2) + \dots + (x - x_0)(x - x_1) \dots \\ \dots (x - x_n) P(x_0, x_1, \dots, x_n). \quad (4)$$

Si $P(x)$ es el polinomio de interpolación para la función $y(x)$, sus valores en los nodos de las redes coincidirán con los valores de la función $y(x)$, y, por consiguiente, coincidirán también las diferencias divididas. Por eso, en lugar de (4) podemos escribir

$$f(x) = y_0 + \sum_{k=1}^n (x - x_0)(x - x_1) \dots \\ \dots (x - x_{k-1}) y(x_0, x_1, \dots, x_k)$$

(polinomio de Newton). Calculadas las diferencias divididas, el polinomio de Newton se calculará con toda la comodidad según el esquema de Horner

$$j(x) = y(x_0) + (x - x_0) [y(x_0, x_1) + (x - x_1) \times \\ \times [y(x_0, x_1, x_2) + \dots]].$$

El cálculo de $f(x)$ para cada x requiere n multiplicaciones y $2n$ operaciones de sumación y sustracción.

Existen también otras fórmulas de interpolación. Entre ellas resulta más aplicable la *interpolación hermitiana*. Aquí el problema se plantea del modo siguiente. Están prefijados n nodos $\{x_i\}$, n valores de la función $\{y_i\}$ y n valores de la

derivada $\{y'_i\}$; se pide hallar tal polinomio de grado máximo $2n - 1$ que se verifique

$$P(x_i) = y_i, \quad P'(x_i) = y'_i, \quad i = 1, 2, \dots, n.$$

Si todos los x_i son distintos, existe la única solución que se halla por un método análogo al de Lagrange.

Se debe tener en cuenta que la aplicación de un polinomio de alto grado puede conducir a los problemas difíciles relacionados con los errores de redondeo.

4. Spline-interpolación. Estudiemos un caso especial de la interpolación polinomial a trozos cuando entre cualesquiera nodos vecinos de la red la función viene interpolada por un polinomio cúbico (*spline-interpolación cúbica*). Los coeficientes de dicho polinomio se determinan en cada intervalo partiendo de las condiciones de conjugación en los nodos:

$$\begin{aligned} f_i &= y_i \\ f'(x_i - 0) &= f'(x_i + 0), \\ f''(x_i - 0) &= f''(x_i + 0), \quad i = 1, 2, \dots, n - 1. \end{aligned}$$

Además, en la frontera se ponen las condiciones para $x = x_0$ y $x = x_n$

$$f''(x_0) = 0, \quad f''(x_n) = 0. \quad (5)$$

Buscaremos el polinomio cúbico en la forma

$$f(x) = a_i + b_i(x - x_{i-1}) + c_i(x - x_{i-1})^2 + d_i(x - x_{i-1})^3, \quad x_{i-1} \leq x \leq x_i. \quad (6)$$

De la condición $f_i = y_i$ tenemos

$$\begin{aligned} f(x_{i-1}) &= a_i = y_{i-1}, \\ f(x_i) &= a_i + b_i h_i + c_i h_i^2 + d_i h_i^3 = y_i, \\ h_i &= x_i - x_{i-1}, \quad i = 1, 2, \dots, n - 1. \end{aligned} \quad (7)$$

Calculemos las derivadas:

$$\begin{aligned} f'(x) &= b_i + 2c_i(x - x_{i-1}) + 3d_i(x - x_{i-1})^2, \\ f''(x) &= 2c_i + 6d_i(x - x_{i-1}), \quad x_{i-1} \leq x \leq x_i \end{aligned}$$

y exijamos su continuidad para $x = x_i$:

$$\begin{aligned} b_{i+1} &= b_i + 2c_i h_i + 3d_i h_i^2, \\ c_{i+1} &= c_i + 3d_i h_i, \quad i = 1, 2, \dots, n-1. \end{aligned} \quad (8)$$

El número total de los coeficientes desconocidos es igual, evidentemente a $4n$, el número de ecuaciones (7) y (8) equivale a $4n - 2$. Dos ecuaciones que faltan las obtenemos de las condiciones (5) para $x = x_0$ y $x = x_n$:

$$c_1 = 0, \quad c_n + 3d_n h_n = 0.$$

Expresando a base de (8) $d_i = (c_{i+1} - c_i)/3h_i$, sustituyendo esta expresión en (7) y excluyendo $a_i = y_{i-1}$, obtendremos

$$\begin{aligned} b_i &= [(y_i - y_{i-1})/h_i] - \frac{1}{3} h_i (c_{i+1} + 2c_i), \quad i = 1, 2, \dots, n-1, \\ b_n &= [(y_n - y_{n-1})/h_n] - \frac{2}{3} h_n c_n. \end{aligned}$$

Ahora, al sustituir las expresiones para b_i , b_{i+1} y d_i en la primera fórmula de (8), obtenemos, después de algunas transformaciones no complejas, una ecuación en diferencias de segundo orden

$$h_i c_i + 2(h_i + h_{i+1}) c_{i+1} + h_{i+1} c_{i+2} = 3 \left(\frac{y_{i+1} - y_i}{h_{i+1}} - \frac{y_i - y_{i-1}}{h_i} \right), \quad i = 1, 2, \dots, n-1, \quad (9)$$

con las condiciones de contorno

$$c_1 = 0, \quad c_{n+1} = 0, \quad (10)$$

y esta ecuación se usa para determinar c_i . La condición $c_{n+1} = 0$ es equivalente a la condición $c_n + 3d_n h_n = 0$ y a la ecuación $c_{i+1} = c_i + d_i h_i$. La ecuación en diferencias (9) con las condiciones (10) se resuelve por el método de factorización.

Se puede introducir la noción de *spline de orden m* como función que es un polinomio de grado m en cada uno de los segmentos de la red y que en todos los nodos interiores de la red satisface las condiciones de continuidad de la función y de las derivadas hasta el orden $m - 1$ inclusive. Habitualmente se usan para la interpolación los casos de $m = 3$

(spline cúbico analizado más arriba) y de $m = 1$ (*spline lineal*, correspondiente a la aproximación de la gráfica de la función $y(x)$ por una quebrada que pasa a través de los puntos (x_i, y_i)).

5. Aplicación de la interpolación. La interpolación se aplica en varios problemas relacionados con los cálculos. Indiquemos aquí algunos de estos problemas.

La elaboración de un experimento físico consistente en la construcción de las fórmulas aproximadas para las magnitudes características según datos tabulares obtenidos en los experimentos.

La construcción de las fórmulas aproximadas a base de los datos de un experimento de cálculo. En este caso surgen problemas no típicos de interpolación, ya que, corrientemente, se escriben fórmulas cuya estructura sea cuanto más simple.

La subtabulación, o sea, el espesamiento de las tablas se usa en aquellos casos cuando el cálculo inmediato de las funciones resulta difícil, o cuando se tienen pocos datos experimentales. A la máquina electrónica se introduce una tabla pequeña, mientras que los valores de la función indispensables en los cálculos se hallan, cuando sea necesario, según la fórmula de interpolación.

La interpolación se aplica también en el problema de *interpolación inversa*: está dada la tabla $y_i = y(x_i)$; se pide hallar x_i como función de y_i . A título de ejemplo de interpolación inversa puede servir el problema de búsqueda de las raíces de una ecuación.

Las fórmulas de interpolación se emplean también al calcular integrales y al escribir aproximaciones de diferencias para las ecuaciones diferenciales a base de las identidades integrales. El aparato matemático de cualquier ordenador contiene programas estándar de interpolación.

6. Aproximación media cuadrática. Hasta ahora hemos analizado la construcción de los polinomios de interpolación $y(x)$ que coinciden con los valores de la función de partida $f(x)$ en cierto conjunto de nodos sobre la red ω :

$$y(x_i) = f(x_i), \quad x_i \in \omega.$$

La función $y(x)$ aproxima la función $f(x)$ en el intervalo de la red.

Sea $L_2 [a, b]$ un espacio de funciones reales con producto escalar

$$(f, \varphi) = \int_a^b f(x) \varphi(x) dx$$

y norma

$$\|f\|_{L_2} = \sqrt{(f, f)}.$$

Examinemos el problema general sobre la aproximación de las funciones $f(x)$ mediante las funciones pertenecientes a L_2 , sustituyendo la exigencia $y_1 = f_1$ por la condición del mínimo de la norma: $\|f - y\|_{L_2}$ o de la pequeñez de la misma: $\|f - y\|_{L_2} < \varepsilon$, donde $\varepsilon > 0$ es la exactitud prefijada.

La búsqueda de $\inf \|f - y\|_{L_2}$ es el problema de encontrar la *mejor aproximación media cuadrática*. A título de $y(x)$ tomemos el *polinomio generalizado*

$$y(x) = \sum_{k=0}^n c_k \varphi_k(x),$$

donde $\{\varphi_k(x)\}$ es una familia de funciones ortonormalizadas en $[a, b]$

$$(\varphi_k, \varphi_m) = \delta_{km}, \quad \delta_{km} = \begin{cases} 1, & k = m, \\ 0, & k \neq m, \end{cases}$$

y c_k son unos coeficientes arbitrarios. Entonces, el problema de encontrar la mejor aproximación media cuadrática se reduce a la búsqueda del mínimo de la función de $n+1$ variables c_0, c_1, \dots, c_n :

$$\min_{(c_k)} \left\| f(x) - \sum_{k=0}^n c_k \varphi_k(x) \right\| = F(c_0, c_1, \dots, c_n).$$

Calculemos la *desviación media cuadrática*

$$\|f - y\|^2 = \|f\|^2 - 2(f, y) + \|y\|^2.$$

Sustituyendo aquí las expresiones

$$(f, y) = \sum_{k=0}^n c_k (f, \varphi_k) = \sum_{k=0}^n c_k f_k, \quad f_k = (f, \varphi_k),$$

$$\|y\|^2 = \sum_{k=0}^n c_k^2,$$

obtendremos

$$\|f - y\|^2 = \|f\|^2 + \sum_{k=0}^n (c_k - f_k)^2 - \sum_{k=0}^n f_k^2.$$

De aquí se ve que el mínimo del error se consigue para $c_k = f_k$, es decir, en la función

$$\bar{y}(x) = \bar{y}_n(x) = \sum_{k=0}^n f_k \varphi_k(x).$$

En este caso

$$\|f - \bar{y}_n(x)\|^2 = \|f\|^2 - \sum_{k=0}^n f_k^2. \quad (11)$$

De este modo, la mejor aproximación media cuadrática existe y es única. Ella lleva al problema sobre el cálculo de las integrales para determinar $c_k = f_k = (f, \varphi_k)$.

Si las funciones $\{\varphi_k\}$ forman un sistema ortonormalizado completo, se verifica la igualdad de Parseval—Steklov

$$\sum_{k=0}^{\infty} f_k^2 = \int_a^b f^2(x) dx = \|f\|^2. \quad (12)$$

Al comparar (11) con (12), encontramos

$$\|f - \bar{y}_n\|^2 = \sum_{k=n+1}^{\infty} f_k^2,$$

es decir, $\|f - \bar{y}_n\| \rightarrow 0$ cuando $n \rightarrow \infty$; la mejor aproximación media cuadrática converge hacia $f(x)$ y queda posible la aproximación con cualquier exactitud: $\|f - \bar{y}_n\| \leq \varepsilon$, siempre que $n \geq N(\varepsilon)$ (n es bastante grande).

OBSERVACIÓN. Todos los razonamientos están en vigor, si el producto escalar se toma con el peso $\rho(x) > 0$:

$$(f, \varphi) = \int_a^b f(x) \varphi(x) \rho(x) dx.$$

Son posibles también otros criterios de la aproximación, cuando la desviación $f - y$ se minimiza en otra norma, por ejemplo, en la norma del espacio C (aproximación uniforme).

Realizándose la mejor aproximación uniforme, nosotros buscamos la función $y(x)$ en la cual se consigue

$$\min_{(y)} \max_{a \leq x \leq b} |f(x) - y(x)|.$$

Sin embargo, hasta ahora no se ha encontrado un método que permita encontrar los coeficientes de la mejor aproximación uniforme (por el número finito de operaciones) para una función, definida en el segmento $[a, b]$. Se pueden indicar, además, otros planteamientos de los problemas de aproximación: en un conjunto discreto, en una totalidad de segmentos y otros. Se estudian también los métodos de aproximación no lineal, por ejemplo, con ayuda de las funciones racionales. Lo último resulta efectivo al elaborar los resultados de los experimentos.

§ 2. Integración numérica

1. Planteamiento del problema. Fórmulas de integración numérica (de cuadratura) más simples. El objetivo de la *integración numérica* es hallar el valor aproximado de la integral

$$J[f] = \int_a^b f(x) dx, \quad (1)$$

donde $f(x)$ es una función prefijada. En el segmento $[a, b]$ se introduce una red $\bar{\omega} = \{x_i: x_0 = a < x_1 < \dots < x_i < x_{i+1} < \dots < x_N = b\}$ y como el valor aproximado de la integral se considera el número

$$J_N[f] = \sum_{i=0}^N c_i f(x_i), \quad (2)$$

donde $f(x_i)$ son los valores de la función $f(x)$ en los *nodos* $x = x_i$ y c_i , *factores ponderales (de peso)* que sólo dependen de los nodos, pero no son dependientes de la elección de $f(x)$. La fórmula (2) se denomina de *cuadratura*, o de *integración numérica*.

El objeto de la integración numérica con ayuda de cuadraturas consiste en búsqueda de tales nodos $\{x_i\}$ y pesos $\{c_i\}$

que el error de la fórmula de cuadratura

$$D[f] = \sum_{i=0}^N c_i f(x_i) - \int_a^b f(x) dx = J_N[f] - J[f]$$

sea mínimo para las funciones pertenecientes a la clase dada (la magnitud de $D[f]$ depende del grado de suavidad de $f(x)$). Al construir la fórmula de cuadratura, la integral (1) se representa, corrientemente, en forma de una suma de las integrales del tipo

$$\int_{\alpha}^{\beta} f(x) dx,$$

cada una de las cuales se reduce a la integral estándar por el segmento de longitud unidad:

$$L[f] = \int_0^1 f(s) ds \quad (3)$$

mediante una sustitución

$$x = \alpha + (\beta - \alpha)s, \quad (4)$$

$$f(x) = f(\alpha + (\beta - \alpha)s) = \bar{f}(s), \quad (5)$$

de modo que

$$\int_{\alpha}^{\beta} f(x) dx = \kappa \int_0^1 \bar{f}(s) ds = \kappa L[\bar{f}], \quad \kappa = \beta - \alpha$$

(la raya por arriba de $f(s)$ se omitirá). Convengamos en considerar que $\bar{\omega}$ es una red uniforme. En este caso podemos escribir

$$J[f] = \sum_{i=1}^N J_i,$$

$$J_i = \int_{x_{i-1}}^{x_i} f(x) dx = h \int_0^1 f(x_{i-1} + hs) ds.$$

Si $N = 2i_0$ es un número par, tenemos

$$J[f] = \sum_{i=1}^{i_0} J_{2i-1},$$

$$J_{2i-1} = \int_{x_{2i-2}}^{x_{2i}} f(x) dx = 2h \int_0^1 f(x_{2i-2} + 2hs) ds,$$

etc.

Así pues, el problema se reduce a la construcción de la fórmula de cuadratura para la integral (3) que se calcula por un segmento unidad. Escojamos en el segmento $0 \leq s \leq 1$ los nodos $0 \leq s_0 < s_1 < \dots < s_m \leq 1$ (molde de la fórmula de cuadratura) y a la integral (3) le pondremos en correspondencia la fórmula

$$\Lambda(f) = \sum_{h=0}^m p_h f(s_h). \quad (6)$$

Veamos las fórmulas de cuadratura más simples:

— *fórmula del rectángulo* (el molde contiene un nodo):

$$m=0, \quad p_0=1, \quad s_0=\frac{1}{2}, \quad \Lambda(f) = f\left(\frac{1}{2}\right);$$

— *fórmula del trapecio* (dos nodos):

$$m=1, \quad p_0=\frac{1}{2}, \quad p_1=\frac{1}{2}; \quad s_0=0, \quad s_1=1$$

$$\Lambda(f) = \frac{1}{2} (f(0) + f(1));$$

— *fórmula de Simpson* (tres nodos):

$$m=2, \quad p_0=p_2=\frac{1}{6}, \quad p_1=\frac{4}{6}, \quad s_0=0, \quad s_1=\frac{1}{2}, \quad s_2=1,$$

$$\Lambda(f) = \frac{1}{6} \left(f(0) + 4f\left(\frac{1}{2}\right) + f(1) \right)$$

y otras. En la práctica se emplean, como regla, las fórmulas con un número pequeño de nodos del molde.

Escribamos ahora las fórmulas correspondientes para la integral (1) en una red uniforme $\{x_i = ih\}$ de paso h . Teniendo presente la sustitución (4) y (5), obtendremos

— fórmula del rectángulo:

$$J_N[f] = \sum_{i=0}^{N-1} hf(x_{i+1/2}), \quad x_{i+1/2} = x_i + \frac{1}{2}h; \quad (7)$$

— fórmula del trapecio:

$$J_N[f] = \sum_{i=0}^N c_i f(x_i) h, \quad c_0 = c_N = \frac{1}{2}, \quad c_i = 1, \\ i = 1, 2, \dots, N-1; \quad (8)$$

— fórmula de Simpson:

$$J_N[f] = \sum_{i=0}^N c_i f(x_i) \bar{h} = \frac{h}{3} (f_0 + 4f_1 + 2f_2 + 4f_3 + \dots \\ \dots + 2f_{N-2} + 4f_{N-1} + f_N) \text{ para } N = 2i_0. \quad (9)$$

2. Construcción de las fórmulas de cuadratura. En virtud de lo dicho más arriba, exponer el problema será suficiente para la integral tipo (3), a la que se le pone en correspondencia la fórmula de cuadratura

$$\int_0^1 f(s) ds \approx \sum_{h=0}^m p_h f(s_h). \quad (10)$$

En el caso general los nodos y los pesos son desconocidos y han de ser determinados.

Examinemos al principio un caso en que los nodos están prefijados y se requiere hallar los pesos de la fórmula de cuadratura $\{p_h\}$. Hagamos uso del requisito: la fórmula (10) debe ser exacta para cualquier polinomio $P_r(s)$ de grado $r \leq m$:

$$\Lambda [P_r] = L [P_r], \quad r \leq m. \quad (11)$$

Para que un polinomio de grado r satisfaga (11), basta por exigir que la fórmula de cuadratura sea exacta para cualquier monomio s^σ de grado σ ($\sigma = 0, 1, \dots, r$). Teniendo presente que $L[s^\sigma] = 1/(\sigma + 1)$, obtenemos de (11) $m + 1$

se ve que la fórmula (10) es exacta para un polinomio de grado m , si los factores ponderales p_k se determinan según la fórmula

$$p_k = \int_0^1 l_k^{(m)}(s) ds. \quad (12)$$

Las fórmulas de este tipo se llaman *fórmulas de cuadratura de Cotes*.

Aduzcamos como ejemplos de las fórmulas de cuadratura dos fórmulas más:

en el molde tetrapuntual, $s_k = k/3$ ($k = 0, 1, 2, 3$), $m = 3$:

$$\Lambda(f) = \frac{1}{8} \left(f(0) + 3f\left(\frac{1}{3}\right) + 3f\left(\frac{2}{3}\right) + f(1) \right),$$

$$p_0 = p_3 = \frac{1}{8}, \quad p_1 = p_2 = \frac{3}{8},$$

en el molde pentapuntual, $s_k = k/4$ ($k = 0, 1, 2, 3, 4$), $m = 4$:

$$\Lambda(f) = \frac{1}{90} \left(7f(0) + 32f\left(\frac{1}{4}\right) + 12f\left(\frac{1}{2}\right) + \right.$$

$$\left. + 32f\left(\frac{3}{4}\right) + 7f(1) \right),$$

$$p_0 = p_4 = \frac{7}{90}, \quad p_1 = p_3 = \frac{32}{90}, \quad p_2 = \frac{12}{90}.$$

Los moldes de las cinco fórmulas de cuadratura aducidas más arriba constan de los nodos simétricos con relación al centro $s = 1/2$ del segmento $0 \leq s \leq 1$.

3. Fórmula de Taylor con término residual en la forma integral. Al investigar el error de la fórmula de cuadratura nos hará falta la fórmula de Taylor con término residual en la forma integral:

$$f(s) = f(0) + sf'(0) + \frac{s^2}{2} f''(0) + \dots + \frac{s^n}{n!} f^{(n)}(0) + R_{n+1}(s), \quad (13)$$

$$R_{n+1}(s) = \int_0^s \frac{(s-t)^n}{n!} f^{(n+1)}(t) dt,$$

Dicha fórmula puede ser demostrada por inducción respecto de n . Para $n = 0$ es justa:

$$f(s) = f(0) + R_1(s), \quad R_1(s) = \int_0^s f'(t) dt.$$

Admitamos que es justa para n . Integrando por partes, obtenemos una correlación

$$\begin{aligned} \int_0^s \frac{(s-t)^n}{n!} f^{(n+1)}(t) dt &= \\ &= -\frac{(s-t)^{n+1}}{(n+1)!} f^{(n+1)}(t) \Big|_0^s + \int_0^s \frac{(s-t)^{n+1}}{(n+1)!} f^{(n+2)}(t) dt = \\ &= \frac{s^{n+1}}{(n+1)!} f^{(n+1)}(0) + \int_0^s \frac{(s-t)^{n+1}}{(n+1)!} f^{(n+2)}(t) dt, \quad (14) \end{aligned}$$

la cual demuestra la fórmula (13) precisamente para $n + 1$. Introduciendo la función

$$K_n(\xi) = \begin{cases} \xi^n/n! & \text{para } \xi \geq 0, \\ 0 & \text{para } \xi < 0, \end{cases} \quad (15)$$

escribamos la fórmula para el término residual R_{n+1} en la forma

$$R_{n+1}(s) = \int_0^1 K_n(s-t) f^{(n+1)}(t) dt. \quad (16)$$

4. Fórmula para el error de la fórmula de cuadratura. Pasemos a la deducción de una fórmula para el error de la fórmula de cuadratura

$$\Delta(f) = \Lambda[f] - L[f] \quad (17)$$

en la clase $C^{(n+1)}$ de funciones que tienen la $(n+1)$ -ésima derivada continua en el segmento $0 \leq s \leq 1$: $f(s) \in C^{(n+1)}[0, 1]$. En este caso sirve la fórmula (13) o bien

$$f(s) = P_n(s) + R_{n+1}(s), \quad P_n(s) = \sum_{\sigma=0}^n \frac{s^\sigma}{\sigma!} f^{(\sigma)}(0). \quad (18)$$

De lo expuesto anteriormente (véase el p. 2) está claro que para un polinomio $P_n(s)$ de grado n la fórmula (10) es exacta en dos casos: para $n \leq m + 1 = n_0$, si m es par y la fórmula es simétrica; para $n \leq m = n_0$ en todos los demás casos. Supondremos por ahora que

$$\Lambda [P_n] = L [P_n], \quad \text{es decir, } n \leq n_0. \quad (19)$$

Volvamos ahora a la diferencia $\Delta(f)$ y sustituyamos $f = P_n + R_{n+1}$ en (17). Tomando en consideración (16) y (19), obtendremos

$$\begin{aligned} \Delta(f) &= \Lambda [f] - L [f] = \\ &= (\Lambda [P_n] - L [P_n]) + (\Lambda [R_{n+1}] - L [R_{n+1}]) = \\ &= \Lambda [R_{n+1}] - L [R_{n+1}] = \sum_{k=0}^m p_k \int_0^1 K_n(s_k - t) f^{(n+1)}(t) dt - \\ &\quad - \int_0^1 \int_0^1 K_n(s - t) f^{(n+1)}(t) dt ds = \\ &= \int_0^1 \left[\sum_{k=0}^m p_k K_n(s_k - t) - \int_0^1 K_n(s - t) ds \right] f^{(n+1)}(t) dt. \end{aligned}$$

Haciendo uso de la expresión (15) para $K_n(s - t)$, hallamos

$$\int_0^1 K_n(s - t) ds = \int_t^1 \frac{(s-t)^n}{n!} ds = \frac{(1-t)^{n+1}}{(n+1)!}.$$

De resultas, la fórmula para el error toma la forma

$$\Delta(f) = \int_0^1 F_{n+1}(t) f^{(n+1)}(t) dt, \quad (20)$$

donde

$$F_{n+1}(t) = \sum_{k=0}^m p_k K_n(s_k - t) - \frac{(1-t)^{n+1}}{(n+1)!}. \quad (21)$$

De aquí se desprende la estimación para el error

$$|\Delta(f)| \leq M_{n+1} c_{n+1} \quad (22)$$

para $|f^{(n+1)}(t)| \leq M_{n+1}$, donde $M_{n+1} > 0$ es una constante, y para

$$c_{n+1} = \int_0^1 |F_{n+1}(t)| dt.$$

Si $F_{n+1}(t)$ no cambia de signo en el segmento $0 \leq s \leq 1$, entonces, en virtud del teorema del valor medio, tenemos

$$\Delta(f) = f^{(n+1)}(\xi) \int_0^1 F_{n+1}(t) dt, \quad \xi \in [0, 1].$$

5. Estimación del error de las fórmulas concretas. Nuestro objetivo es obtener la estimación del error $\Delta(\bar{f}) = = \Lambda(\bar{f}) - L(\bar{f})$ de la fórmula de cuadratura para la integral estándar (3). Al pasar a las fórmulas para las integrales (1) y (3), se debe tener en cuenta que

$$\frac{d^\sigma \bar{f}(s)}{ds^\sigma} = \kappa^\sigma \frac{d^\sigma f(x)}{dx^\sigma},$$

$\bar{f}(s) = f(x)$, $x = \alpha + (\beta - \alpha)s$, $dx = \kappa ds$, $\kappa = \beta - \alpha$.
Por eso, para el error

$$d[f] = \sum_{k=0}^m \kappa p_k f(x_k) - \int_{\alpha}^{\beta} f(x) dx = \kappa \Delta(\bar{f})$$

es justa, en virtud de (22), la fórmula

$$|d[f]| \leq c_{n+1} \kappa^{n+2} \max_{x \in [\alpha, \beta]} |f^{(n+1)}(x)|,$$

$$c_{n+1} = \int_0^1 |F_{n+1}(t)| dt.$$

Para el cálculo del error $J_N[f] - J[f]$ es necesario, evidentemente, sumar sobre la red los errores $|D[f]|$.

Veamos las fórmulas de cuadratura más simples.

1) FÓRMULA DEL RECTÁNGULO: $m = 0$, $p_0 = 1$, $s_0 = 1/2$,
 $\Lambda(\bar{f}) = \bar{f}(1/2)$. Debido a la fórmula (20) tenemos

$$\Delta_1(\bar{f}) = \int_0^1 F_2(t) \bar{f}''(t) dt, \quad F_2(t) = K_1 \left(\frac{1}{2} - t \right) - \frac{(1-t)^2}{2},$$

es decir, $F_2(t) = -(1-t)^2/2 < 0$ para $t > 1/2$, $F_2(t) = (1/2 - t) - (1-t)^2/2 = -t^2/2 < 0$ para $t < 1/2$, es decir, $F(t) < 0$ es una función de signo constante y

$$\Delta_1(\bar{f}) = \bar{f}''(\eta) \int_0^1 F_2(t) dt = -\frac{\bar{f}''(\eta)}{24}, \quad \eta \in (0, 1).$$

De aquí se infiere que

$$d_i[f] = hf(x_{i-1/2}) - \int_{x_{i-1}}^{x_i} f(x) dx = -\frac{h^3}{24} f''(\xi_i),$$

$$\xi_i \in [x_{i-1}, x_i]. \quad (23)$$

Sumando según $i = 1, 2, \dots, N$, y teniendo presente que la media aritmética es igual a

$$\sum_{i=1}^N hf''(\xi_i) = \frac{b-a}{N} \sum_{i=1}^N f''(\xi_i) = f''(\xi^*) (b-a), \quad \xi^* \in [a, b],$$

obtenemos para el error la fórmula del rectángulo:

$$D_N(f) = -\frac{h^3}{24} f''(\xi^*) (b-a).$$

Si $f(x)$ tiene derivadas continuas por lo menos de cuarto orden, $f(x) \in C^{(n)}$ ($n \geq 4$), podemos anotar el desarrollo asintótico para el error:

$$D_N(f) = \alpha_2 h^2 + \alpha_4 h^4, \quad (24)$$

donde

$$\alpha_2 = -\frac{1}{24} \int_a^b f''(x) dx = -\frac{1}{24} [f'(b) - f'(a)].$$

En efecto, al sustituir en (20) la expresión

$$f''(t) = \bar{f}''\left(\frac{1}{2}\right) + \left(t - \frac{1}{2}\right) \bar{f}'''\left(\frac{1}{2}\right) +$$

$$+ \frac{1}{2} \left(t - \frac{1}{2}\right)^2 \bar{f}^{(IV)}(\eta), \quad \eta \in (0, 1)$$

hallemos, después de ciertos cálculos no complejos,

$$\Delta_1(\bar{f}) = -\frac{1}{24} \bar{f}'' \left(\frac{1}{2} \right) + \frac{1}{960} \bar{f}^{IV}(\eta), \quad \eta \in (0, 1).$$

De aquí se deduce que

$$D_N(f) = -\frac{h^2}{24} \sum_{i=1}^N h f_{i-1/2}'' + \frac{h^4}{960} \sum_{i=1}^N h f^{IV}(\xi_i).$$

Al tomar en consideración que, en virtud de (23),

$$\sum_{i=1}^N h f_{i-1/2}'' = \int_a^b f''(x) dx - \frac{h^2}{24} f^{IV}(\xi^*) \cdot (b-a), \quad \xi^* \in [a, b],$$

obtenemos el desarrollo (24).

De (24) se ve que la fórmula del rectángulo tiene el cuarto grado de precisión: $D_N(f) = O(h^4)$, si la función $f(x)$ satisface la condición $f'(a) = f'(b)$. Si se conocen $f'(a)$ y $f'(b)$, podemos poner $f(x) = \varphi(x) + \alpha x + \beta x^2$, donde $\varphi(x)$ satisface la condición $\varphi'(a) = \varphi'(b)$, siempre que α y β se elijan del modo siguiente

$$\alpha = \frac{b f'(a) - a f'(b)}{b-a}, \quad \beta = \frac{f'(b) - f'(a)}{2(b-a)}.$$

Entonces

$$\begin{aligned} \int_a^b f(x) dx &= \int_a^b \varphi(x) dx + c, \\ c &= \frac{1}{2} \alpha (b^2 - a^2) + \frac{1}{6} \beta (b^3 - a^3). \end{aligned}$$

La integral de $\varphi(x)$ se calcula según la fórmula del rectángulo con la exactitud de $O(h^4)$.

2) FÓRMULA DEL TRAPEZIO: $m = 1$, $p_0 = p_1 = 1/2$, $s_0 = 0$, $s_1 = 1$,

$$\Lambda(\bar{f}) = \frac{1}{2} (\bar{f}(0) + \bar{f}(1)).$$

La función $F_2(t) = \frac{1}{2} t(1-t) > 0$ es de signo constante, por lo cual queda válida la estimación

$$D_N(f) = \frac{h^2}{12} f''(\xi^*) \cdot (b-a), \quad \xi^* \in [a, b],$$

es decir, el coeficiente de h^2 en la expresión para el error de la fórmula del trapecio es dos veces mayor que para la fórmula del rectángulo. Reiterando los razonamientos, análogos a los citados más arriba, nos convencemos de que es justa la fórmula

$$D_N(f) = -2\alpha_2 h^2 + \alpha_4 h^4 \quad \text{para } f \in C^{(n)}, \quad n \geq 4,$$

donde α_2 se determina de acuerdo con (24), $\alpha_4 = O(1)$.

3. FÓRMULA DE SIMPSON: $m = 2$, $s_0 = 0$, $s_1 = 1/2$, $s_2 = 1$, $p_0 = p_2 = 1/2$, $p_1 = 4/6$.

$$\Lambda(\bar{f}) = \frac{1}{6} \left(\bar{f}(0) + 4\bar{f}\left(\frac{1}{2}\right) + \bar{f}(1) \right).$$

Por cuanto la fórmula de Simpson es exacta para un polinomio de tercer grado, entonces $n = 3$ y calculamos:

$$\Delta_3(\bar{f}) = \int_0^1 F_4(t) \bar{f}^{IV}(t) dt,$$

$$F_4(t) = \frac{1}{6} (K_3(0-t) + K_3(1-t)) + \\ + \frac{4}{6} K_3\left(\frac{1}{2}-t\right) - \frac{(1-t)^4}{24}.$$

De aquí encontramos

$$F_4(t) = \frac{1}{72} (2t^3 - 3t^4), \quad t < \frac{1}{2};$$

$$F_4(t) = \frac{1}{72} (2(1-t)^3 - 3(1-t)^4), \quad t > \frac{1}{2},$$

$$F_4(t) > 0 \quad \text{para todos los } t \in (0, 1),$$

y, por consiguiente,

$$\int_0^1 F_4(t) dt = \frac{1}{2880}$$

de suerte que es exacta la fórmula

$$\Delta_3(\bar{f}) = \frac{1}{2880} \bar{f}^{IV}(\eta), \quad \eta \in (0, 1).$$

Pasando a las integrales respecto de x y teniendo presente que $\kappa = 2h$, $\bar{f}^{IV}(\eta) = (2h)^4 f^{IV}(\xi_1)$, obtendremos

$$D_N(f) = \sum_{i=0}^{i_0-1} 2h \left\{ \frac{f_{i-1} + 4f_i + f_{i+1}}{6} - \frac{1}{2h} \int_{x_{i-1}}^{x_{i+1}} f(x) dx \right\} = \\ = \frac{b-a}{180} h^6 f^{IV}(\xi^*), \quad \xi^* \in [a, b],$$

donde $N = 2i_0$, $h = 1/N$.

Si es que $f(x) \in C^{(n)}$ ($n \geq 6$), entonces podemos obtener un desarrollo de la forma

$$D_N(f) = \alpha_4 h^4 + \alpha_6 h^6, \quad \alpha_6 = O(1), \\ \alpha_4 = \frac{1}{180} \int_0^1 f^{IV}(x) dx = \frac{1}{180} (f''(1) - f''(0)).$$

6. Aumento del orden de exactitud. Método de Runge.

Para las fórmulas de cuadratura (por la analogía con lo anterior) se puede obtener un desarrollo asintótico de la forma

$$D_N(f) = J_N(f) - J(f) = \alpha_2 h^2 + \alpha_4 h^4 + \alpha_6 h^6 + \dots,$$

si $f(x)$ es una función suficientemente suave. En este caso $|\alpha_{k+2}|$ es considerablemente inferior a $|\alpha_k|$ ($k = 2, 4$), razón por la cual el aumento del orden de exactitud de la fórmula de cuadratura resulta muy importante.

Realicemos los cálculos sobre dos redes uniformes con los pasos h_1 y h_2 , respectivamente, y hallemos las expresiones

$$J^{h_1}[f] = J_{N_1}[f] \quad \text{y} \quad J^{h_2}[f] = J_{N_2}[f], \quad h_1 N_1 = h_2 N_2 = b - a.$$

Exijamos que el error para su combinación lineal

$$\tilde{D}^h(f) = D^{h_1}(f) + (1 - \sigma) D^{h_2}(f)$$

sea una magnitud de orden superior en comparación con D^{h_1} y D^{h_2} . Si para $D^h = D_N$ tiene lugar la fórmula del tipo

$$D^h = J^h(f) - J(f) = \alpha_p h^p + \alpha_q h^q + \dots, \quad q > p,$$

entonces para $\tilde{D}^h = (\sigma J^{h_1} [f] + (1 - \sigma) J^{h_2} [f]) - J [f]$ obtendremos

$$\tilde{D}^h (f) = \alpha_p (\sigma h_1^p + (1 - \sigma) h_2^p) + \alpha_q (\sigma h_1^q + (1 - \sigma) h_2^q) + \dots$$

Elijamos el parámetro σ , partiendo de la condición $\sigma h_1^p + (1 - \sigma) h_2^p = 0$:

$$\sigma = h_2^p / (h_2^p - h_1^p).$$

En este caso tendremos

$$\tilde{D}^h (f) = \alpha_q (\sigma h_1^q + (1 - \sigma) h_2^q) + \dots = O(h^q), \quad h = \max(h_1, h_2),$$

con la particularidad de que $\sigma h_1^q + (1 - \sigma) h_2^q < 0$. Así, por ejemplo, si $p = 2$, $q = 4$, entonces $\tilde{D}^h (f) = -\alpha_4 h_1^2 h_2^2 + \dots = O(h^4)$. De este modo, al realizar los cálculos sobre dos redes con los pasos h_1 y $h_2 \neq h_1$, hemos aumentado el orden de exactitud en 2 (en $q - p$) para $\tilde{J} = \sigma J^{h_1} + (1 - \sigma) J^{h_2}$.

Observemos que combinando la fórmula del trapecio $J_{\text{trap}}^{2h} [f]$ y la del rectángulo $J_{\text{rect}}^{2h} [f]$, ambas con paso $2h$, obtendremos la fórmula de Simpson J_{Simp}^h con paso h :

$$\begin{aligned} J_{\text{Simp}}^h [f] &= \frac{1}{3} J_{\text{trap}}^{2h} [f] + J_{\text{rect}}^{2h} [f] = \\ &= \frac{h}{6} (f_0 + 4f_1 + 2f_2 + \dots + 2f_{2N-2} + 4f_{2N-1} + f_{2N}), \end{aligned}$$

donde $h = (b - a)/(2N)$.

El método de cálculo sobre varias redes se aplica para el aumento del orden de exactitud incluso en aquel caso cuando se desconoce el orden del término principal del error (proceso de Aitken). Supongamos que para el error tiene lugar la representación

$$D^h (f) = \alpha_p h^p + \alpha_q h^q + \dots, \quad q > p,$$

de suerte que

$$J^h [f] = J [f] + \alpha_p h^p + \alpha_q h^q + \dots$$

Realicemos los cálculos sobre tres redes: $h_1 = h$, $h_2 = \rho h$, $h_3 = \rho^2 h$ ($0 < \rho < 1$). Determinemos, al principio, p , me-

nospreciando el término $O(h^q)$. Formemos una razón

$$A = \frac{J^{h_1}[f] - J^{h_2}[f]}{J^{h_2}[f] - J^{h_3}[f]} \approx \frac{h_1^p - h_2^p}{h_2^p - h_3^p} = \frac{1 - \rho^p}{\rho^p(1 - \rho^p)} = \left(\frac{1}{\rho}\right)^p$$

y hallemos

$$p \approx \ln A / \ln \frac{1}{\rho}.$$

Sabiendo el valor aproximado de p , se puede, empleando el método de Runge expuesto más arriba, aumentar el orden de exactitud. Con este fin formemos una combinación $\tilde{J}^h = \sigma J^{h_1} + (1 - \sigma) J^{h_2}$, y elijamos σ de una manera tal que se verifique $\sigma h_1^p + (1 - \sigma) h_2^p = (\sigma + (1 - \sigma) \rho^p) h^p = 0$, es decir, $\sigma = \rho^p / (\rho^p - 1) = 1 / (1 - A)$. Entonces, para el error $\tilde{D}^h = \tilde{J}^h - J$ obtenemos

$$\tilde{D}^h(f) = O(h^q).$$

Todos estos razonamientos tienen sentido, por supuesto, si la función $f(x)$ tiene una suavidad correspondiente.

7. Fórmulas de cuadratura de otro tipo. Sin perturbar la generalidad de razonamientos podemos considerar

$$J[f] = \int_0^1 f(x) dx. \quad (25)$$

Hasta ahora se han analizado las fórmulas de cuadratura con los nodos prefijados $\{x_k\}$:

$$J_N[f] = \sum_{k=0}^N c_k f(x_k). \quad (26)$$

Las fórmulas citadas son exactas para todos los polinomios de grado N . Si se consideran desconocidos no sólo c_k , sino también los nodos x_k , podemos exigir que la fórmula de cuadratura (26) sea exacta para todos los polinomios de grado $2N - 1$. La fórmula de esta índole lleva el nombre de Gauss. Exigiendo que para los monomios $1, x, x^2, \dots$

\dots, x^m, \dots, x^N la fórmula sea exacta, es decir, que

$$J_N [x^m] = \sum_{h=0}^N c_h x_h^m = \int_0^1 x^m dx = \frac{x^{m+1}}{m+1} \Big|_0^1 = \\ = \frac{1}{m+1}, \quad m = 0, 1, \dots, 2N-1,$$

obtendremos $2N + 2$ ecuaciones para los nodos y pesos

$$\begin{aligned} c_0 + c_1 + \dots + c_N &= 1 \\ c_0 x_0 + c_1 x_1 + \dots + c_N x_N &= 1/2, \\ \dots & \\ c_0 x_0^m + c_1 x_1^m + \dots + c_N x_N^m &= 1/(m+1), \\ \dots & \\ c_0 x_0^{2N+1} + c_1 x_1^{2N+1} + \dots + c_N x_N^{2N+1} &= 1/(N+1). \end{aligned}$$

El número total de las incógnitas es igual a $2N + 2$, es decir, $N + 1$ nodos y $N + 1$ factores ponderales desconocidos. El número de ecuaciones es también igual a $2N + 2$. Se puede demostrar que el sistema escrito de ecuaciones tiene la solución.

Aduzcamos la fórmula de Gauss más sencilla para $N = 2$:

$$J_N [f] = \frac{5}{18} f(x_0) + \frac{8}{18} f(x_1) + \frac{5}{18} f(x_2),$$

donde

$$x_0 = \frac{1 - \sqrt{0,6}}{2}, \quad x_2 = \frac{1 + \sqrt{0,6}}{2}, \quad x_1 = \frac{1}{2}.$$

Las fórmulas de Gauss proporcionan buena precisión con el número reducido de nodos.

De un ejemplo más sirve la fórmula de cuadratura de *Chébishev* en la que se eligen los mejores nodos bajo la suposición de que todos los pesos son iguales. En este caso

$$J_N [f] = \frac{1}{N} \sum_{h=1}^N f(x_h).$$

Exigiendo que la fórmula sea exacta para $f(x) = x, x^2, \dots, x^N$, obtendremos N ecuaciones para determinar x_1, x_2, \dots

... , x_N :

$$x_1^m + x_2^m + \dots + x_N^m = \frac{1}{m+1}, \quad m = 1, 2, \dots, N.$$

Estas ecuaciones tienen soluciones para $m = 1, 2, \dots, 7, 9$, y para $m = 8$ y $m \geq 10$ no tienen raíces reales. Cuando $m = 3$, la fórmula de Chébishev tiene por expresión

$$\int_0^1 f(x) dx \approx J_3[f] = \frac{1}{3} \left[f\left(\frac{1}{2} - \frac{1}{4}\sqrt{2}\right) + f\left(\frac{1}{2}\right) + f\left(\frac{1}{2} + \frac{1}{4}\sqrt{2}\right) \right].$$

Para ella el coeficiente de $\|f^{IV}\|_C$ en la estimación del error es dos veces menor que para la fórmula de Simpson.

OBSERVACIONES. En ciertos casos al cálculo de las integrales le debe preceder su transformación, teniendo en cuenta los rasgos específicos de la función subintegral. Ejemplos:

1) $f(x) = \frac{1}{2\sqrt{x}} f_0(x)$, $f_0(0) \neq 0$, es decir, $f(x)$ tiene

una singularidad cuando $x = 0$. Esta singularidad se elimina por el cambio de variable:

$$\begin{aligned} \int_0^1 f(x) dx &= \int_0^1 \frac{f_0(x)}{2\sqrt{x}} dx = \int_0^1 f_0(x) d\sqrt{x} \\ &= \int_0^1 f(t^2) dt, \quad t = \sqrt{x}. \end{aligned}$$

2) La función subintegral tiene el carácter exponencial: $f(x) \approx ce^{ax}$, es decir, la función $\ln f(x)$ es lineal. Representemos $f(x)$ como $f(x) = \exp\{\ln f(x)\}$, interpolemos $\ln f(x)$ linealmente en el segmento $[x_{l-1}, x_l]$:

$$\ln f(x) = \frac{x_l - x}{x_l - x_{l-1}} \ln f_{l-1} + \frac{x - x_{l-1}}{x_l - x_{l-1}} \ln f_l,$$

e integremos después respecto de x entre x_{l-1} y x_l . Esta fórmula resulta útil en los cálculos prácticos.

3) Si $f(x)$ es una función rápidamente oscilante, de modo que puede ser escrita en la forma $f(x) = y(x) \cos \omega x$, donde

la frecuencia $\omega \gg 1$ es grande, entonces calculando una integral se puede recurrir al procedimiento siguiente. Primero integramos por partes:

$$\begin{aligned} \int_{x_{i-1}}^{x_i} f(x) dx &= \int_{x_{i-1}}^{x_i} y(x) \cos \omega x dx = \\ &= \frac{1}{\omega} y \operatorname{sen} \omega x \Big|_{x_{i-1}}^{x_i} - \frac{1}{\omega} \int_{x_{i-1}}^{x_i} y'(x) \operatorname{sen} \omega x dx. \end{aligned}$$

Si $y(x)$ es lineal en $[x_{i-1}, x_i]$, entonces la integral del segundo miembro se calcula en la forma explícita. Si $y(x)$ es un polinomio de grado n , la integración por partes se realiza n veces.

Resolución numérica de los sistemas de ecuaciones algebraicas lineales

En este capítulo se estudian los métodos de resolución numérica de los sistemas de ecuaciones algebraicas lineales, es decir, los métodos numéricos del álgebra lineal. Existen dos tipos de métodos: directos e iterativos. Analizamos, ante todo, el método de eliminación de Gauss para los sistemas del tipo general y las variantes: método de factorización y métodos de factorización matricial para los sistemas del tipo especial (con matrices tridiagonal o tridiagonal por bloques). Estos son los *métodos directos*. Su eficiencia depende del orden del sistema y de la estructura de la matriz.

Al estudiar los métodos *iterativos*, consideramos todo sistema de ecuaciones como una ecuación operacional de la primera especie $Au = f$, y exponemos la teoría general de los métodos iterativos para ecuaciones operacionales con suposiciones mínimas respecto del operador A . La teoría general permite demostrar la convergencia de las iteraciones para el método de Seidel y para el método de relajación superior con restricciones mínimas para el operador A . Se han analizado dos clases de los métodos: 1) para el caso en que se conocen las fronteras $\gamma_1 > 0$ y $\gamma_2 \geq \gamma_1$ del espectro del operador A en cierto espacio energético H_D ; 2) para el caso en que las fronteras γ_1 y γ_2 no se conocen. Es de gran eficacia el método triangular alternado que se estudia en el § 5.

§ 1. Sistemas de ecuaciones algebraicas lineales

1. **Sistemas de ecuaciones.** El problema fundamental del álgebra lineal consiste en la resolución del sistema de ecuaciones

$$Au = f, \quad (1)$$

donde $u = (u^{(1)}, \dots, u^{(N)})$ es el vector buscado, $f = (f^{(1)}, f^{(2)}, \dots, f^{(N)})$ es un vector conocido de dimensión N , $A = (a_{ij})$ ($i, j = 1, 2, \dots, N$) es una matriz cuadrada de dimensión $N \times N$ con elementos a_{ij} .

Se supondrá que la matriz A es regular, $\det A \neq 0$, de modo que la ecuación $Au = 0$ tiene sólo una solución trivial, y el sistema (1) tiene la única solución

$$u = A^{-1}f.$$

En el curso de álgebra lineal la solución del sistema (1) se expresa, corrientemente, según las fórmulas de Cramer como una razón de los determinantes. Dichas fórmulas no sirven para la resolución numérica del sistema (1), puesto que requieren el cálculo de $N + 1$ determinantes, lo que, a su vez, exige un gran número de operaciones aritméticas (hasta $N!$). Si incluso escogemos el mejor método, para el cálculo de un solo determinante se necesitará aproximadamente tanto tiempo que se requiere para la resolución de un sistema de ecuaciones lineales por los métodos numéricos modernos. Además, hemos de tener en cuenta, que los cálculos según las fórmulas de Cramer conducen con frecuencia a los grandes errores de redondeo.

La peculiaridad de la mayoría de los métodos numéricos para (1) consiste en que se abandona la idea de buscar la matriz inversa. El requisito principal que se levanta ante el método de resolución es el mínimo de operaciones aritméticas suficientes para la búsqueda de una solución aproximada con la precisión prefijada $\epsilon > 0$ (eficiencia del método numérico).

2. Casos particulares de los sistemas. No es difícil resolver el sistema (1) en los casos particulares que van abajo. Sea A una matriz *diagonal*, es decir, $a_{ij} = 0$, $j \neq i$, $a_{ii} \neq 0$ ($i, j = 1, 2, \dots, N$). Entonces, el sistema tiene por expresión

$$a_{ii}u^{(i)} = f^{(i)},$$

de donde encontramos

$$u^{(i)} = f^{(i)}/a_{ii}, \quad i = 1, 2, \dots, N.$$

Esta es la ecuación de segundo orden que ha sido analizado en el cap. I, donde para su resolución se aplicó el método de factorización.

3. Ecuación operacional de primera especie. Se sabe que toda matriz $A = (a_{ij})$ ($i, j = 1, 2, \dots, N$) define un operador lineal A que aplica el espacio H^N en sí mismo: $Au \in H^N$ para cualquier $u \in H^N$, o bien $A: H^N \rightarrow H^N$. Viceversa e todo operador A (en cierta base ξ_1, \dots, ξ_N) le corresponde una matriz $A = (a_{ij})$ de dimensión $N \times N$, donde a_{ij} es el componente del vector $A\xi_j$. Por eso, la ecuación (1) puede considerarse como una *ecuación operacional de primera especie*

$$Au = f, \quad u, f \in H^N,$$

con el operador $A: H^N \rightarrow H^N$.

Con el fin de recalcar la equivalencia de los problemas (1) y (3), dejamos invariable la designación A tanto para la matriz como para el operador. Omitiremos el índice N en H^N y escribiremos simplemente H . El paso de (1) a la ecuación operacional resulta cómodo para la exposición de la teoría de métodos iterativos. En este caso no se emplea ninguna información concreta sobre la estructura de la matriz A .

En el espacio H introduzcamos un producto escalar (\cdot, \cdot) , y una norma $\|u\| = \sqrt{(u, u)}$. Supondremos que el operador A es autoconjugado y positivo: $A = A^* > 0$. Analizaremos también el espacio energético H_D con el producto escalar $(u, v)_D = (Du, v)$ y la norma $\|u\|_D = \sqrt{(Du, u)}$, donde D es un operador lineal positivo y autoconjugado: $D: H \rightarrow H$, $D = D^* > 0$.

Denotemos con (ξ_s, λ_s) ($s = 1, 2, \dots, N$) los vectores propios y los valores propios del operador A :

$$A\xi_s = \lambda_s \xi_s, \quad (\xi_s, \xi_m) = \delta_{sm}, \\ s, m = 1, 2, \dots, N.$$

Por cuanto $A > 0$, se tiene $\lambda_s > 0$, y podemos considerar que $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_N$, y, por consiguiente, se verifica la desigualdad

$$\lambda_1 E \leq A \leq \lambda_N E, \quad \lambda_1 = \min_s \lambda_s, \quad \lambda_N = \max_s \lambda_s.$$

La razón λ_N/λ_1 lleva el nombre de *número convenido*.

En la práctica resulta más conveniente emplear la razón inversa, es decir, el parámetro $\xi = \lambda_1/\lambda_N$, el cual se denominará *medida convenida*. En lo sucesivo se mostrará que de dicho parámetro depende la convergencia de las iteraciones. El parámetro ξ para las ecuaciones en diferencias que aproximan las ecuaciones de la física matemática (por ejemplo, la ecuación de Laplace) es pequeño: $\xi \approx 10^{-2}-10^{-4}$ (el número convenido es grande).

De la fórmula $u = A^{-1}f$ se ve que

$$\|u\| \leq \|A^{-1}\| \|f\|, \quad \|A^{-1}\| = 1/\lambda_1.$$

Esta desigualdad expresa la estabilidad de la solución del problema (1) respecto del segundo miembro. Si $\|A^{-1}\| = 1/\lambda_1$ es muy grande, puede suceder que el problema (3) no sea correcta, es decir, inestable con relación a los errores en la prefijación del segundo miembro, incluso a los errores de redondeo.

4. Métodos directos e iterativos. Los métodos numéricos de la resolución del sistema (1) se subdividen convencionalmente en dos grupos y se distinguen métodos directos y los iterativos (se tienen, por supuesto, los métodos mixtos). Los métodos *directos* permiten obtener, después de un número finito de operaciones, una solución exacta del sistema de ecuaciones, siempre que la información de entrada (el segundo miembro de la ecuación f y los elementos a_{ij} de la matriz A) viene dada con toda la exactitud y los cálculos se realizan sin redondeo. El ejemplo más simple del método directo es el de factorización. Por supuesto, los métodos directos dan también la solución con cierta precisión, la que depende de los errores de redondeo, es decir, del ordenador y del carácter de la estabilidad de cálculo, lo que depende, a su vez, del método mismo.

El método iterativo permite hallar la solución aproximada del sistema construyendo una sucesión de aproximaciones (iteraciones), a partir de cierta aproximación inicial. La propia solución aproximada es el resultado de los cálculos obtenido después de haberse realizado un número finito de iteraciones.

La elección de tal o cual método numérico depende de varias circunstancias: de los programas disponibles, del tipo de los cálculos y de la matriz A , etc. Expliquemos las

palabras «tipo de los cálculos». Son posibles distintos planteamientos del problema:

- 1) hallar la solución de un problema concreto (1);
- 2) hallar la solución de varias variantes del problema (1) con una misma matriz A y segundos miembros diferentes de f . Puede ocurrir que un método, no optimal para un problema (1), resulte muy eficaz para el cálculo multivariante.

En el cálculo multivariante se puede disminuir el número medio de operaciones para una variante, si se conservan ciertas magnitudes y no se calculan de nuevo para cada variante, lo que depende, naturalmente, del ordenador y del volumen de su memoria de acceso rápido.

De aquí está claro que la elección de un algoritmo debe depender del tipo de los cálculos, del volumen de la memoria de acceso rápido del ordenador y, naturalmente, del orden del sistema. La calidad de un algoritmo se determina por el tiempo de máquina que se exige para hallar la solución del sistema (1). Se elige, naturalmente, un método, para el cual el tiempo de resolución es mínimo en comparación con los otros métodos. No obstante, el tiempo de cálculo depende de varios factores, entre cuales pueden citarse el número de operaciones aritméticas y lógicas que son necesarias para obtener la solución con la exactitud prefijada, la velocidad de funcionamiento y el volumen de la memoria de acceso rápido del ordenador, la calidad del programa. Al estimar teóricamente la calidad de los algoritmos su comparación se realiza por el número $Q(\varepsilon)$ de operaciones aritméticas suficientes para hallar la solución del problema con la exactitud prefijada $\varepsilon > 0$.

§ 2. Métodos directos

1. **Método de Gauss.** Hay varias variantes de cálculo del método de Gauss basado en la idea de eliminación sucesiva. El proceso de resolución del sistema de ecuaciones algebraicas lineales $Ax = f$, o

$$\sum_{j=1}^N a_{ij}x_j = f_i, \quad i = 1, 2, \dots, N, \quad (1)$$

por el método de Gauss, consta de dos etapas.

PRIMERA ETAPA (*procedimiento directo*). El sistema (1) se reduce a la forma triangular

$$x + B^+x = \varphi, \quad (2)$$

donde $x = (x_1, \dots, x_N)$ y $\varphi = (\varphi_1, \dots, \varphi_N)$ son los vectores desconocido y conocido, respectivamente, B^+ es la matriz triangular superior.

SEGUNDA ETAPA (*procedimiento inverso*). Las incógnitas x_N, x_{N-1}, \dots, x_1 se determinan por las fórmulas (2) del § 1.

Pasemos a la exposición detallada del método. El primer paso del método de Gauss consiste en la eliminación de la incógnita x_1 de todas las ecuaciones con excepción de la primera. Supongamos que $a_{11} \neq 0$, dividamos la primera ecuación (1) ($i = 1$) por a_{11} y escribamos el sistema (1) en la forma

$$x_1 + b_{12}x_2 + \dots + b_{1N}x_N = \varphi_1, \quad b_{1j} = a_{1j}/a_{11}, \\ 2 \leq j \leq N, \quad \varphi_1 = f_1/a_{11}, \quad (3)$$

$$a_{i1}x_1 + a_{i2}x_2 + \dots + a_{iN}x_N = f_i, \quad i = 2, 3, \dots, N. \quad (4)$$

Multipliquemos la ecuación (3) por a_{i1} , donde i es cualquiera de los números $i = 2, 3, \dots, N$, y sustrayamos la ecuación obtenida de la i -ésima ecuación (4):

$$(a_{i2} - a_{i1}b_{12})x_2 + \dots + (a_{iN} - a_{i1}b_{1N})x_N = f_i - a_{i1}\varphi_1, \\ i = 2, 3, \dots, N.$$

Introduciendo las designaciones

$$a_{ij}^{(1)} = a_{ij} - a_{i1}b_{1j}, \quad f_i^{(1)} = f_i - a_{i1}\varphi_1, \\ i, j = 2, 3, \dots, N, \quad (5)$$

reescribamos el sistema obtenido de ecuaciones (que es equivalente al sistema (1)) en la forma

$$x_1 + b_{12}x_2 + \dots + b_{1N}x_N = \varphi_1, \\ a_{i2}^{(1)}x_2 + \dots + a_{iN}^{(1)}x_N = f_i^{(1)}, \quad i = 2, 3, \dots, N.$$

La primera columna de la matriz de este sistema se compone de ceros, a excepción del primer elemento para $i = 1, j = 1$, que es igual a uno.

El *paso segundo* consiste en la eliminación x_2 del sistema

$$\begin{aligned} a_{22}^{(1)}x_2 + \dots + a_{2N}^{(1)}x_N &= f_2^{(1)}, \\ \dots & \\ \varphi_{N2}^{(1)}x_2 + \dots + a_{NN}^{(1)}x_N &= f_N^{(1)}. \end{aligned} \quad (6)$$

Con este objeto dividamos la primera ecuación por $a_{22}^{(1)}$:

$$\begin{aligned} x_2 + b_{23}x_3 + \dots + b_{2N}x_N &= \varphi_2, \\ \varphi_2 = f_2^{(1)}/a_{22}^{(1)}, \quad b_{2j} &= a_{2j}^{(1)}/a_{22}^{(1)}, \quad j = 3, \dots, N, \end{aligned}$$

multipliquémosla después por $(-a_{i2}^{(1)})$ y sumemos con la ecuación

$$a_{i2}^{(1)}x_2 + a_{i3}^{(1)}x_3 + \dots + a_{iN}^{(1)}x_N = f_i^{(1)}, \quad i = 3, 4, \dots, N.$$

De resultas obtendremos un sistema

$$\begin{aligned} x_2 + b_{23}x_3 + \dots + b_{2N}x_N &= \varphi_2, \\ a_{i3}^{(2)}x_3 + \dots + a_{iN}^{(2)}x_N &= f_i^{(2)}, \quad i = 3, 4, \dots, N, \end{aligned} \quad (7)$$

$$\begin{aligned} a_{ij}^{(2)} &= a_{ij}^{(1)} - a_{i2}^{(1)}b_{2j}, \quad f_i^{(2)} = f_i^{(1)} - a_{i2}^{(1)}\varphi_2, \\ i &= 3, 4, \dots, N. \end{aligned} \quad (8)$$

Para x_3, x_4, \dots, x_N tenemos un sistema de $(N - 2)$ -ésimo orden análogo al sistema (6) de $(N - 1)$ -ésimo orden para x_2, x_3, \dots, x_N .

Continuando los razonamientos, obtendremos tras el $(N - 1)$ -ésimo *paso* (es decir, al haber excluido x_1, x_2, \dots, x_{N-1})

$$a_{NN}^{(N-1)}x_N = f_N^{(N-1)}, \quad \text{o bien } x_N = \varphi_N, \quad \varphi_N = f_N^{(N-1)}/a_{NN}^{(N-1)}. \quad (9)$$

Llegamos en fin al sistema (2) con la matriz triangular superior

$$\begin{aligned} x_1 + b_{12}x_2 + b_{13}x_3 + \dots + b_{1N}x_N &= \varphi_1, \\ x_2 + b_{23}x_3 + \dots + b_{2N}x_N &= \varphi_2, \\ \dots & \\ x_{N-1} + b_{N-1, N}x_N &= \varphi_{N-1}, \\ x_N &= \varphi_N. \end{aligned} \quad (10)$$

El procedimiento inverso del método de Gauss consiste en determinar todos los x_i pertenecientes al sistema (10) con la matriz triangular superior. No es difícil mostrar que el método de Gauss expuesto más arriba puede aplicarse solamente en aquel caso en que todos los menores principales son distintos de cero.

Contaremos el número de multiplicaciones y divisiones en el método de Gauss. Veamos primero el procedimiento directo. En el primer paso se requieren $Q_1 = N^2$ divisiones y multiplicaciones, el segundo paso exige $Q_2 = (N - 1)^2$ operaciones, etc. En total se deben hacer N pasos del procedimiento directo realizando para ello

$$\sum_{k=1}^N (N - k + 1)^2 = \sum_{s=1}^N s^2 = \frac{N(N+1)(2N+1)}{6}$$

multiplicaciones y divisiones. Es evidente que en el procedimiento inverso se deben realizar $N(N - 1)/2$ multiplicaciones. De este modo, para resolver el sistema de ecuaciones (1) se necesitan $Q = N(N^2 + 3N - 1)/3$ operaciones de multiplicación y división. Se necesitarán también aproximadamente el mismo número de las operaciones de suma-ción.

Demos a conocer un ejemplo de aplicación del método de Gauss. Examinemos un sistema de tres ecuaciones ($N = 3$)

$$2x_1 + 4x_2 + 3x_3 = 4, \quad (11)$$

$$3x_1 + x_2 - 2x_3 = -2, \quad (12)$$

$$4x_1 + 11x_2 + 7x_3 = 7. \quad (13)$$

PROCEDIMIENTO DIRECTO. PRIMER PASO. Dividamos la primera ecuación por $a_{11} = 2$:

$$x_1 + 2x_2 + 1,5x_3 = 2. \quad (14)$$

Multipliquemos (14) por -3 y sumemos con (12), a continuación multipliquemos (14) por -4 y sumemos con (13):

$$-5x_2 - 6,5x_3 = -8, \quad (15)$$

$$3x_2 + x_3 = 1. \quad (16)$$

Se ha obtenido el sistema de segundo orden.

SEGUNDO PASO. Dividamos (15) por -5 :

$$x_2 + 1,3x_3 = 1,6. \quad (17)$$

Multipliquemos (17) por -3 y sumemos con (16):

$$-2,9x_3 = -5,8. \quad (18)$$

TERCER PASO. Dividamos (18) por $-2,9$:

$$x_3 = 2.$$

De resultas obtenemos un sistema

$$\begin{aligned} x_1 + 2x_2 + 1,5x_3 &= 2, \\ x_2 + 1,3x_3 &= 1,6, \\ x_3 &= 2 \end{aligned}$$

con la matriz triangular superior

$$\begin{bmatrix} 1 & 2 & 1,5 \\ 0 & 1 & 1,3 \\ 0 & 0 & 1 \end{bmatrix}.$$

PROCEDIMIENTO INVERSO. Del sistema hallamos sucesivamente: $x_3 = 2$, $x_2 = 1,6 - 1,3 \cdot x_3 = 1,6 - 1,3 \cdot 2 = -1$, $x_1 = 2 - 2x_2 - 1,5 \cdot x_3 = 1$. De este modo, queda determinada la solución del sistema (11)–(13):

$$x_1 = 1, \quad x_2 = -1, \quad x_3 = 2.$$

2. Método de la raíz cuadrada. Este método se emplea para los sistemas

$$Au = f \quad (19)$$

con matriz hermitiana A (simétrica, en el caso real). La matriz A se desarrolla en un producto

$$A = S^*DS, \quad (20)$$

donde S y D son matrices superior triangular y diagonal, respectivamente. La resolución de la ecuación $Au = f$ se reduce a la resolución de dos sistemas

$$S^*Dy = f, \quad Su = y. \quad (21)$$

Con el fin de obtener el desarrollo (20), designamos $S = (s_{ij})$, $D = (d_{ii}\delta_{ij})$ y hallamos

$$(DS)_{ij} = \sum_{k=1}^N d_{ik}t_{kj} = d_{ii}s_{ij}, \quad (S^*DS)_{ij} = \sum_{k=1}^N \bar{s}_{ki}d_{kk}s_{kj},$$

puesto que $S^* = (\bar{s}_{ij})$, donde la raya significa una conjugación compleja.

De resultas obtenemos una ecuación

$$\sum_{k=1}^N \bar{s}_{ki}d_{kk}s_{kj} = a_{ij}. \quad (22)$$

El sistema de ecuaciones (22) se puede resolver de manera recurrente. Por cuanto S es una matriz triangular superior, entonces $s_{ki} = 0$ para $k > i$, $\bar{s}_{ih} = 0$ para $k < i$, y, por lo tanto,

$$\begin{aligned} \sum_{k=1}^N \bar{s}_{ki}s_{kj}d_{kk} &= \sum_{k=1}^{i-1} \bar{s}_{ki}s_{kj}d_{kk} + \bar{s}_{ii}s_{ij}d_{ii} + \sum_{k=i+1}^N \bar{s}_{ki}s_{kj}d_{kk} = \\ &= \sum_{k=1}^{i-1} \bar{s}_{ki}s_{kj}d_{kk} + s_{ii}s_{ij}d_{ii} = a_{ij}. \end{aligned}$$

Para $i = j$ tenemos

$$|s_{ii}|^2 d_{ii} = a_{ii} - \sum_{k=1}^{i-1} |s_{ki}|^2 d_{kk}. \quad (23)$$

Al escoger

$$d_{ii} = \text{sign}(a_{ii} - \sum_{k=1}^{i-1} |s_{ki}|^2 d_{kk}), \quad (24)$$

hallemos

$$s_{ii} = \sqrt{|a_{ii} - \sum_{k=1}^{i-1} |s_{ki}|^2 d_{kk}|}. \quad (25)$$

Cuando $i < j$, obtenemos

$$s_{ij} = \frac{a_{ij} - \sum_{k=1}^{i-1} \bar{s}_{ki}s_{kj}d_{kk}}{\bar{s}_{ii}d_{ii}}. \quad (26)$$

Suponiendo $i = 1, 2, \dots$, encontramos sucesivamente $s_{11} = \sqrt{|a_{11}|}$, $d_{11} = \text{sign } a_{11}$, $s_{22} = \sqrt{|a_{22} - d_{11}| |s_{12}|^2}$, ...

El determinante de la matriz es, evidentemente, igual a

$$\det A = \prod_{i=1}^N d_{ii} s_{ii}^2.$$

El método de la raíz cuadrada requiere aproximadamente $N^3/3$ operaciones aritméticas, es decir, cuando N es grande, el método es dos veces más rápido en comparación con el método de Gauss y ocupa dos veces menos células de la memoria. Esta circunstancia se debe a que el método emplea la información sobre la simetría de la matriz.

3. Relación del método de Gauss con el desarrollo de la matriz en factores. Sea dada una matriz regular A de dimensión $N \times N$. Representémosla en forma de un producto

$$A = BC, \quad A = (a_{ij}), \quad B = (b_{ij}), \quad C = (c_{ij}) \quad (27)$$

donde B y C son las matrices triangulares de la forma

$$B = \begin{bmatrix} b_{11} & 0 & \dots & 0 \\ b_{21} & b_{22} & 0 & \dots & 0 \\ b_{31} & b_{32} & b_{33} & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ b_{N1} & b_{N2} & b_{N3} & \dots & b_{NN} \end{bmatrix},$$

$$C = \begin{bmatrix} 1 & c_{12} & c_{13} & \dots & c_{1N} \\ 0 & 1 & c_{23} & \dots & c_{2N} \\ 0 & 0 & 1 & \dots & c_{3N} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix},$$

es decir, $b_{ik} = 0$ para $k > i$, $c_{ik} = 0$ para $k < i$, $c_{ii} = 1$. De (27) se infiere que

$$a_{ij} = \sum_{k=1}^N b_{ik} c_{kj}.$$

Transformemos esta suma de dos modos:

$$\sum_{h=1}^N b_{ik}c_{hj} = \sum_{h=1}^{i-1} b_{ik}c_{hj} + b_{ii}c_{ij} + \sum_{h=i+1}^N b_{ik}c_{hj} =$$

$$= \sum_{h=1}^{i-1} b_{ik}c_{hj} + b_{ii}c_{ij}$$

$$\sum_{h=1}^N b_{ik}c_{hj} = \sum_{h=1}^{j-1} b_{ik}c_{hj} + b_{ij}c_{jj} + \sum_{h=j+1}^N b_{ik}c_{hj} =$$

$$= \sum_{h=1}^{j-1} b_{ik}c_{hj} + b_{ij}c_{jj}$$

De aquí encontramos

$$b_{ij} = a_{ij} - \sum_{h=1}^{j-1} b_{ik}c_{hj} \quad \text{para } i \geq j, \quad b_{ii} = a_{ii}, \quad c_{ii} = 1,$$

$$c_{ij} = \frac{1}{b_{ij}} \left[a_{ij} - \sum_{h=1}^{i-1} b_{ik}c_{hj} \right] \quad \text{para } i < j.$$

Las matrices B y C quedan determinadas.

La resolución de la ecuación $Au = BCu = f$ se reduce a la resolución sucesiva de las ecuaciones

$$B\varphi = f, \quad Cu = \varphi.$$

La construcción de las matrices B y C , como también la búsqueda de $\varphi = B^{-1}f$ corresponden al procedimiento directo, mientras que la resolución de la ecuación

$$Cu = \varphi$$

corresponde al procedimiento inverso del método de Gauss.

§ 3. Métodos iterativos

1. Método de iteraciones para resolver el sistema de ecuaciones algebraicas lineales. Una atención especial se prestará en este capítulo a los métodos iterativos, puesto que dichos métodos son de amplio uso en la resolución de las ecuaciones en diferencias de la física matemática cuyos operadores están en correspondencia con las matrices de cinta A de orden superior.

Pasemos a la descripción general del *método de iteraciones* para un sistema de ecuaciones algebraicas lineales

$$Au = f. \quad (1)$$

Con el fin de resolverlo se elige cierta aproximación inicial $y_0 \in H$ y se hallan sucesivamente soluciones aproximadas (iteraciones) de la ecuación (1). El valor de una *iteración* y_{k+1} se expresa en términos de las iteraciones precedentes conocidas y_k, y_{k-1}, \dots . Si, al calcular y_{k+1} , se utiliza sólo una iteración precedente y_k , entonces el método iterativo se denomina *de un paso (de dos capas)*; si, en cambio, y_{k+1} se expresa en términos de dos iteraciones, y_k e y_{k-1} , el método se llama *de dos pasos (o de tres capas)*. En esta obra se analizarán, principalmente, los métodos de un paso. Conviengamos en considerar que $A: H \rightarrow H$ es un operador lineal en un espacio de dimensión finita H con un producto escalar (\cdot, \cdot) .

Un papel importante lo desempeña la inscripción de los métodos iterativos en la forma unificada (canónica). Cualquier método iterativo de dos capas puede ser escrito en la siguiente forma canónica:

$$B \frac{y_{k+1} - y_k}{\tau_{k+1}} + Ay_k = f, \quad k=0, 1, \dots, \text{ para todos } y_0 \in H, \quad (2)$$

donde $A: H \rightarrow H$ es el operador de la ecuación de partida (1), $B: H \rightarrow H$, operador lineal que cuenta con su inversa B^{-1} , k es el número de la iteración; $\tau_1, \tau_2, \dots, \tau_{k+1}$ son todos los parámetros de iteración, $\tau_{k+1} > 0$. El operador B puede depender, en el caso general, del número k ; para simplificar la exposición, suponemos siempre que B no depende de k .

Si $B = E$ es un operador unidad, entonces el método

$$\frac{y_{k+1} - y_k}{\tau_{k+1}} + Ay_k = f, \quad k=0, 1, \dots, \text{ para todos los } y_0 \in H, \quad (3)$$

se denominará *explícito*: y_{k+1} se halla por una fórmula explícita

$$y_{k+1} = y_k - \tau_{k+1} (Ay_k - f).$$

Generalmente, cuando $B \neq E$, el método (2) se llama iterativo *implícito*: para determinar y_{k+1} hace falta resolver la ecuación

$$By_{k+1} = By_k - \tau_{k+1} (Ay_k - f) = F_k, \quad k = 0, 1, \dots \quad (4)$$

Es natural exigir que el volumen de los cálculos para resolver el sistema $By_{k+1} = F_k$ sea inferior al volumen de los cálculos para la resolución directa del sistema $Au = f$.

La exactitud del método iterativo (2) se caracteriza por la magnitud del error $z_k = y_k - u$, es decir, por la diferencia entre la solución y_k de la ecuación (2) y la solución exacta u del sistema inicial de ecuaciones algebraicas lineales. La sustitución $y_k = z_k + u$ en (2) lleva a una ecuación homogénea para el error:

$$B \frac{z_{k+1} - z_k}{\tau_{k+1}} + Az_k = 0, \quad k = 0, 1, \dots, \quad z_0 = y_0 - u. \quad (5)$$

Suele decirse que un método iterativo *converge* en H_D , si

$$\lim_{k \rightarrow \infty} \|z_k\|_D = 0, \quad \text{donde } \|z\|_D \sqrt{(Dz, z)}, \quad D = D^* > 0,$$

$$D: H \rightarrow H.$$

En el caso general se fija cierto error (relativo) $\varepsilon > 0$ con el que se debe hallar la solución aproximada y_n , los cálculos se dan por terminados cuando queda cumplida la condición

$$\|y_n - u\|_D \leq \varepsilon \|y_0 - u\|_D. \quad (6)$$

Si $n = n(\varepsilon)$ es el mínimo de los números, para los cuales se verifica (6), entonces el número total de operaciones aritméticas que han de realizarse para hallar la solución aproximada de la ecuación (1) es igual a $Q_n(\varepsilon) = n(\varepsilon) q_0$, donde q_0 es el número de operaciones que se realizan para hallar una iteración, es decir, para resolver la ecuación (4). El problema consiste en minimizar $Q_n(\varepsilon)$ eligiendo de modo adecuado B y los parámetros $\{\tau_k\}$. Comencemos por analizar los métodos iterativos más simples.

2. Método de la iteración simple. Para la resolución del sistema de ecuaciones (1) puede emplearse el *método de la*

iteración simple

$$y_{k+1}^{(i)} = y_k^{(i)} - \tau \left(\sum_{j=1}^N a_{ij} y_k^{(j)} - f^{(i)} \right), \quad i = 1, 2, \dots, N, \quad (7)$$

donde $\tau > 0$ es el parámetro de iteración. Escribamos (7) en la forma operacional:

$$\frac{y_{k+1} - y_k}{\tau} + Ay_k = f, \quad k = 0, 1, \dots, \text{ para cualquier } y_0 \in H. \quad (8)$$

Al comparar con (3) vemos que el método de la iteración simple se da mediante un esquema explícito de dos capas con el parámetro constante $\tau_k \equiv \tau$.

Existen también otras variantes del método de la iteración simple, por ejemplo, una que sigue:

$$y_{k+1}^{(i)} = \frac{1}{a_{ii}} \left(\sum_{j \neq i}^{1+N} a_{ij} y_k^{(j)} - f^{(i)} \right).$$

Al sustituir aquí

$$\sum_{j \neq i}^{1+N} a_{ij} y_k^{(j)} = \sum_{j=1}^N a_{ij} y_k^{(j)} - a_{ii} y_k^{(i)} = (Ay_k)^{(i)} - (Dy_k)^{(i)},$$

donde $D = (a_{ii} \delta_{ij})$ es una matriz diagonal, obtenemos

$$y_{k+1}^{(i)} = y_k^{(i)} - \frac{1}{a_{ii}} \left(\sum_{j=1}^N a_{ij} y_k^{(j)} - f^{(i)} \right),$$

o bien, en la forma canónica

$$D \frac{y_{k+1} - y_k}{\tau} + Ay_k = f, \quad k = 0, 1, \dots, \quad \tau = 1.$$

Aunque este esquema es formalmente implícito ($B = D \neq E$), no obstante $D = (a_{ii} \delta_{ij})$ es una matriz diagonal, por lo cual y_{k+1} se determina según las fórmulas explícitas.

3. Método de Seidel. Es de amplio uso (en particular, cuando es insuficiente la información sobre la matriz A) el

método iterativo de Seidel en una de las siguientes formas:

$$\sum_{j=1}^i a_{1j} y_{k+1}^{(j)} + \sum_{j=i+1}^N a_{1j} y_k^{(j)} = f^{(i)}, \quad a_{11} \neq 0, \quad i = 1, 2, \dots, N, \quad (9)$$

$$\sum_{j=1}^i a_{1j} y_k^{(j)} + \sum_{j=i+1}^N a_{1j} y_{k+1}^{(j)} = f^{(i)}, \quad i = 1, 2, \dots, N. \quad (10)$$

Los componentes del vector y_{k+1} se hallan sucesivamente de ambas fórmulas. Así, por ejemplo, de (9) determinamos sucesivamente $y_{k+1}^{(1)}, y_{k+1}^{(2)}, \dots, y_{k+1}^{(N)}$:

$$y_{k+1}^{(1)} = \frac{1}{a_{11}} \left(f^{(1)} - \sum_{j=2}^N a_{1j} y_k^{(j)} \right),$$

$$y_{k+1}^{(i)} = \frac{1}{a_{ii}} \left(f^{(i)} - \sum_{j=i+1}^N a_{ij} y_k^{(j)} - \sum_{j=1}^{i-1} a_{ij} y_{k+1}^{(j)} \right),$$

$$i = 2, \dots, N.$$

Haciendo uso de (10) encontramos sucesivamente para $i = N, N-1, \dots, 1$

$$y_{k+1}^{(N)} = \frac{1}{a_{NN}} \left(f^{(N)} - \sum_{i=1}^{N-1} a_{Ni} y_k^{(i)} \right),$$

$$y_{k+1}^{(i)} = \frac{1}{a_{ii}} \left(f^{(i)} - \sum_{j=1}^{i-1} a_{ij} y_k^{(j)} - \sum_{j=i+1}^N a_{ij} y_{k+1}^{(j)} \right),$$

$$i = N-1, \dots, 1.$$

Escribamos este método en la forma matricial (operacional). Con este fin representemos la matriz A como una suma

$$A = A^- + D + A^+,$$

donde $D = (a_{ii} \delta_{ij})$ es una matriz diagonal de dimensión $N \times N$, $A^- = (a_{ij}^-)$ es la matriz triangular (subdiagonal) inferior con la diagonal principal llenada de ceros, $a_{ij}^- = 0$ para $j \geq i$, $a_{ij}^- = a_{ij}$ para $j < i$, $A^+ = (a_{ij}^+)$ es la matriz triangular (sobrediagonal) superior con la diagonal principal

llenada de ceros, $a_{ij}^+ = 0$ para $j \leq i$, $a_{ij}^+ = a_{ij}$ para $j > i$. De la definición de A^- , D , A^+ se desprende que

$$Dy^{(i)} = a_{ii}y^{(i)}, \quad A^-y^{(i)} = \sum_{j=1}^{i-1} a_{ij}y^{(j)},$$

$$A^+y^{(i)} = \sum_{j=i+1}^N a_{ij}y^{(j)}, \quad (A^+ + D)y^{(i)} = \sum_{j=i}^N a_{ij}y^{(j)}.$$

Por esto, la ecuación (10) puede anotarse en la forma

$$((A^+ + D)y_{k+1})^{(i)} + (A^-y_k)^{(i)} = f^{(i)}, \quad i = 1, 2, \dots, N,$$

o bien, en la forma vectorial,

$$(A^+ + D)y_{k+1} + A^-y_k = f.$$

Realizadas unas transformaciones evidentes

$$(A^+ + D)y_{k+1} + A^-y_k = (A^+ + D)(y_{k+1} - y_k) + \\ + (A^- + (A^+ + D))y_k = (A^+ + D)(y_{k+1} - y_k) + Ay_k$$

escribamos el método de Seidel (10) en la forma canónica:

$$(D + A^+)(y_{k+1} - y_k) + Ay_k = f, \quad k = 0, 1, 2, \dots \quad (11)$$

Comparando con (2) vemos que el método de Seidel (10) corresponde a

$$B = D + A^+, \quad \tau \equiv 1,$$

es decir, el esquema (11) es implícito. Sin embargo, por cuanto $B = D + A^+$ es una matriz triangular, entonces la iteración y_{k+1} se determina por las fórmulas explícitas. Análogamente se escribe la otra variante del método de Seidel:

$$(D + A^-)(y_{k+1} - y_k) + Ay_k = f, \quad k = 0, 1, \dots, \quad (12)$$

cuando $B = D + A^-$ es una matriz triangular inferior. A continuación (en el p. 5) se mostrará que el método de Seidel converge, si A es una matriz simétrica definida positiva.

4. Método de relajación superior. Con el objeto de acelerar un proceso iterativo, se puede reducir el método de Seidel al de relajación superior, introduciendo el parámetro

de iteración ω de modo tal que se verifique

$$(D + \omega A^-) \frac{y_{k+1} - y_k}{\omega} + Ay_k = f, \\ k = 0, 1, \dots, \text{ para cualquier } y_0 \in H. \quad (13)$$

Al comparar con (2) vemos que

$$B = D + \omega A^-, \quad \tau = \omega.$$

Transformemos la ecuación (13) a la forma de cálculo. Teniendo presente que

$$(D + \omega A^-) \frac{y_{k+1} - y_k}{\omega} + Ay_k = \left(A^- + \frac{1}{\omega} D \right) y_{k+1} + \\ + \left(A - A^- - \frac{D}{\omega} \right) y_k = \left(A^- + \frac{1}{\omega} D \right) y_{k+1} + \\ + \left(A^* + \left(1 - \frac{1}{\omega} \right) D \right) y_k,$$

tenemos

$$\left(A^- + \frac{1}{\omega} D \right) y_{k+1} + \left(A^* + \left(1 - \frac{1}{\omega} \right) D \right) y_k = f.$$

De aquí encontramos

$$y_{k+1}^{(i)} = y_k^{(i)} + \frac{\omega}{a_{ii}} \left[f^{(i)} - \sum_{j=1}^{i-1} a_{ij} y_{k+1}^{(j)} - \sum_{j=i}^N a_{ij} y_k^{(j)} \right], \\ i = 1, 2, \dots, N.$$

Cuando $\omega = 1$, obtenemos la fórmula del método de Seidel.

La velocidad de convergencia del método de relajación superior depende del parámetro ω . En el p. 5 se mostrará que para la convergencia del método se ha de exigir que $0 < \omega < 2$.

5. Convergencia de los métodos iterativos estacionarios.

El método de Seidel y el de relajación superior sirven de ejemplo de los esquemas implícitos de la forma

$$B \frac{y_{k+1} - y_k}{\tau} + Ay_k = f, \quad k = 0, 1, \dots, \\ \text{para cualquier } y_0 \in H, \quad (14)$$

con el operador no autoconjugado B que tiene su inverso B^{-1} . El método (14) lleva el nombre de iterativo *estacionario*, puesto que B y τ no dependen del número de iteración. Para que exista el operador inverso B^{-1} , es suficiente exigir que el operador B sea positivo. Sea $B = D + \omega A^{-}$. Por cuanto $A = A^* > 0$, entonces $(A^{-}y, y) = (A^{+}y, y)$, $(A^{+})^* = A^{-}$, y, por consiguiente, $(Ay, y) = (Dy, y) + 2(A^{-}y, y)$, es decir,

$$(A^{-}y, y) = \frac{1}{2} ((A - D)y, y).$$

Sustituyendo esta expresión en la fórmula $(By, y) = (Dy, y) + \omega (A^{-}y, y)$, hallamos

$$(By, y) = \left(1 - \frac{1}{2}\omega\right) (Dy, y) + \omega (Ay, y) > 0,$$

siempre que $0 < \omega < 2$.

Para el error $z_k = y_k - u$ obtenemos una ecuación homogénea

$$B \frac{z_{k+1} - z_k}{\tau} + Az_k = 0, \quad k = 0, 1, 2, \dots, \quad z_0 = y_0 - u. \quad (15)$$

TEOREMA 1. Sea A un operador autoconjugado y positivo y suponemos cumplida la condición

$$B > \frac{\tau}{2} A. \quad (16)$$

En este caso el método de iteraciones (14) converge en H_A , es decir,

$$\|z_k\|_A = \|y_k - u\|_A \rightarrow 0 \quad \text{cuando} \quad k \rightarrow \infty.$$

DEMOSTRACION. Nos hará falta una identidad energética

$$2\tau \left(\left(B - \frac{\tau}{2} A \right) \frac{z_{k+1} - z_k}{\tau}, \frac{z_{k+1} - z_k}{\tau} \right) + \|z_{k+1}\|_A^2 = \|z_k\|_A^2, \quad (17)$$

donde $\|z\|_A^2 = (Az, z)$. Transformemos primero la ecuación (15) a la forma

$$\left(B - \frac{\tau}{2} A \right) \frac{z_{k+1} - z_k}{\tau} + \frac{1}{2} A (z_k + z_{k+1}) = 0, \quad (18)$$

sustituyendo con este fin $z_k = \frac{1}{2}(z_{k+1} + z_k) - \frac{\tau}{2} \frac{(z_{k+1} - z_k)}{\tau}$.

Al multiplicar (18) escalarmente por $2\tau \left(\frac{z_{k+1} - z_k}{\tau} \right) = 2(z_{k+1} - z_k)$ y teniendo presente que $(Az_{k+1}, z_k) = (z_{k+1}, Az_k)$, puesto que $A = A^*$ y $(A(z_k + z_{k+1}), z_{k+1} - z_k) = (Az_{k+1}, z_{k+1}) - (Az_k, z_k) + (Az_k, z_{k+1}) - (Az_{k+1}, z_k) = (Az_{k+1}, z_{k+1}) - (Az_k, z_k)$, obtenemos (17).

Supongamos cumplida la condición $B > \tau A/2$. Entonces, el primer sumando en el miembro izquierdo de la identidad (17) es no negativo y $\|z_{k+1}\|_\lambda^2 \leq \|z_k\|_\lambda^2$. De aquí se deduce que $0 \leq \|z_{k+1}\|_\lambda \leq \|z_k\|_\lambda \leq \dots \leq \|z_0\|_\lambda$, es decir, la sucesión $\{\|z_k\|_\lambda\}$ no es creciente y está acotada inferiormente por cero. Por ello, en virtud del teorema de Weierstrass, $\{\|z_k\|_\lambda\}$ converge para $k \rightarrow \infty$. Demostremos que $\lim_{k \rightarrow \infty} \|z_k\|_\lambda = 0$.

El operador $P = B - \frac{\tau}{2}A$ es positivo, y $P_0 = B_0 - \frac{\tau}{2}A = \frac{1}{2}(P + P^*)$, definido positivo, es decir, existe tal número $\delta > 0$ (véase el cap. I, § 4), que

$$(Py, y) = (P_0y, y) \geq \delta \|y\|^2 \text{ para cualquier } y \in H.$$

Por eso, de la identidad (17) obtenemos una desigualdad

$$\frac{2\delta}{\tau} \|z_{k+1} - z_k\|^2 + \|z_{k+1}\|_\lambda^2 \leq \|z_k\|_\lambda^2. \quad (*)$$

Dado que $\{\|z_k\|_\lambda\}$ es convergente, de aquí se infiere que existe

$$\lim_{k \rightarrow \infty} \|z_{k+1} - z_k\| = 0. \quad (19)$$

Luego, de la ecuación (15) encontramos

$$Az_k = -\frac{1}{\tau} B(z_{k+1} - z_k), \quad z_k = -\frac{1}{\tau} A^{-1}B(z_{k+1} - z_k),$$

$$(Az_k, z_k) = \frac{1}{\tau^2} (A^{-1}B(z_{k+1} - z_k), B(z_{k+1} - z_k)),$$

$$\|z_k\|_\lambda^2 \leq \frac{1}{\tau^2} \|A^{-1}\| \|B\|^2 \|z_{k+1} - z_k\|^2. \quad (**)$$

De aquí precisamente concluimos que $\lim_{k \rightarrow \infty} \|z_k\|_A = 0$.

OBSERVACION. De las desigualdades (*) y (**) proviene que el método de iteraciones (14) converge en las condiciones (16) con la velocidad de una progresión geométrica, $\|z_{k+1}\|_A^2 \leq \rho^2 \|z_k\|_A^2$, donde $\rho^2 = 1 - \frac{2\delta\tau}{\|A^{-1}\| \|B\|^2} < 1$.

Apliquemos el teorema 1 para demostrar la convergencia de los métodos iterativos estudiados en los pp. 2-4.

MÉTODO DE LA ITERACIÓN SIMPLE, $B = E$. Al tomar en consideración que $E \geq \frac{1}{\|A\|} A$, tenemos

$$B - \frac{\tau}{2} A = B - \frac{\tau}{2} A \geq \left(\frac{1}{\|A\|} - \frac{\tau}{2} \right) A > 0$$

para $\frac{1}{\|A\|} - \frac{\tau}{2} > 0$. El método de la iteración simple converge para todos los valores de τ , que satisfacen la desigualdad $\tau < 2/\|A\|$.

MÉTODO DE SEIDEL, $B = D + A^{-}$, $\tau = 1$. En este caso

$$\begin{aligned} B - \frac{1}{2} A &= D + A^{-} - \frac{1}{2} (A^{-} + A^{+} + D) = \frac{D}{2} + \frac{1}{2} (A^{-} - A^{+}), \\ \left((B - \frac{1}{2} A) y, y \right) &= \frac{1}{2} (Dy, y) + \frac{1}{2} ((A^{-} - A^{+})y, y) = \\ &= \frac{1}{2} (Dy, y) > 0, \end{aligned}$$

siempre que $D > 0$.

OBSERVACION. La desigualdad $D > 0$ proviene de la condición $A > 0$. En efecto, supongamos que $A > 0$ y $\xi = (\xi^1, 0, \dots, 0)$; entonces $(A\xi, \xi) = (D\xi, \xi) = a_{11} (\xi^1)^2 > 0$, es decir, $a_{11} > 0$. De un modo análogo nos convencemos de que $a_{ii} > 0$, y, por consiguiente, $D > 0$. Así pues, el método de Seidel es siempre convergente, si A es un operador autoconjugado positivo.

Para estimar la velocidad de convergencia se deben estipular las suposiciones más fuertes. Citemos el siguiente

TEOREMA 2. El método de Seidel converge con la velocidad de una progresión geométrica de razón $q < 1$, si $A = (a_{ij}) = A^* > 0$, y

$$\sum_{j \neq i}^{1+N} |a_{ij}| \leq q |a_{ii}|, \quad i = 1, 2, \dots, N, \quad q < 1. \quad (20)$$

En efecto, para el error $z_h = y_h - u$ tenemos

$$a_{ii} z_{h+1}^{(i)} = - \sum_{j < i} a_{ij} z_{h+1}^{(j)} - \sum_{j > i} a_{ij} z_h^{(j)},$$

$$|a_{ii}| |z_{h+1}^{(i)}| \leq \sum_{j < i} |a_{ij}| |z_{h+1}^{(j)}| + \sum_{j > i} |a_{ij}| |z_h^{(j)}|.$$

Supongamos que el máx $|z_{h+1}^{(i)}|$ se alcanza para cierto $i = i_0$, de modo que

$$\|z_{h+1}\|_C = |z_{h+1}^{(i_0)}|, \quad |a_{i_0 i_0}| \cdot \|z_{h+1}\|_C \leq \sum_{j < i_0} |a_{i_0 j}| \cdot \|z_{h+1}\|_C + \sum_{j > i_0} |a_{i_0 j}| \|z_h\|_C,$$

$$\|z_{h+1}\|_C \leq \left[\sum_{j > i_0} |a_{i_0 j}| / (|a_{i_0 i_0}| - \sum_{j < i_0} |a_{i_0 j}|) \right] \|z_h\|_C.$$

En virtud de la condición (20) tenemos

$$\sum_{j > i_0} |a_{i_0 j}| \leq q |a_{i_0 i_0}| - \sum_{j < i_0} |a_{i_0 j}| < q (|a_{i_0 i_0}| - \sum_{j < i_0} |a_{i_0 j}|),$$

y, por consiguiente,

$$\|z_{h+1}\|_C \leq q \|z_h\|_C \leq q^{h+1} \|z_0\|_C,$$

lo que se trataba de demostrar.

La condición (20) significa que $A = (a_{ij})$ es una matriz con preponderancia diagonal.

MÉTODO DE RELAJACIÓN SUPERIOR. $B = D + \omega A^-$, $\tau = \omega$. Hallemos la diferencia

$$\begin{aligned} B - \frac{\tau}{2} A &= D + \omega A^- - \frac{\omega}{2} (A^- + A^+ + D) = \\ &= \left(1 - \frac{\omega}{2}\right) D + \frac{\omega}{2} (A^- - A^+) \end{aligned}$$

y calculemos

$$\left(\left(B - \frac{\tau}{2} A \right) y, y \right) = \left(1 - \frac{\omega}{2} \right) (Dy, y) > 0 \text{ para } 0 < \omega < 2.$$

De este modo el método de relajación superior converge para cualesquiera valores de $\omega \in (0, 2)$, si $A = A^* > 0$.

6. Velocidad de convergencia del método implícito de iteración simple. El propio hecho de convergencia de las iteraciones no es bastante para poder juzgar sobre la aplica-

bilidad en la práctica de tal o cual método iterativo. Se necesita la información sobre la velocidad de la convergencia del método, es decir, de hecho, sobre el número de iteraciones $n = n_0(\varepsilon)$ que sean suficientes para la resolución del problema con una exactitud prefijada $\varepsilon > 0$. El número de iteraciones $n_0(\varepsilon)$ depende del parámetro τ , el que debe precisamente escogerse a partir de la condición del número mínimo de iteraciones $n = n(\varepsilon)$, con el cual se cumple la condición $\|y_n - u\|_D \leq \varepsilon \|y_0 - u\|_D$, donde D es un operador, $D = D^* > 0$.

Analizaremos aquí un esquema estacionario implícito (esquema implícito de la iteración simple)

$$B \frac{y_{k+1} - y_k}{\tau} + Ay_k = f, \quad k=0, 1, \dots, \quad \text{para cualesquiera } y_0 \in H, \quad (21)$$

donde A y B son los operadores autoconjugados positivos.

Los métodos de Seidel y de relajación superior no pertenecen a esta familia de esquemas, puesto que para ellos el operador B no es autoconjugado. Para la corrección

$$w_k = B^{-1}r_k, \quad r_k = Ay_k - f$$

se verifica (al igual que para el error $z_k = y_k - u$) una ecuación homogénea

$$B \frac{w_{k+1} - w_k}{\tau} + Aw_k = 0, \quad k=0, 1, \dots, \quad w_0 = B^{-1}(Ay_0 - f) \quad (22)$$

donde $r_k = Ay_k - f$ es un defecto, $w_k = B^{-1}r_k$ es la corrección. Efectivamente, de (21) encontramos

$$y_{k+1} = y_k - \tau B^{-1}(Ay_k - f) = y_k - \tau w_k,$$

$$Ay_{k+1} - f = Ay_k - f - \tau Aw_k, \quad r_{k+1} = r_k - \tau Aw_k.$$

Por cuanto $r_k = B(B^{-1}r_k) = Bw_k$, de aquí proviene (22).

Supondremos cumplidas las desigualdades operacionales

$$\gamma_1 B \leq A \leq \gamma_2 B, \quad \gamma_1 > 0, \quad \gamma_2 \geq \gamma_1 > 0, \quad (23)$$

o bien

$$\gamma_1 (Bx, x) \leq (Ax, x) \leq \gamma_2 (Bx, x) \quad \text{para cualesquiera } x \in H, \quad (24)$$

donde las constantes γ_1, γ_2 son conocidas,

9) **TEOREMA 3.** *Supongamos cumplidas las condiciones (23), (24). En este caso el número mínimo de iteraciones según el método (21) se alcanza para*

$$\tau = \tau_0 = \frac{2}{\gamma_1 + \gamma_2}. \quad (25)$$

Además, se verifica la desigualdad

$$\|Ay_n - f\|_{B^{-1}} \leq \rho_0^n \|Ay_0 - f\|_{B^{-1}}, \quad n = 1, 2, \dots, \quad (26)$$

$$\rho_0 = (1 - \xi)/(1 + \xi), \quad \xi = \gamma_1/\gamma_2. \quad (27)$$

DEMOSTRACIÓN. Con el fin de resolver el problema (22) hagamos uso de la siguiente estimación (la demostración de la estimación se aduce en el cap. V)

$$\|w_n\|_B \leq \rho^n \|w_0\|_B \quad \text{para} \quad \tau \leq \tau_0, \quad (28)$$

donde $\rho = 1 - \tau\gamma_1$. El valor mínimo de ρ (para el cual el número de iteraciones es mínimo) se alcanza si $\tau = \tau_0$: $\rho \geq \rho_0 = 1 - \tau_0\gamma_1 = (1 - \xi)/(1 + \xi)$. Nos queda tomar en consideración que $\|w_n\|_B = \|B^{-1}r_n\|_B = \|r_n\|_{B^{-1}}$. El teorema está demostrado.

Exigiendo que sea $\rho_0^n \leq \varepsilon$, o bien $(1/\rho_0)^n \geq 1/\varepsilon$, obtendremos la estimación para el número de iteraciones:

$$n \geq \ln(1/\varepsilon)/\ln(1/\rho_0). \quad (29)$$

OBSERVACIÓN. Una función $\varphi(\xi) = \ln(1 + \xi)/(1 - \xi) - 2\xi$ es positiva para cualesquiera $0 < \xi < 1$, puesto que $\varphi'(\xi) = 2\xi^2/(1 - \xi^2) > 0$, $\varphi(0) = 0$; por esto, $1/\ln(1/\rho_0) < 1/(2\xi)$ y la condición (29) queda cumplida, siempre que

$$n \geq n_0(\varepsilon) = (1/(2\xi)) \ln 1/\varepsilon, \quad \xi = \gamma_1/\gamma_2 \quad (30)$$

($n_0(\varepsilon)$ no es, en el caso general, entero). La condición (30) resulta más cómoda para las estimaciones. La estimación $\rho_0^n \leq \varepsilon$ es, evidentemente, verídica, si $n_0(\varepsilon) \leq n < n_0(\varepsilon) + 1$. Por esta razón, a título de n es suficiente tomar la parte entera del número $n_0(\varepsilon) + 1$.

7. Problema modelo. La comparación de los diferentes métodos iterativos se realizara a base del siguiente problema modelo

$$\frac{v_{i-1} - 2v_i + v_{i+1}}{h^2} = -\tilde{f}_i, \quad i = 1, 2, \dots, N-1,$$

$$v_0 = \mu_1, \quad v_N = \mu_2, \quad h = \frac{1}{N}, \quad (31)$$

el cual es un esquema de diferencias para el problema de contorno

$$\frac{d^2 u}{dx^2} = -\tilde{f}(x), \quad 0 < x < 1, \quad u(0) = \mu_1, \quad u(1) = \mu_2.$$

Escribamos el sistema de ecuaciones primeramente en la forma matricial:

$$Av = f, \quad (32)$$

donde $v = (v^{(1)}, v^{(2)}, \dots, v^{(N-1)})$ es un vector de dimensión $N - 1$, y A es una matriz tridiagonal de dimensión $(N - 1) \times (N - 1)$:

$$A = -\frac{1}{h^2} \begin{bmatrix} -2 & 1 & 0 & \dots & & 0 \\ & 1 & -2 & 1 & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & & & \dots & 1 & -2 & 1 \\ 0 & & & \dots & 0 & 1 & -2 \end{bmatrix}.$$

El segundo miembro de la ecuación (32) cuenta con los componentes $f_i = \tilde{f}_i$ para $i = 2, 3, \dots, N - 2$, $\tilde{f}_1 = \tilde{f}_1 + \mu_1/h^2$, $f_{N-1} = \tilde{f}_{N-1} + \mu_2/h^2$. A la matriz A le corresponde el operador A que actúa en el espacio $H = \Omega$ de funciones reticulares definidas en los nodos interiores de la red $\omega_h = \{x_i = ih, 0 < i < N\}$. Sea $\Lambda v = v_{\bar{x}}$, \bar{v} una función reticular que está definida sobre la red $\bar{\omega}_h = \{x_i = ih, 0 \leq i \leq N\}$ y que se reduce a cero en la frontera cuando $i = 0, N$. En este caso podemos escribir

$$Av = -\Lambda \bar{v}, \quad v \in \Omega = H, \quad v \in \bar{\Omega}.$$

Introduzcamos en $H = \Omega$ (como lo hacemos habitualmente) un producto escalar

$$(y, v) = \sum_{i=1}^{N-1} y_i v_i h$$

y hagamos uso de las fórmulas (17), (56) del § 4, cap. I, en virtud de las cuales

$$(Av, w) = (v, Aw), \quad \text{es decir, } A = A^*,$$

$$(Av, v) \geq \delta \|v\|^2, \quad \delta = \frac{4}{h^2} \operatorname{sen}^2 \frac{\pi h}{2}, \quad A \geq \delta E.$$

Ahora tenemos

$$\|A\| = \Delta = \frac{4}{h^2} \cos^2 \frac{\pi h}{2}.$$

Estimemos el número de iteraciones para el esquema explícito de una iteración simple en el caso del problema modelo. Se tiene $B = E$, $\delta E \leq A \leq \Delta E$, es decir,

$$\gamma_1 = \delta, \quad \gamma_2 = \Delta, \quad \xi = \frac{\gamma_1}{\gamma_2} = \operatorname{tg}^2 \frac{\pi h}{2} \approx \frac{\pi^2 h^2}{4}.$$

Para el número de iteraciones tenemos

$$n(\varepsilon) \geq n_0(\varepsilon) = \frac{\ln 1/\varepsilon}{2\xi} \approx \frac{2}{10h^2} \ln \frac{1}{\varepsilon}.$$

Prefijemos $\varepsilon = \frac{1}{2} \cdot 10^{-4} \approx \varepsilon^{-10}$, entonces $n_0(\varepsilon) \approx \frac{2}{h^2} = 2N^2$.

En particular, el número de iteraciones:

$$n_0(\varepsilon) \approx 200 \text{ para } N = 10$$

$$n_0(\varepsilon) \approx 20\,000 \text{ para } N = 100.$$

El método de iteración simple depende fuertemente del número de ecuaciones N ($n_0(\varepsilon) \approx N^2$). Abajo se exponen los métodos (véanse los §§ 4, 5), para los cuales la dependencia citada (n en función de N) será más débil ($n_0(\varepsilon) \approx N$ y $n_0(\varepsilon) \approx \sqrt{N}$).

El problema (31) es un problema tipo, puesto que una ecuación en diferencias análoga simula la ecuación de Laplace en diferencias para los casos bidimensional y tridimensional, y el número de iteraciones no depende prácticamente del número de mediciones (depende sólo de h).

8. Esquema de tres capas. Si y_{k+1} se calcula mediante dos iteraciones precedentes y_k e y_{k-1} , entonces el método iterativo se denomina *de dos pasos* (o *de tres capas*). Demos un ejemplo del esquema iterativo de tres capas. El esquema explícito de tres capas con parámetros constantes se anota, corrientemente, en la forma

$$y_{k+1} = (1 + \alpha)(E - \tau_0 A) y_k - \alpha y_{k-1} + (1 + \alpha) \tau_0 f, \\ k = 1, 2, \dots \quad (33)$$

La primera iteración se calcula según el método explícito de iteración simple:

$$y_1 = (E - \tau_0 A) y_0 + \tau_0 f \text{ para cualesquiera } y_0 \in H, \quad (34)$$

donde

$$\tau_0 = \frac{2}{\gamma_1 + \gamma_2}, \quad \alpha = \rho_1^2, \quad \rho_1 = \frac{1 - \sqrt{\xi}}{1 + \sqrt{\xi}}, \quad \xi = \frac{\gamma_1}{\gamma_2}, \quad (35)$$

$\gamma_1, \gamma_2 > 0$ son las fronteras del espectro del operador $A = A^*$: $\gamma_1 E \leq A \leq \gamma_2 E$.

Se puede mostrar que para el método (33), (34) el número de iteraciones se halla de la condición

$$q_n = \rho_1^n \left(1 + \frac{1 - \rho_1^2}{1 + \rho_1^2} n \right) \leq \varepsilon.$$

De aquí se ve que

$$n_0(\varepsilon) \approx \frac{c_0}{2\sqrt{\xi}} \ln \frac{1}{\varepsilon}, \quad 1 < c_0 < 2. \quad (36)$$

Para el problema modelo $\sqrt{\xi} \approx \pi h/2$ y

$$n_0(\varepsilon) = \frac{c_0}{\pi h} \ln \frac{1}{\varepsilon} \approx c_0 \frac{0,32}{h} \ln \frac{1}{\varepsilon} \approx c_0 \cdot 3,2N \text{ para } \varepsilon \approx e^{-10}.$$

El número de iteraciones:

$$n_0(\varepsilon) \approx 32 \div 60 \text{ para } N = 10,$$

$$n_0(\varepsilon) \approx 320 \div 620 \text{ para } N = 100,$$

es decir, considerablemente menos que para la iteración simple.

El esquema implícito de tres capas tiene por expresión

$$B y_{k+1} = (1 + \alpha)(B - \tau_0 A) y_k - \alpha B y_{k-1} + (1 + \alpha) \tau_0 f,$$

$$k = 1, 2, \dots,$$

$$B y_1 = B y_0 - \tau_0 A y_0 + \tau_0 f \text{ para cualesquiera } y_0 \in H.$$

Si $B = B^* > 0$ y se cumplen las desigualdades (23), (24) mientras que α, τ_0 se calculan según las fórmulas (35), entonces la estimación (36) para el número de iteraciones queda justa en este caso también.

§ 4. Esquema iterativo de dos capas con parámetros de Chébishev

1. Planteamiento del problema. Sea dada una ecuación

$$Au = f, \quad A: H \rightarrow H. \quad (1)$$

Veamos un esquema iterativo con parámetros variables $\{\tau_k\}$:

$$B \frac{y_{k+1} - y_k}{\tau_{k+1}} + Ay_k = f, \quad k = 0, 1, 2, \dots, \\ \text{para cualquier } y_0 \in H. \quad (2)$$

La ecuación homogénea

$$B \frac{z_{k+1} - z_k}{\tau_{k+1}} + Az_k = 0, \quad k = 0, 1, 2, \dots, \quad z_0 = y_0 - u, \quad (3)$$

es satisfecha no sólo por el error $z_k = y_k - u$, sino también por la corrección $w_k = B^{-1}(Ay_k - f)$ ($k = 0, 1, \dots$) con la condición inicial $w_0 = B^{-1}(Ay_0 - f)$. La condición para que terminen las iteraciones tiene por expresión

$$\|z_n\|_D \leq \varepsilon \|z_0\|_D, \quad \text{o bien} \quad \|w_n\|_D \leq \varepsilon \|w_0\|_D. \quad (4)$$

De (3) se ve que

$$z_{k+1} = S_{k+1}z_k, \quad S_{k+1} = E - \tau_{k+1}B^{-1}A, \quad (5)$$

donde S_{k+1} es el operador de paso de la capa k a la capa $k+1$. Eliminando z_k, z_{k-1}, \dots, z_1 , encontramos para $k = n-1$:

$$z_n = T_n z_0, \quad T_n = S_n S_{n-1} \dots S_2 S_1,$$

donde T_n es el operador de resolución del esquema (3). De aquí proviene que

$$\|z_n\|_D \leq q_n \|z_0\|_D, \quad q_n = \|T_n\|_D. \quad (6)$$

La condición para que terminen las iteraciones queda cumplida, si

$$q_n = \|T_n\|_D \leq \varepsilon. \quad (7)$$

Para estimar el número de iteraciones $n = n(\varepsilon)$ se debe obtener la desigualdad (7).

Estudiemos el esquema explícito (2)

$$\frac{y_{k+1} - y_k}{\tau_{k+1}} + Ay_k = f, \quad k=0, 1, 2, \dots, \quad (8)$$

con la particularidad de que consideramos prefijado cualquier $y_0 \in H$, y elijamos los parámetros $\tau_1, \tau_2, \dots, \tau_n$ partiendo de la condición de mín $n(\varepsilon)$. Se supone, además, que

$$A = A^* > 0, \quad \gamma_1 E \leq A \leq \gamma_2 E, \quad \gamma_1 > 0.$$

Para el residuo $r_k = Ay_k - f$ se verifica una ecuación homogénea

$$\frac{r_{k+1} - r_k}{\tau_{k+1}} + Ar_k = 0, \quad k=0, 1, 2, \dots, \quad r_0 = Ay_0 - f,$$

o bien

$$r_{k+1} = S_{k+1} r_k, \quad S_{k+1} = E - \tau_{k+1} A.$$

De aquí hallemos

$$r_n = T_n r_0, \quad T_n = S_1 S_2 \dots S_n.$$

El operador de resolución T_n es un polinomio de grado n respecto de A :

$$T_n = P_n(A) = (E - \tau_1 A)(E - \tau_2 A) \dots (E - \tau_n A)$$

con los coeficientes que sólo dependen de $\tau_1, \tau_2, \dots, \tau_n$.

Para determinar $\tau_1, \tau_2, \dots, \tau_n$ obtenemos la estimación

$$\|r_n\| \leq \|P_n(A)\| \|r_0\|.$$

Es menester hallar tales $\tau_1, \tau_2, \dots, \tau_n$, para los cuales $\|P_n(A)\|$ es mínima y, después, estimar dicha norma a través de las constantes γ_1 y γ_2 . Demos aquí sin demostración la solución de este problema.

$$\text{Designemos con } \mathfrak{M}_n = \left\{ -\cos \frac{2i-1}{2n} \pi, \quad i=1, 2, \dots, n \right\}$$

el conjunto de ceros del polinomio de Chébishev $T_n(x) = \cos(n \arccos x)$ en el segmento $-1 \leq x \leq 1$, y con $\{\mu_k\}$, una sucesión cualquiera de estos ceros, $\mu_k = \mathfrak{M}_n$. El número mínimo de iteraciones se alcanza para los valores de los

parámetros

$$\begin{aligned} \tau_k &= \frac{\tau_0}{1 + \rho_0 \mu_k}, & k &= 1, 2, \dots, n, \\ \tau_0 &= \frac{2}{\gamma_1 + \gamma_2}, & \rho_0 &= \frac{1 - \xi}{1 + \xi}, & \xi &= \frac{\gamma_1}{\gamma_2}. \end{aligned} \quad (9)$$

En este caso es válida la estimación

$$\|Ay_k - f\| \leq q_n \|Ay_0 - f\|, \quad q_n = \frac{2\rho_1^n}{1 + \rho_1^{2n}}, \quad \rho_1 = \frac{1 - \sqrt{\xi}}{1 + \sqrt{\xi}}. \quad (10)$$

El esquema (8) con parámetros iterativos (9) lleva el nombre de *esquema iterativo de Chébishev*.

El requisito $q_n \leq \varepsilon$, ó $2\rho_1^n \leq \varepsilon(1 + \rho_1^{2n})$ queda cumplido, si $\rho_1^n \leq \varepsilon/2$, o bien

$$n(\varepsilon) \geq \ln \frac{2}{\varepsilon} / \ln \frac{1}{\rho_1}. \quad (11)$$

Al observar (compárese con el § 3, p. 6) que $\ln \frac{1}{\rho_1} = \ln \frac{1 + \sqrt{\xi}}{1 - \sqrt{\xi}} > 2\sqrt{\xi}$, sustituyamos (11) por el requisito más fuerte:

$$n(\varepsilon) > n_0(\varepsilon) = \frac{1}{4\sqrt{\xi}} \ln \frac{2}{\varepsilon}, \quad (12)$$

que sea más cómodo para la comprobación. De (12) se deduce, evidentemente, (10), y $q_n \leq \varepsilon$.

Comparemos, según el número de iteraciones, el esquema (8) con la totalidad indicada de parámetros y el método de iteración simple, recurriendo con este fin al ejemplo del problema modelo citado en el § 3. En este caso $\xi \approx \pi^2 h^2/4$, $\sqrt{\xi} = \pi h/2$.

Para el método de iteración simple

$$n_0^{(1)}(\varepsilon) \approx 2/h^2 \quad \text{para} \quad \varepsilon = 10^{-4}.$$

Para el esquema de Chébishev

$$n_0^{(2)}(\varepsilon) = 3,4/h \quad \text{para} \quad \varepsilon = 10^{-4}.$$

De aquí se ve que

$$n_0^{(2)} \approx 34, \quad n_0^{(1)} \approx 200 \quad \text{para } N = 10 \quad (h = 1/10).$$

$$n_0^{(2)} \approx 340, \quad n_0^{(1)} \approx 20\,000 \quad \text{para } N = 100 \quad (h = 1/100).$$

2. Argumentación de la elección óptimal de los parámetros. Demostremos la estimación (10) en el caso de parámetros iterativos (9). Nos hace falta hallar el mín $\|P_n(A)\|$.
(τ_k)

El polinomio

$$P_n(A) = \prod_{k=1}^n (E - \tau_k A) =$$

$$= c_0 + c_1 A + \dots + c_k A^k + \dots + c_n A^n,$$

$$c_0 = 1, \quad P_n(0) = 1$$

es un operador autoconjugado. Sean ξ_s, λ_s ($s = 1, 2, \dots, N$) funciones propias y valores propios, respectivamente, del operador A :

$$A\xi_s = \lambda_s \xi_s, \quad (\xi_s, \xi_m) = \delta_{sm}, \quad s, m = 1, 2, \dots, N.$$

El operador A^k tiene las mismas funciones propias y los mismos valores propios λ_s^k :

$$A^k \xi_s = \lambda_s^k \xi_s. \quad (13)$$

Multiplicando (13) por c_k y sumando según $k = 0, 1, \dots, n$ ($c_0 = 1$), obtendremos

$$P_n(A) \xi_s = \sum_{k=0}^n c_k A^k \xi_s = \sum_{k=0}^n c_k \lambda_s^k \xi_s = P_n(\lambda_s) \xi_s.$$

Al cotejar esto con $P_n(A) \xi_s = \lambda_s (P_n(A)) \xi_s$, vemos que

$$\lambda (P_n(A)) = P_n(\lambda(A)).$$

Los valores propios del polinomio operacional $P_n(A)$ se definen como polinomios $P_n(\lambda)$ de los valores propios correspondientes del operador A , mientras que las funciones propias son las mismas que tiene el operador A . Siendo el operador $P_n(A)$ autoconjugado, su norma es igual a un valor propio cuyo módulo es máximo

$$\|P_n(A)\| = \max_{1 \leq s \leq N} P_n(\lambda_s).$$

Los valores propios λ_s del operador A se disponen en el segmento $[\gamma_1, \gamma_2]$: $\gamma_1 \leq \lambda_s \leq \gamma_2$. Es evidente que

$$\max_{1 \leq s \leq N} |P_n(\lambda_s)| \leq \max_{\gamma_1 \leq x \leq \gamma_2} |P_n(x)|,$$

donde el argumento continuo x toma todos los valores en el segmento $[\gamma_1, \gamma_2]$, y, por consiguiente, el problema del mínimo de $\|P_n(A)\|$ se reduce al de mini-máx del polinomio $P_n(x)$, es decir, al problema de determinar $\min_{(\tau_k)} \max_{\gamma \leq x_1 \leq \gamma_2} |P_n(x)|$.

Apliquemos el segmento $[\gamma_1, \gamma_2]$ sobre el segmento $[-1, 1]$, suponiendo

$$x = \frac{1}{2} [(\gamma_2 - \gamma_1)t + \gamma_2 + \gamma_1], \quad -1 \leq t \leq 1,$$

$$\text{para } \gamma_1 \leq x \leq \gamma_2. \quad (14)$$

Entonces, $P_n(t) = \tilde{P}_n(t)$. La condición de normalización $P_n(0) = 1$ toma la forma

$$\tilde{P}_n(t_0) = 1, \quad t_0 = -1/\rho_0. \quad (15)$$

Así pues, se requiere hallar un polinomio cuya desviación de cero en el segmento $-1 \leq t \leq 1$ sea mínima, de modo que el máx $|\tilde{P}_n(t)|$ sea mínimo para la condición complementaria de la normalización (15). El polinomio buscado es

$$\tilde{P}_n(t) = \frac{T_n(t)}{T_n(t_0)}, \quad (16)$$

donde $T_n(t)$ es el *polinomio de Chébishev*,

$$T_n(t) = \cos(n \arccos t) \quad \text{para } |t| \leq 1, \quad (17)$$

$$T_n(t) = \frac{1}{2} [(t + \sqrt{t^2 - 1})^n + ((t - \sqrt{t^2 - 1})^n)] \\ \text{para } |t| > 1. \quad (18)$$

El polinomio de Chébishev tiene ceros

$$t_i = \cos \frac{2i-1}{2n} \pi, \quad i = 1, 2, \dots, n. \quad (19)$$

El polinomio $P_n(x) = (1 - \tau_1 x)(1 - \tau_2 x) \dots (1 - \tau_n x)$ tiene ceros $x_i = 1/\tau_i$.

Exigiendo que las raíces de estos polinomios coincidan y teniendo presente la relación (14) entre x y t , obtenemos $2 = [(\gamma_1 + \gamma_2) + (\gamma_2 - \gamma_1) t_i] \tau_i$, de donde se infiere

$$\tau_i = 2/[\gamma_2 + \gamma_1 + (\gamma_2 - \gamma_1) t_i], \quad i = 1, 2, \dots, n. \quad (20)$$

Esta fórmula queda en vigor, cualquiera que sea el método de poner en orden los ceros del polinomio de Chébishev; por ejemplo, en lugar de (19) podemos hacer $t_i = -\cos \frac{2i-1}{2n} \pi$. Al tener esto en cuenta llegamos a la fórmula (9). Se ha de notar que si $n = 1$, obtenemos $\tau_i = \tau_0$ que es un parámetro optimal del método de iteración simple.

Así pues, los parámetros $\tau_1, \tau_2, \dots, \tau_n$ se han determinado según (9). Hallemos ahora

$$\begin{aligned} q_n &= \max_{\gamma_1 < x < \gamma_2} |P_n(x)| = \max_{-1 \leq t \leq 1} |\tilde{P}_n(t)| = \\ &= \max_{-1 \leq t \leq 1} \left| \frac{T_n(t)}{T_n(t_0)} \right| = \frac{1}{|T_n(t_0)|}, \end{aligned}$$

puesto que $\max_{-1 \leq t \leq 1} |T_n(t)| = 1$. Tenemos $|t_0| > 1$; por ello, para $|T_n(t_0)|$ se usará la fórmula (18) con $t = t_0$. Transformemos las expresiones que figuran en esta fórmula:

$$\begin{aligned} |t_0| \pm \sqrt{t_0^2 - 1} &= \frac{1}{\rho_0} \pm \sqrt{\frac{1}{\rho_0^2} - 1} = \frac{1}{\rho_0} [1 \pm \sqrt{1 - \rho_0^2}] = \\ &= \frac{1}{\rho_0} \left(1 \pm \frac{2\sqrt{\xi}}{1 - \xi} \right) = \frac{1}{\rho_0} (1 \pm \sqrt{\xi})^2 / (1 - \xi) = \\ &= (1 \pm \sqrt{\xi})^2 / (1 - \xi) = (1 \pm \sqrt{\xi}) / (1 \mp \sqrt{\xi}), \end{aligned}$$

de modo que $|t_0| + \sqrt{t_0^2 - 1} = \frac{1}{\rho_1}$, $|t_0| - \sqrt{t_0^2 - 1} = \rho_1$, y

$$|T_n(t_0)| = \frac{1}{2} \left(\frac{1}{\rho_1^n} + \rho_1^n \right) = \frac{1 + \rho_1^{2n}}{2\rho_1^n} = \frac{1}{q_n}.$$

La estimación (10) queda demostrada.

3. Estabilidad computacional y ordenación de los parámetros. El método iterativo (8) con parámetros de Chébishev $\{\tau_k\}$ se denomina a veces *método de Richardson*. Se conoce desde hace tiempo, no obstante casi no se empleaba en la práctica hasta el último tiempo por su *inestabilidad*

computacional. Expliquemos esta noción con un ejemplo. Tomemos un sistema de ecuaciones

$$u(i-1) - 2u(i) + u(i+1) = 0, \quad i = 1, 2, \dots, N-1, \\ u(0) = 1, \quad u(N) = 0. \quad (21)$$

Su solución es $u(i) = 1 - x_i$, $x_i = ih$, $h = 1/N$. Buscaremos la solución de este problema usando el método iterativo de Chébishev para $N = 20$. El valor de $n_0(\varepsilon)$ podemos calcular. Puede resultar no entero. Elegimos el número entero próximo $n \geq n_0$. Para N y ε dados tenemos $n(\varepsilon) = 64$. Al conocer

$$\gamma_1 = \frac{4}{h^2} \operatorname{sen}^2 \frac{\pi h}{2}, \quad \gamma_2 = \frac{4}{h^2} \cos^2 \frac{\pi h}{2}, \\ h = \frac{1}{N}, \quad \xi = \operatorname{tg}^2 \frac{\pi h}{2} \approx \frac{\pi^2 h^2}{4} \approx 0,006,$$

se puede calcular τ_k según la fórmula (20). A título de la aproximación inicial se toma una función

$$y_0^{(1)} = \begin{cases} 1, & i = 0. \\ 0, & i > 0, \end{cases}$$

Resulta que para el método (8), (9) no es igual en que orden se toman los ceros μ_k del polinomio de Chébishev. He aquí dos modos de numerar los ceros:

$$\alpha_1) \mu_k = \cos \frac{2k-1}{2n} \pi, \quad k = 1, 2, \dots, n, \\ t_1 = \cos \frac{\pi}{2n}, \quad t_n = -\cos \frac{\pi}{2n}, \\ \alpha_2) \mu_k = -\cos \frac{2k-1}{2n} \pi.$$

Los resultados de los cálculos se aducen en la tabla 1.

Para los menores valores de N y n puede resultar que el aumento de los valores intermedios de y_k no lleva al parem, sin embargo tiene lugar la acumulación de los errores de redondeo, y tras n iteraciones no se cumplen las condiciones en que terminan las iteraciones ($\|Ay_k - f\| \leq \varepsilon \|Ay_0 - f\|$).

Estas dos peculiaridades del proceso computacional, a saber, el aumento de los valores intermedios que conduce al parem y la acumulación de los errores de redondeo, se carac-

TABLA 1

Surtido α_1		Surtido α_2	
k	$\Delta_k = \max_{x_i} v_k(x_i) - v_{k-1}(x_i) $	k	$\Delta_k = \max_{x_i} v_k(x_i) - v_{k-1}(x_i)$
53	0,12	1	39,6
55	27	2	$2,6 \cdot 10^3$
57	$1,9 \cdot 10^4$	4	$8,2 \cdot 10^8$
59	$3,7 \cdot 10^7$	7	$3,3 \cdot 10^{11}$
60	$2,6 \cdot 10^9$	9	$1,2 \cdot 10^{14}$
61	$2,5 \cdot 10^{11}$	11	$1,9 \cdot 10^{16}$
62	$3,3 \cdot 10^{13}$	12	Parem
63	$5 \cdot 10^{15}$		
64	Parem		

terizan por un término: *inestabilidad computacional*. La causa de inestabilidad computacional del método de Chébishev radica en que las normas $\|S_{k+1}\|$ del operador de paso $S_{k+1} = E - \tau_{k+1}A$ son para ciertas iteraciones superiores a uno, mientras que el proceso computacional es real, es decir, se tienen acotaciones por debajo y por arriba para los números admisibles (se tienen cero de máquina e infinidad de máquina) y en cada etapa de los cálculos surgen los errores de redondeo.

Calculemos la norma para $S_k = E - \tau_k A$. Por cuanto $S_k^* = S_k$, entonces $\|S_{k+1}\| = \sup_{\|x\|=1} |(S_{k+1}x, x)|$. De la condición $\gamma_1 E \leq A \leq \gamma_2 E$ proviene $(\tau_{k+1}\gamma_1 - 1)E \leq \tau_{k+1}A - E \leq (\tau_{k+1}\gamma_2 - 1)E$. Sustituyendo aquí la expresión para τ_{k+1} y teniendo presente que $1 - \tau_0\gamma_1 = \tau_0\gamma_2 - 1 = \rho_0$, obtenemos

$$-\frac{\rho_0(1-\mu_k)}{1+\rho_0\mu_k} E \leq \tau_{k+1}A - E \leq \frac{\rho_0(1+\mu_k)}{1+\rho_0\mu_k} E.$$

De aquí encontramos

$$\|S_{k+1}\| = \|\tau_{k+1}A - E\| = \begin{cases} \frac{\rho_0(1+\mu_k)}{1+\rho_0\mu_k} & \text{para } \mu_k > 0, \\ \frac{\rho_0(1-\mu_k)}{1+\rho_0\mu_k} & \text{para } \mu_k < 0 \end{cases}$$

de suerte que $\|S_{k+1}\| < 1$ para todos los $\mu_k > 0$, y $\|S_{k+1}\| > 1$ para $\mu_k < -(1 - \rho_0)/(2\rho_0)$. Por cuanto

$$-\cos \frac{\pi}{2n} \leq \mu_k \leq -\cos \frac{(2n-1)}{2n} \pi = \cos \frac{\pi}{2n},$$

$$k = 1, 2, \dots, n,$$

entonces, para la mayor cantidad de números k tenemos $\|S_k\| > 1$, y si se emplean muchos parámetros τ_k seguidos, para los cuales $\|S_k\| > 1$, tienen lugar la acumulación del error de redondeo y el crecimiento de las aproximaciones iterativas, lo que conduce a la inestabilidad computacional.

Con el fin de debilitar el efecto mencionado, es natural tratar de alternar los parámetros τ_k , para los cuales $\|S_k\| > 1$ con aquellos, para los cuales $\|S_k\| < 1$. En este procedimiento se realiza precisamente la construcción de una sucesión de parámetros $\{\tau_k\}$, para la cual la convergencia de las iteraciones tiene un carácter monótono y la inestabilidad computacional está ausente. Existe una regla para tal ordenación de los ceros $t_i = -\cos \frac{2i-1}{2n} \cos \pi$ del polinomio de Chébishev y, consecuentemente, también de los parámetros $\{\tau_k\}$ (para n cualquiera), para la que tiene lugar la estabilidad computacional.

Demos a conocer dicha regla para el caso en que n es una potencia del número 2, $n = 2^p$, $p > 0$ es un número entero¹⁾. Designemos el conjunto de ceros t_i , ordenado según esta regla, mediante

$$\mathfrak{M}_n^* = \left\{ -\cos \beta_i, \beta_i = \frac{\pi}{2n} \theta_i^{(n)}, i = 1, 2, \dots, n \right\}, n = 2^p,$$

donde $\theta_i^{(n)}$ es uno de los números impares $1, 3, 5, \dots, 2n - 1$. El problema se reduce, pues, a la ordenación del conjunto de n números impares: $\theta_n = \{\theta_1^{(n)}, \theta_2^{(n)}, \dots, \theta_n^{(n)}\}$. Partiendo del conjunto $\theta_1 = \{1\}$, construyamos un conjunto $\theta_{(n)}^* = \theta_{2^p}^*$ según las fórmulas

$$\theta_{2i-1}^{(2m)} = \theta_i^{(m)},$$

$$\theta_{2i}^{(2m)} = 4m - \theta_{2i-1}^{(2m)}, \quad i = 1, 2, \dots, m; \quad m = 1, 2, \dots, 2^{p-1},$$

¹⁾ La regla de ordenación de $\{\tau_k\}$ para n cualquiera se da, por ejemplo, en [6, 9].

si se conocen $\theta_1^{(m)}$. La sucesión correspondiente de parámetros $\{\tau_k^*\}$ se denominará *surtido estable*. Sea, por ejemplo, $n = 16 = 2^4$. Encontramos sucesivamente $\theta_1 = \{1\}$, $\theta_2 = \{1, 3\}$, $\theta_4 = \{1, 7, 3, 5\}$, $\theta_8 = \{1, 15, 7, 9, 3, 13, 5, 11\}$, $\theta_{16} = \{1, 31, 15, 17, 7, 25, 9, 23, 3, 29, 13, 19, 5, 27, 11, 21\}$. Al pasar de θ_m a θ_{2m} es suficiente poner tras cada $\theta_{2^{m-1}}$ un número igual a $4m - \theta_{2^{m-1}}$ (la numeración corresponde a θ_m). La sucesión «estable» θ_n^* no depende del problema. La convergencia de las iteraciones para este surtido de parámetros $\{\tau_k^*\}$ lleva un carácter no monótono, pero las oscilaciones aquí no son de gran amplitud y se amortiguan al fin y al cabo.

He aquí los resultados de cálculos para el problema (21) según el esquema (8), (9) con el surtido estable de parámetros $\{\tau_k^*\}$:

k	1	4	8	16	24	32	48	50	62
Δk	39,6	4,7	1,1	0,2	0,1	0,04	$1,5 \cdot 10^{-3}$	$6,7 \cdot 10^{-3}$	$8,7 \cdot 10^{-3}$

4. Esquemas implícitos. El método de Seidel y el de relajación superior convergen más rápidamente que el método explícito de iteración simple, razón por la cual se justifica el paso a los esquemas implícitos. ¿Cómo se debe elegir el operador B ? Es fundamental el requisito general del mínimo de operaciones $Q(\varepsilon)$ necesarias para hallar solución con una exactitud $\varepsilon > 0$; dicho requisito se reduce a dos exigencias: 1) del número mínimo de iteraciones, el cual depende tanto de B como de la elección de $\{\tau_k\}$; 2) del número mínimo de operaciones para la resolución de la ecuación

$$By_{k+1} = F_k$$

(carácter económico del operador B). De ejemplo puede servir un operador triangular correspondiente a la matriz triangular.

Mostremos ahora que los resultados obtenidos más arriba para un esquema explícito pueden extenderse a un esquema implícito. Veamos un esquema implícito

$$B \frac{y_{k+1} - y_k}{\tau_{k+1}} + Ay_k = f, \quad k = 0, 1, \dots, \text{ para todo } y_0 \in H, \quad (22)$$

donde $A = A^* > 0$, $B = B^* > 0$, y

$$\gamma_1 B \leq A \leq \gamma_2 B, \quad \gamma_1 > 0. \quad (23)$$

Eligiendo los parámetros de iteración $\{\tau_k^*\}$ según las fórmulas (9) y ordenándolos en concordancia con el punto precedente, obtendremos para la solución del problema (22) una estimación

$$\|Ay_n - f\|_{B^{-1}} \leq q_n \|Ay_0 - f\|_{B^{-1}}, \quad q_n = \frac{2\rho_1^n}{1 + \rho_1^{2n}},$$

$$\rho_1 = \frac{1 - \sqrt{\xi}}{1 + \sqrt{\xi}}, \quad \xi = \frac{\gamma_1}{\gamma_2}, \quad (24)$$

donde γ_1 y γ_2 son los números que figuran en (23). Para el número de iteraciones $n = n(\epsilon)$ son válidas las estimaciones (11) y (12). Para convencernos de esto, basta reducir el problema (22) a un problema equivalente para el esquema explícito

$$\frac{x_{k+1} - x_k}{\tau_{k+1}} + Cx_k = 0, \quad k = 0, 1, \dots, x_0 = B^{1/2}w_0, \quad (25)$$

donde $x_k = B^{1/2}w_k$, $C = B^{-1/2}AB^{-1/2}$ es un operador positivo autoconjugado con las fronteras del espectro γ_1 y γ_2 :

$$\gamma_1 E \leq C \leq \gamma_2 E. \quad (26)$$

En efecto, por cuanto $B = B^* > 0$, existe, pues, $B^{1/2} = (B^{1/2})^* > 0$. Aplicando el operador $B^{-1/2}$ a la ecuación (22), obtenemos (25) para $x_k = B^{1/2}w_k$. El modo inverso de razonamientos es evidente. Queda por demostrar la equivalencia de las desigualdades (23) y (26). Estudiemos una funcional

$$J = ((A - \gamma B)y, y) = (Ay, y) - \gamma (By, y) =$$

$$= (AB^{-1/2}(B^{1/2}y), B^{-1/2}(B^{1/2}y)) - \gamma (B^{1/2}y, B^{1/2}y) =$$

$$= (Cx, x) - \gamma (x, x) = ((C - \gamma E)x, x),$$

donde $x = B^{1/2}y$. Como y (y, por tanto, también x) es un vector arbitrario de H , entonces de la igualdad

$$J = ((A - \gamma B)y, y) = ((C - \gamma E)x, x) \quad (27)$$

se deduce que los operadores $A - \gamma B$ y $C - \gamma E$ son de signos iguales. Si, por ejemplo, $A - \gamma_1 B \geq 0$, entonces para $\gamma = \gamma_1$ la igualdad (27) nos da $C - \gamma_1 E \geq 0$, etc.

Para el esquema explícito tenemos una estimación $\|x_n\| \leq q_n \|x_0\|$. Al sustituir aquí $x_h = B^{1/2}w_h = B^{-1/2}r_h$, $r_h = Ay_h - f$, obtenemos la estimación (24).

Para los métodos de Seidel y de relajación superior $B \neq B^*$, por lo cual la totalidad de parámetros de Chébishev no puede ser empleado.

§ 5. Método alternado triangular

1. Método alternado triangular. Analizaremos un esquema iterativo implícito

$$B \frac{y_{k+1} - y_k}{\tau_{k+1}} + Ay_k = f, \quad k=0, 1, \dots \quad (1)$$

Si el operador B representa un producto del número finito de operadores económicos, será también económico. Así, es económico el operador $B = B_1 B_2$ que es igual al producto de los operadores triangulares B_1 y B_2 .

Veamos el así llamado método *alternado triangular* (1), para el cual el operador B tiene por expresión

$$B = (D + \omega R_1) D^{-1} (D + \omega R_2), \quad (2)$$

donde $D = D^* > 0$, $R_1^* = R_2$, $R_1 + R_2 = R$, $R = R^* > 0$, $\omega > 0$ es el parámetro.

Probemos que el operador B es positivo y autoconjugado, es decir, que el esquema (1) con el operador (2) pertenece a la familia inicial de esquemas (2) del § 3, razón por la cual pueden aprovecharse todos los resultados de la teoría general obtenidos anteriormente. En efecto,

$$\begin{aligned} (By, v) &= ((D + \omega R_1) D^{-1} (D + \omega R_2) y, v) = \\ &= ((D + \omega R_2) y, D^{-1} (D + \omega R_1) v) = \\ &= (y, (D + \omega R_1) D^{-1} (D + \omega R_2) v), \end{aligned}$$

y, por consiguiente, $(By, v) = (y, Bv)$, es decir, $B = B^*$. Luego, encontramos $(By, y) = ((D + \omega R_2) y, D^{-1} \times \times D^{-1} (D + \omega R_1) y) = \|(D + \omega R_2) y\|_{D^{-1}}^2 > 0$, es decir, $B = B^* > 0$.

Al operador R le corresponde una matriz $R = (r_{ij})$. A título de las matrices R_1 y R_2 pueden intervenir las matri-

ces triangulares inferior y superior, es decir,

$$R_1 = (r_{ij}), \quad r_{ij} = \begin{cases} r_{ii}/2, & j = i, \\ r_{ij}, & j < i, \\ 0, & j > i; \end{cases}$$

$$R_2 = (r'_{ij}), \quad r'_{ij} = \begin{cases} r_{ii}/2, & j = i, \\ r_{ij}, & j > i, \\ 0, & j < i. \end{cases}$$

Si R es una matriz simétrica, $r_{ji} = r_{ij}$, entonces R_1 y R_2 son recíprocamente conjugados, $R_2 = R_1^*$.

A título de $D = (d_{ij})$ tomemos una matriz diagonal. Entonces, $D + \omega R_1$ es una matriz triangular inferior y $D + \omega R_2$, una matriz triangular superior. De este modo, el proceso de la iteración se reduce a la inversión alternada de las matrices triangulares inferior y superior (de aquí proviene la denominación del método). Efectivamente, para cada iteración se debe resolver una ecuación

$$B y_{k+1} = (D + \omega R_1) D^{-1} (D + \omega R_2) y_{k+1} = F_k. \quad (3)$$

Al denotar $D^{-1} (D + \omega R_2) y_{k+1} = \bar{y}_{k+1}$, obtenemos

$$(D + \omega R_1) \bar{y}_{k+1} = F_k, \quad (D + \omega R_2) y_{k+1} = D \bar{y}_{k+1},$$

$$k = 0, 1, \dots \quad (4)$$

Observando que $(R_1 y, y) = (R_2 y, y) = (R y, y)/2$, encontramos

$$((D + \omega R_1) y, y) = (D y, y) + \omega (R_1 y, y) =$$

$$= \left(\left(D + \frac{\omega}{2} R \right) y, y \right) = ((D + \omega R_2) y, y) > 0,$$

puesto que $D > 0$, $\omega > 0$, y $R > 0$.

De aquí proviene la existencia de los operadores inversos $(D + \omega R_1)^{-1}$, $(D + \omega R_2)^{-1}$, es decir, la resolubilidad de los problemas (4).

2. Elección del parámetro ω . Para poder emplear la teoría general, se deben hallar primeramente los parámetros γ_1 y γ_2 que figuran en las desigualdades

$$\gamma_1 B \leq A \leq \gamma_2 B, \quad (5)$$

las cuales siempre se verifican gracias a que los operadores A y B son acotados y positivos. Empecemos con la determinación del parámetro $\omega > 0$.

LEMA. Supongamos que el operador B se determina según la fórmula (2), donde

$$R_2^* = R_1, \quad R_1 + R_2 = R, \quad R = R^* > 0$$

y que R satisface las condiciones

$$R \geq \delta D, \quad \delta > 0, \quad R_1 D^{-1} R_2 \leq \frac{\Delta}{4} R, \quad \Delta > 0, \quad (6)$$

En este caso será válida la estimación

$$\dot{\gamma}_1 B \leq R \leq \dot{\gamma}_2 B, \quad \dot{\gamma}_1 = \frac{\delta}{1 + \omega\delta + 0,25\omega^2\delta\Delta}, \quad \dot{\gamma}_2 = \frac{1}{2\omega}. \quad (7)$$

La razón $\xi = \dot{\gamma}_1(\omega)/\dot{\gamma}_2(\omega)$ tiene el valor máximo cuando

$$\omega = \dot{\omega} = 2/\sqrt{\delta\Delta};$$

con la particularidad de que

$$\xi = \frac{2\sqrt{\eta}}{1+\sqrt{\eta}}, \quad \eta = \frac{\delta}{\Delta}, \quad \dot{\gamma}_1 = \frac{\delta}{2(1+\sqrt{\eta})}, \quad \dot{\gamma}_2 = \frac{\delta}{4\sqrt{\eta}}. \quad (9)$$

DEMOSTRACION. Las desigualdades (6) significan que

$$(Ry, y) \geq \delta (Dy, y), \quad (D^{-1}R_2y, R_2y) \leq \frac{\Delta}{4} (Ry, y)$$

para cualesquiera $y \in H$.

Realizadas las transformaciones

$$\begin{aligned} B &= (D + \omega R_1) D^{-1} (D + \omega R_2) = \\ &= D - \omega (R_1 + R_2) + \omega^2 R_1 D^{-1} R_2 + 2\omega (R_1 + R_2) = \\ &= (D - \omega R_1) D^{-1} (D - \omega R_2) + 2\omega R, \end{aligned}$$

obtenemos

$$\begin{aligned} (By, y) &= (D^{-1} (D - \omega R_2) y, (D - \omega R_2) y) + 2\omega (Ry, y) \\ &= \| (D - \omega R_2) y \|_{B^{-1}}^2 + 2\omega (Ry, y) \geq 2\omega (Ry, y), \end{aligned}$$

de suerte que

$$B \geq 2\omega R, \quad \text{o bien } R \leq \frac{1}{2\omega} B, \quad \dot{\gamma}_2 = \frac{1}{2\omega}.$$

Obtendremos ahora para B la estimación por arriba. Tomando en consideración (6), hallemos

$$B = D + \omega R + \omega^2 R_1 D^{-1} R_2 \frac{1}{\delta} R + \omega R + \frac{\omega^3 \Delta}{4} R \leq \leq \frac{1}{\delta} \left(1 + \omega \delta + \frac{\omega^3 \delta \Delta}{4} \right) R,$$

$$R_1 \geq \dot{\gamma}_1 B, \quad \dot{\gamma}_1 = \delta \left(1 + \omega \delta + \frac{\omega^3 \delta \Delta}{4} \right)^{-1}.$$

El número de iteraciones necesarias para la resolución de la ecuación $Ry = f$ depende de la razón

$$\xi(\omega) = \dot{\gamma}_1 / \dot{\gamma}_2 = 2\omega\delta (1 + \omega\delta + \omega^3\delta\Delta/4)^{-1}.$$

Elijamos ω de la condición de que $\xi(\omega)$ sea máxima. Igualando a cero la derivada $\xi'(\omega) = 2\delta(1 - \omega^3\delta\Delta/4)(1 + \omega\delta + \omega^3\delta\Delta/4)^{-2}$, encontramos $\omega = \dot{\omega} = 2/\sqrt{\delta\Delta}$; en este caso $\xi''(\dot{\omega}) < 0$. Al sustituir este valor de ω en las fórmulas para $\dot{\gamma}_1, \dot{\gamma}_2, \xi(\omega)$, obtenemos la fórmula (9). El lema queda demostrado.

3. Velocidad de convergencia.

TEOREMA. Supongamos que el operador $A = A^* > 0$ se representa como una suma $A = A_1 + A_2$, $A_2 = A_1^*$, y que se cumplen las condiciones

$$A \geq \delta D, \quad A_1 D^{-1} A_2 \leq \frac{\Delta}{4} A, \quad \delta > 0, \quad \Delta > 0. \quad (10)$$

Entonces para el método alternado triangular (1) con

$$B = (D + \omega A_1) D^{-1} (D + \omega A_2), \quad D = D^* > 0, \quad (11)$$

con el parámetro $\omega = 2/\sqrt{\delta\Delta}$ y una totalidad de parámetros de Chébishev

$$\tau_k^* = \frac{\tau_0}{1 + \rho_0 \mu_k^*}, \quad \tau_0 = \frac{2}{\gamma_1 + \gamma_2}, \quad \rho_0 = \frac{1 - \xi}{1 + \xi},$$

$$\xi = \frac{\gamma_1}{\gamma_2} = \frac{2\sqrt{\eta}}{1 + \sqrt{\eta}} \quad (12)$$

donde

$$\gamma_1 = \frac{\delta}{2(1 + \sqrt{\eta})}, \quad \gamma_2 = \frac{\delta}{4\sqrt{\eta}}, \quad \eta = \frac{\delta}{\Delta}, \quad \mu_k^* \in \overline{\mathfrak{M}}_n^*, \quad (13)$$

son suficientes $n(\varepsilon)$ iteraciones:

$$n_0(\varepsilon) \leq n(\varepsilon) < n_0(\varepsilon) + 1, \quad n_0(\varepsilon) < \ln \frac{2}{\varepsilon} / (2\sqrt{2}\sqrt{\eta}), \quad (14)$$

y en este caso se cumple la siguiente estimación

$$\|Ay_n - f\|_{B^{-1}} \leq \varepsilon \|Ay_0 - f\|_{B^{-1}}. \quad (15)$$

DEMOSTRACION. Hagamos uso del lema precedente suponiendo $R = A$, $R_1 = A_1$, $R_2 = A_2$, y también de la estimación (24) del § 4:

$$\|Ay_n - f\|_{B^{-1}} \leq q_n \|Ay_0 - f\|_{B^{-1}}$$

con

$$q_n = \frac{2\rho_1^n}{1 + \rho_1^{2n}}, \quad \rho_1 = \frac{1 - \sqrt{\xi}}{1 + \sqrt{\xi}}.$$

En el § 3 se ha obtenido la estimación para el número de iteraciones $n = n(\varepsilon)$:

$$n_0(\varepsilon) \leq n(\varepsilon) < n_0(\varepsilon) + 1, \quad \text{donde } n_0(\varepsilon) < \ln \frac{2}{\varepsilon} / (2\sqrt{\xi}).$$

Al sustituir aquí $\xi = 2\sqrt{\eta}/(1 + \sqrt{\eta})$, obtendremos (15).

4. Ejemplo de aplicación del método alternado triangular.
Analicemos un problema modelo

$$u_{xx, i} = \frac{u_{i-1} - 2u_i + u_{i+1}}{h^2} = -f_i, \quad i = 1, 2, \dots, N-1, \\ u_0 = 0, \quad u_N = 0. \quad (16)$$

Sea $H = \Omega$ un espacio de funciones reticulares definidas en los nodos interiores $i = 1, 2, \dots, N-1$ de la red ω_h ; introduzcamos un producto escalar

$$(y, v) = \sum_{i=1}^{N-1} y_i v_i h.$$

El operador $Ay = -\bar{y}_{xx}$ es autoconjugado y definido positivo:

$$A \geq \delta E, \quad \delta = \frac{4}{h^2} \operatorname{sen}^2 \frac{\pi h}{2}.$$

Introduzcamos los operadores $Dy = y$ ($D = E$) y

$$A_1 y = R_1 y = \frac{y_{x,i}}{h} = \frac{y_i - y_{i-1}}{h^2},$$

$$A_2 y = R_2 y = -\frac{y_{x,i}}{h} = -\frac{y_{i+1} - y_i}{h^2}, \quad A_1 + A_2 = A.$$

Las iteraciones $(y_i)_k = y_k(i)$ se hallan según las fórmulas

$$(E + \omega A_1)(\bar{y}_i)_{k+1} = \left(\bar{y}_i + \omega \frac{\bar{y}_i - \bar{y}_{i-1}}{h^2} \right)_{k+1} = F_k(i),$$

$$\bar{y}_{k+1}(i) = \frac{\omega \bar{y}_{k+1}(i-1) + h^2 F_k(i)}{h^2 + \omega},$$

$$(E + \omega A_2)y_{k+1}(i) =$$

$$= y_{k+1}(i) - \frac{\omega}{h^2} (y_{k+1}(i+1) - y_{k+1}(i)) = \bar{y}_{k+1}(i),$$

en definitiva tenemos

$$y_{k+1}(i) = \frac{\omega y_{k+1}(i+1) + h^2 \bar{y}_{k+1}(i)}{\omega + h^2},$$

$$i = N-1, N-2, \dots, 2, 1.$$

Los valores de $\bar{y}_{k+1}(i)$ se hallan sucesivamente al mover de izquierda a derecha (de $i-1$ a i) y los de $y_{k+1}(i)$, de derecha a izquierda (de $i+1$ a i); y en este caso se toman en consideración las condiciones de contorno

$$\bar{y}_{k+1}(0) = 0, \quad y_{k+1}(N) = 0.$$

Las fórmulas del tipo semejante se denominan *fórmulas de cómputo móvil*.

De las igualdades $y_{\bar{x},i+1} = y_{x,i}$ se desprende que $A_1^* = A_2$. En efecto, por cuanto $v_1 = v_0 + hv_{\bar{x},1} = hv_{x,1}$,

entonces

$$\begin{aligned} (A_2 y, v) &= - \sum_{i=1}^{N-1} y_{x, i} v_i = - y_1 v_1 \frac{1}{h} - \sum_{i=1}^{N-1} y_{x+i} x_{x, i} = \\ &= y_1 v_{\bar{x}, 1} + \sum_{i=2}^N y_i v_{\bar{x}, i} = \sum_{i=1}^{N-1} y_i v_{\bar{x}, i} = h \sum_{i=1}^{N-1} y_i \frac{v_{\bar{x}, i}}{h} = (y, A_1 v), \end{aligned}$$

es decir, $A_1 = A_2^*$.

Calculemos la constante Δ :

$$\begin{aligned} (A_1 A_2 y, y) &= (A_2 y, A_2 y) = \\ &= \frac{1}{h^2} \sum_{i=1}^{N-1} (y_{x, i})^2 h = \frac{1}{h^2} \sum_{i=2}^{N-1} (y_{\bar{x}, i})^2 h \leq \\ &\leq \frac{1}{h^2} \sum_{i=1}^N h (y_{\bar{x}, i})^2 = \frac{1}{h^2} \sum_{i=1}^{N-1} h (A y)_i y_i = \frac{1}{h^2} (A y, y), \end{aligned}$$

de donde se infiere $\Delta = 4/h^2$. Así que,

$$\eta = \frac{\delta}{\Delta} = \operatorname{sen}^2 \frac{\pi h}{2} \approx \frac{\pi^2 h^2}{4}, \quad \sqrt{\eta} \approx \frac{\pi h}{2},$$

$$\xi = 2\sqrt{\eta}/(1 + \sqrt{\eta}) \approx 2\sqrt{\eta} \approx \pi h, \quad \sqrt{\xi} \approx \sqrt{\pi h},$$

de suerte que

$$n_0(\varepsilon) \approx \frac{1}{2\sqrt{\pi h}} \ln \frac{2}{\varepsilon}.$$

Si $\varepsilon = 10^{-4}$, entonces $n_0(\varepsilon) \approx 3/\sqrt{h}$.

El resultado es:

$$n_0(\varepsilon) \approx 10 \text{ para } h = 1/10 \text{ (} N = 10 \text{),}$$

$$n_0(\varepsilon) \approx 30 \text{ para } h = 1/100 \text{ (} N = 100 \text{).}$$

Recordemos que para $N = 100$ se deben realizar 20 000 iteraciones por el método de iteración simple y 340 iteraciones, por el esquema explícito de Chébishev. De este modo, el método alternado triangular ha resultado mejor entre los métodos estudiados.

§ 6. Métodos iterativos de tipo variacional

1. Método de los residuos mínimos. Al estudiar los métodos iterativos siempre se suponía hasta ahora que las constantes γ_1 y γ_2 , es decir, las fronteras del espectro del operador A en H o en H_B están conocidas. Pero, ¿qué se debe hacer, si tal información está ausente? En este caso pueden emplearse los métodos que no utilizan los parámetros γ_1 y γ_2 en la forma explícita. Estos son los métodos de tipo variacional. Aquí se analizarán los métodos de los residuos mínimos, del descenso más rápido y de los gradientes conjugados.

Empecemos con el método de los residuos mínimos para un esquema explícito

$$\frac{y_{k+1} - y_k}{\tau_{k+1}} + Ay_k = f, \quad k=0, 1, \dots, \quad \text{Para todo } y_0 \in H. \quad (1)$$

Para el residuo $r_k = Ay_k - f$ tenemos una ecuación homogénea

$$\frac{r_{k+1} - r_k}{\tau_{k+1}} + Ar_k = 0, \quad k=0, 1, \dots, \quad r_0 = Ay_0 - f. \quad (2)$$

El parámetro τ_{k+1} se escogerá partiendo de la condición de que sea mínimo el residuo r_{k+1} según la norma:

$$\begin{aligned} \|r_{k+1}\|^2 &= \|r_k - \tau_{k+1} Ar_k\|^2 = \\ &= \|r_k\|^2 - 2\tau_{k+1} (r_k, Ar_k) + \tau_{k+1}^2 \|Ar_k\|^2 = \varphi(\tau_{k+1}) \end{aligned}$$

Diferenciemos esta expresión respecto de τ_{k+1} , igualemos a cero la derivada $\varphi'(\tau_{k+1})$:

$$\varphi'(\tau_{k+1}) = -2(r_k, Ar_k) + 2\tau_{k+1} \|Ar_k\|^2 = 0$$

y hallemos

$$\tau_{k+1} = \frac{(Ar_k, r_k)}{\|Ar_k\|^2}, \quad k=1, 2, \dots \quad (3)$$

Con este valor de τ_{k+1} la segunda derivada $\varphi''(\tau_{k+1})$ es positiva y, por consiguiente, se alcanza el mín $\|r_{k+1}\|^2$.

Hasta ahora no se suponía que A es un operador autoconjugado. En cambio, si $A = A^* > 0$, entonces son válidas

las estimaciones

$$\|r_{k+1}\| \leq \rho_0 \|r_k\|, \quad \|Ay_n - f\| \leq \rho_0^n \|Ay_0 - f\|,$$

$$\rho_0 = \frac{1-\xi}{1+\xi}, \quad \xi = \frac{\gamma_1}{\gamma_2}, \quad (4)$$

donde γ_1 y γ_2 son las fronteras exactas del espectro del operador A . En efecto, por cuanto, de acuerdo con (3), la norma $\|r_{k+1}\|$ es mínima para τ_{k+1} , entonces para cada $\tau \neq \tau_{k+1}$ ella debe crecer, por lo cual

$$\|r_{k+1}\|^2 = \|r_k - \tau_{k+1}Ar_k\|^2 \leq \|r_k - \tau_0Ar_k\|^2 \leq \|E - \tau_0A\|^2 \|r_k\|^2.$$

Por otra parte se conoce que

$$\|E - \tau_0A\| = \rho_0 \text{ para } \tau_0 = 2/(\gamma_1 + \gamma_2).$$

De aquí precisamente se deduce que $\|r_{k+1}\| \leq \rho_0 \|r_k\|$.

De este modo, el método de los residuos mínimos converge con la misma velocidad que el método de iteración simple (siempre que en este último se empleen los valores exactos de γ_1 y γ_2).

En el caso del método de residuos implícito o del método de correcciones en vez de (1) obtenemos una ecuación para la corrección:

$$B \frac{w_{k+1} - w_k}{\tau_{k+1}} + Aw_k = 0, \quad k = 0, 1, \dots,$$

$$w_k = B^{-1}r_k, \quad w_0 = B^{-1}(Ay_0 - f), \quad (5)$$

donde τ_{k+1} se determina por la fórmula

$$\tau_{k+1} = \frac{(Aw_k, w_k)}{(B^{-1}Aw_k, Aw_k)}, \quad k = 0, 1, \dots \quad (6)$$

En lugar de (4) obtenemos la estimación

$$\|Ay_n - f\|_{B^{-1}} \leq \rho_0^n \|Ay_0 - f\|_{B^{-1}}.$$

2. Método del descenso más rápido. El método explícito del descenso más rápido se diferencia del método de los residuos mínimos sólo en la fórmula para τ_{k+1} :

$$\tau_{k+1} = \frac{(r_k, r_k)}{(Ar_k, r_k)}, \quad k = 0, 1, \dots \quad (7)$$

Esta fórmula se obtiene o bien de la condición del mínimo de la norma $\|z_{k+1}\|_A$ del error $z_{k+1} = y_{k+1} - u$, o bien de la condición de ortogonalidad de los residuos r_k y r_{k+1} . Al multiplicar escalarmente la ecuación $r_{k+1} = r_k - \tau_{k+1}Ar_k$ por r_k , obtenemos $0 = (r_k, r_k) - \tau_{k+1}(Ar_k, r_k)$ de donde se infiere la fórmula (7). Por cuanto $Az_k = Ay_k - Au = r_k$, entonces se tiene

$$\begin{aligned} (Az_{k+1}, z_{k+1}) &= (Az_k - \tau_{k+1}A^2z_k, z_k - \tau_{k+1}Az_k) = \\ &= (r_k - \tau_{k+1}Ar_k, z_k - \tau_{k+1}r_k) = \\ &= (r_k, z_k) - 2\tau_{k+1}(r_k, r_k) + \tau_{k+1}^2(Ar_k, r_k). \end{aligned}$$

Diferenciando $\|z_{k+1}\|_A^2$ respecto de τ_{k+1} e igualando a cero la derivada, obtendremos (7).

Luego, tenemos

$$\begin{aligned} \|z_{k+1}\|_A^2 &= \|(E - \tau_{k+1}A)z_k\|_A^2 \leq \|(E - \tau_0A)z_k\|_A^2 \leq \\ &\leq \|E - \tau_0A\|^2 \|z_k\|_A^2 \leq \rho_0^2 \|z_k\|_A, \end{aligned}$$

es decir,

$$\|z_{k+1}\|_A = \|y_{k+1} - u\|_A \leq \rho_0^n \|y_0 - u\|_A.$$

El método del descenso más rápido converge en H_A con la misma velocidad que el método de iteración simple.

3. Método de los gradientes conjugados. Los métodos de tipo variacional con mayor velocidad de convergencia pueden encontrarse en la clase de esquemas iterativos implícitos de tres capas:

$$\begin{aligned} By_{k+1} &= \alpha_{k+1}(B - \tau_{k+1}A)y_k + (1 - \alpha_{k+1})By_{k-1} + \\ &\quad + \alpha_{k+1}\tau_{k+1}f, \quad k = 1, 2, \dots, \quad (8) \\ By_k &= (B - \tau_1A)y_0 + \tau_1f. \end{aligned}$$

Veamos el *método de gradientes conjugados* que es de amplio uso en la práctica. En este método los parámetros iterativos α_{k+1} y τ_{k+1} se determinan según las fórmulas

$$\begin{aligned} \tau_{k+1} &= \frac{(r_k, w_k)}{(Aw_k, w_k)}, \\ \alpha_{k+1} &= \left(1 - \frac{\tau_{k+1}}{\tau_k} \frac{(r_k, w_k)}{(r_{k-1}, w_{k-1})} \frac{1}{\alpha_k}\right)^{-1}, \quad (9) \end{aligned}$$

donde $k = 0, 1, 2, \dots$, bajo el supuesto de que $A = A^* > 0$, $B = B^* > 0$, $\gamma_1 B \leq A \leq \gamma_2 B$, $\gamma_1 > 0$. Las fórmulas para τ_{k+1} , α_{k+1} se obtienen del requisito del mínimo de la norma del operador de resolución. Con estos valores optimales de los parámetros iterativos queda lícita la estimación

$$\|y_n - u\|_A \leq q_n \|y_0 - u\|_A, \quad q_n = \frac{2\rho_1^n}{1 + \rho_1^{2n}},$$

$$\rho_1 = \frac{1 - \sqrt{\xi}}{1 + \sqrt{\xi}}, \quad \xi = \frac{\gamma_1}{\gamma_2}, \quad (10)$$

es decir, la velocidad de convergencia del método de gradientes conjugados es la misma que la del método iterativo de dos capas con parámetros de Chébishev (en el que se usan γ_1 y γ_2 al calcular los parámetros τ_{k+1}). Por eso, para el número de iteraciones tenemos una estimación

$$n_0(\varepsilon) \leq n(\varepsilon) \leq n_0(\varepsilon) + 1, \quad n_0(\varepsilon) = \frac{1}{2\sqrt{\xi}} \ln \frac{2}{\varepsilon}.$$

A título de operador B podemos tomar el operador factorizado del método alternado triangular

$$B = (D + \omega A_1) D^{-1} (D + \omega A_2),$$

$$A_1 + A_2 = A > 0, \quad A_1^* = A_2, \quad D = D^* > 0.$$

Los cálculos muestran que el número de iteraciones al aplicar el método alternado triangular en conjunto con el de los gradientes conjugados es menor que en el caso de emplear el esquema de Chébishev.

§ 7. Resolución de las ecuaciones no lineales

1. **Métodos iterativos.** Analicemos una ecuación no lineal

$$f(x) = 0, \quad x \in [a, b],$$

donde $f(x)$ es una función continua. La ecuación puede tener una o varias raíces. Se pide; 1) establecer la existencia de las raíces de la ecuación; 2) hallar los valores aproximados de las raíces. Ambos problemas se resuelven a menudo simultáneamente. Para hallar las raíces se emplean los métodos iterativos.

El más elemental es el *método de dicotomía* (división por la mitad). Sea $f(x_0) f(x_1) \leq 0$; entonces, en el segmento $[x_0, x_1]$ se ubica por lo menos una raíz. Determinemos $f(x_2)$, donde $x_2 = (x_0 + x_1)/2$, y elijamos x_3 : aquél de los valores x_0 o x_1 , para el cual se cumple la condición $f(x_2) f(x_3) \leq 0$. El segmento $[x_2, x_3]$ se dividirá por la mitad de nuevo, etc. La división continúa hasta que la longitud del segmento se haga inferior a 2ε , donde ε es la exactitud con la que se debe determinar la raíz. En este caso el centro de dicho segmento nos presta precisamente el valor de la raíz con la exactitud requerida ε . Es evidente que el proceso converge con la velocidad de una progresión geométrica de razón $1/2$. La deficiencia del método consiste en la elección del segmento inicial $[x_0, x_1]$: no está claro de antemano a qué raíz convergerá el proceso (si en $[x_0, x_1]$ hay varias raíces).

El segundo método es el de *iteración simple*. Escribamos la ecuación (1) en la forma

$$x = \varphi(x), \quad (2)$$

donde $\varphi(x)$ puede definirse por uno de los siguientes métodos:

$$\varphi(x) = x - \alpha f(x), \quad \alpha = \text{const},$$

$\varphi(x) = x + \rho(x)f(x)$, $\rho(x)$ es una función arbitraria que no tiene raíces en el segmento $[a, b]$.

El método de iteración simple se determina por la fórmula

$$x_{n+1} = \varphi(x_n), \quad n = 0, 1, 2, \dots, \quad (3)$$

donde n es el número de la iteración, x_0 es la aproximación inicial prefijada arbitrariamente. Se pide hallar aproximadamente la solución (la raíz) $x = x^*$ de la ecuación $x = \varphi(x)$ con un error relativo $\varepsilon > 0$ de un modo tal que para cualquier $n \geq n_0$ se verifique la desigualdad

$$|x_n - x^*| \leq \varepsilon |x_0 - x^*|, \quad n \geq n_0(\varepsilon). \quad (4)$$

Esta condición puede cumplirse, siempre que la sucesión de iteraciones $\{x_n\}$ converja, para $n \rightarrow \infty$, al límite x^* : $\lim_{n \rightarrow \infty} x_n = x^*$. Si (4) tiene lugar, los cálculos pueden ser terminados con $n = n_0$. De aquí se ve que la cuestión más importante en este caso es la de convergencia de las iteracio-

nes, como también de la velocidad de su convergencia, es decir, la cuestión sobre el número mínimo de iteraciones $n_0(\epsilon)$, para el cual queda cumplida la desigualdad (4). Supongamos que en cierto δ -entorno

$$\Delta = (x_0 - \delta, x_0 + \delta), \quad \delta > 0, \quad (5)$$

del punto x_0 la función $\varphi(x)$ satisface la condición de Lipschitz:

$$|\varphi(x'') - \varphi(x')| \leq q |x'' - x'| \quad \text{para cualesquiera} \\ x', x'' \in \Delta \quad (6)$$

con el coeficiente $q < 1$:

$$0 < q < 1 \quad (7)$$

y sea pequeño el residuo inicial $x_0 - \varphi(x_0)$ de modo que

$$|x_0 - \varphi(x_0)| \leq (1 - q)\delta. \quad (8)$$

En este caso son justas las afirmaciones:

— todas las iteraciones x_n ($n = 1, 2, \dots$) pertenecen al intervalo $\Delta : x_n \in \Delta$;

— la sucesión $\{x_n\}$ converge, para $n \rightarrow \infty$, hacia el límite x^* que es la raíz de la ecuación (8);

— la ecuación (2) tiene en Δ una sola raíz.

La condición $x_h \in \Delta$ significa que

$$|x_h - x_0| < \delta. \quad (9)$$

En virtud de (8) tenemos $|x_1 - x_0| = |\varphi(x_0) - x_0| \leq (1 - q)\delta < \delta$, es decir, (9) se cumple para $k = 1$. Demostremos por el método de inducción que (9) se verifica para cualesquiera $k = 1, 2, \dots$. Supongamos que (9) se verifica para $k = 1, 2, \dots, n$; entonces se pueden calcular $\varphi(x_n)$ y $x_{n+1} = \varphi(x_n)$. De (6) se deduce que $|x_{k+1} - x_k| = |\varphi(x_k) - \varphi(x_{k-1})| \leq q |x_k - x_{k-1}|$, es decir,

$$|x_{k+1} - x_k| \leq q |x_k - x_{k-1}|. \quad (10)$$

Aplicando sucesivamente esta desigualdad, encontramos

$$|x_{k+1} - x_k| \leq q^k |x_1 - x_0|, \quad k = 1, 2, \dots, n. \quad (11)$$

Al tomar en consideración que $x_{n+1} - x_0 = (x_{n+1} - x_n) + (x_n - x_{n-1}) + \dots + (x_2 - x_1) + (x_1 - x_0)$, obtendremos

$$|x_{n+1} - x_0| \leq (q^n + q^{n-1} + \dots + q + 1) |x_1 - x_0| = \\ = \frac{1 - q^{n+1}}{1 - q} |x_1 - x_0| < \frac{1}{1 - q} |x_1 - x_0| < \delta,$$

es decir, $x_{n+1} \in \Delta$. En virtud de (8), la desigualdad (9) se verifica para $k = 1$, y, por lo tanto, se verifica también para $k = 2, 3, \dots$

Veamos ahora la diferencia $x_{n+m} - x_n = (x_{n+m} - x_{n+m-1}) + (x_{n+m-1} - x_{n+m-2}) + \dots + (x_{n+2} - x_{n+1}) + (x_{n+1} - x_n)$ y estimémosla:

$$|x_{n+m} - x_n| \leq (q^{m-1} + q^{m-2} + \dots + q + 1) |x_{n+1} - x_n| \leq \\ \leq \frac{1 - q^m}{1 - q} q^n |x_1 - x_0| < q^n \delta,$$

es decir, $|x_{n+m} - x_n| \rightarrow 0$ para $n \rightarrow \infty$ y cualquier $m = 1, 2, \dots$. De aquí, en virtud del criterio de Cauchy, proviene la convergencia de $\{x_n\}$: $\lim_{n \rightarrow \infty} x_n = x^* \in \Delta$. Pasando

ahora en (3) al límite para $n \rightarrow \infty$, nos convencemos de que x^* es una raíz de la ecuación (2): $x^* = \varphi(x^*)$. Esta raíz es única. En efecto, supongamos que existen dos raíces distintas x' y $x'' \neq x'$, de modo que $x' = \varphi(x')$, $x'' = \varphi(x'')$. Entonces, $|x'' - x'| = |\varphi(x'') - \varphi(x')| \leq q |x'' - x'| < |x'' - x'|$, es decir, $|x'' - x'| \leq |x'' - x'|$, lo que no es posible.

Para el error $z_{n+1} = x_{n+1} - x^*$ tenemos

$$|z_{n+1}| = |\varphi(x_n) - \varphi(x^*)| \leq q |x_n - x^*| = \\ = q |z_n| \leq q^{n+1} |z_0|, \quad (12) \\ |z_{n+1}| \leq q^{n+1} |z_0|,$$

es decir, el método de iteración simple converge con la velocidad de una progresión geométrica. El número de iteraciones para el cual queda cumplida la desigualdad (4) se determina de la condición $q^n \leq \varepsilon$, es decir,

$$n \geq \ln \frac{1}{\varepsilon} / \ln \frac{1}{q}.$$

El número mínimo de iteraciones $n_0(\varepsilon)$ para las cuales se cumple (4) es, evidentemente, igual a

$$n_0(\varepsilon) = \left[\ln \frac{1}{\varepsilon} / \ln \frac{1}{q} \right],$$

donde $[a]$ es la parte entera del número $a > 0$.

OBSERVACIONES. Si $\varphi(x)$ tiene derivada en Δ , entonces (6) se cumple en el caso en que

$$|\varphi'(x)| \leq q \text{ para todo } x \in \Delta. \quad (14)$$

2. Método de Newton. El método se determina mediante la fórmula

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}, \quad f'(x_n) \neq 0, \quad n = 0, 1, 2, \dots \quad (15)$$

Esta fórmula se obtiene, si en el desarrollo

$$0 = f(x^*) = f(x_n) + (x^* - x_n) f'(x_n) + \frac{1}{2} (x^* - x_n)^2 f''(\xi),$$

$$\xi = x_n + \theta (x^* - x_n) \quad 0 \leq \theta \leq 1, \quad (16)$$

donde x^* es la solución exacta de la ecuación $f(x) = 0$, se desecha el último término, al sustituir x^* por x_{n+1} :

$$0 = f(x_n) + f'(x_n) (x_{n+1} - x_n).$$

El método de Newton se denomina también *método de tangentes o de linearización*. La interpretación geométrica de este método consiste en que un trozo de la curva $y = f(x)$ para $x \in [x_n, x_{n+1}]$, si $x_n < x_{n+1}$ (o bien para $x \in [x_{n+1}, x_n]$, si $x_n > x_{n+1}$) se sustituye por el segmento de una tangente trazada desde punto $x = x_n$.

Al escribir $f(x) = 0$ en la forma $x = \varphi(x)$, vemos que el método de Newton puede ser considerado como el método de iteración simple (3) con el segundo miembro

$$\varphi(x) = x - f(x)/f'(x). \quad (17)$$

Ilustremos el método de Newton con el ejemplo de extracción de una raíz cuadrada de un número $a > 0$, es decir, de resolución de la ecuación $x^2 = a$ o $f(x) = x^2 - a = 0$. Al aplicar la fórmula (15), obtendremos

$$x_{n+1} = \frac{1}{2} \left(x_n + \frac{a}{x_n} \right), \quad n = 0, 1, \dots$$

Sea, $a = 2$. Eligiendo $x_0 = 1$, hallemos $x_1 = 1,5$, $x_2 = 1,417$, $x_3 = 1,414$, . . ., es decir, la iteración converge muy rápidamente.

Estimemos la velocidad de convergencia de las iteraciones. Supongamos que existe una raíz real x^* de la ecuación (1). Tomemos cierto entorno de la raíz:

$$\Delta_0 = (x^* - \delta_0, x^* + \delta_0), \quad \delta_0 > 0.$$

Convengamos en considerar que la función (17) es dos veces diferenciable en Δ_0 y que su derivada segunda está acotada

$$|\varphi''(x)| \leq 2q, \quad (18)$$

donde $q > 0$ es una constante. Desarrollemos $\varphi(x)$ en una línea de Taylor en el entorno de $x = x^*$:

$$\begin{aligned} \varphi(x) &= \varphi(x^*) + \varphi'(x^*)(x-x^*) + \frac{\varphi''(\xi)}{2}(x-x^*)^2, \\ \xi &= x^* + \theta(x-x^*), \quad 0 \leq \theta \leq 1. \end{aligned} \quad (19)$$

Calculando a continuación

$$\varphi'(x) = ff''(f')^2 = -f(1/f')', \quad \psi''(x) = -\left(f\left(\frac{1}{f'}\right)'\right)'$$

y observando que $\varphi'(x^*) = 0$ cuando $f'(x^*) \neq 0$, obtendremos

$$\varphi(x_n) = \varphi(x^*) + \frac{(x_n - x^*)^2}{2} \varphi''(\xi). \quad (20)$$

Para el error $z_{n+1} = x_{n+1} - x^*$ obtendremos una fórmula:

$$z_{n+1} = x_{n+1} - x^* = \varphi(x_n) - \varphi(x^*) = \frac{1}{2}(x_n - x^*)^2 \varphi''(\xi),$$

$$z_{n+1} = \frac{1}{2} \varphi''(\xi) z_n^2.$$

De aquí y de (20) se desprende

$$|z_{n+1}| \leq qz_n^2. \quad (21)$$

Denotando $v_n = q|z_n|$, obtenemos $v_{n+1} \leq v_n^2 \leq v_{n-1}^{2^2} \leq \dots \leq \dots \leq v_1^{2^n} \leq v_0^{2^{n+1}}$, y, por consiguiente,

$$|z_{n+1}| \leq \frac{1}{q} (q|z_0|)^{2^{n+1}}. \quad (22)$$

De aquí se ve que las iteraciones (15) convergen hacia la raíz x^* para $n \rightarrow \infty$, si

$$q |z_0| < 1 \text{ o } |z_0| = |x_0 - x^*| < 1/q, \quad (23)$$

es decir, la aproximación inicial se dispone en el entorno $\Delta_0 = (x^* - 1/q, x^* + 1/q)$ con $\delta_0 = 1/q$ de la raíz $x = x^*$ de la ecuación (1). En este caso el método de Newton converge, como suele decirse, con la *velocidad cuadrática* (el método de iteración simple converge con la velocidad de una progresión geométrica).

La condición para que terminen las iteraciones $|z_n| \leq \varepsilon |z_0|$ (como se infiere de (22)) o $|z_n| \leq (q |z_0|)^{2^{n-1}} \times |z_0|$ se cumple, si $n \geq n_0(\varepsilon)$, donde

$$n_0(\varepsilon) = \left[\ln \left(1 + \ln \frac{1}{\varepsilon} / \ln \frac{1}{q |z_0|} \right) / \ln 2 \right]. \quad (24)$$

Es evidente que si la aproximación inicial se dispone en el entorno pequeño de x^* , entonces todas las iteraciones posteriores quedarán dentro de este entorno Δ_0 . En efecto, sea $|x_0 - x^*| \leq \delta_0$, con la particularidad de que $q\delta_0 < 1$. Tendremos, pues, $|x_1 - x^*| \leq q |x_0 - x^*|^2$, $q\delta_0^2 < \delta_0$, $|x_2 - x^*| \leq q |x_1 - x^*|^2 \leq q\delta_0 < \delta_0$, etc., de suerte que $|x_n - x^*| \leq \delta_0$ para cualquier $n = 1, 2, \dots$

OBSERVACIONES. 1. No nos detenemos en la demostración de la existencia de la raíz $x = x^*$.

2. La convergencia cuadrática del método de Newton puede establecerse también para las restricciones más débiles impuestas sobre $f(x)$:

$$|f'(x)| \geq M_1 > 0, \quad |f''(x)| \leq M_2 \text{ para todo } x \in \Delta_0. \quad (25)$$

Haciendo uso de (15) y (16), obtendremos para el error $z_{n+1} = x_{n+1} - x^*$ una expresión

$$z_{n+1} = \frac{f''(\xi)}{2f'(\xi)} z_n^2,$$

de la cual, en virtud de las condiciones (25), proviene una desigualdad

$$|z_{n+1}| \leq q |z_n|^2, \quad q = M_2 / (2M_1),$$

que coincide con (21) (la diferencia consiste sólo en q). Los razonamientos ulteriores nos llevan a (22), (23), y (24).

3. Método continuo de Newton. La solución de la ecuación $f(x) = 0$ puede considerarse como un límite, para $t \rightarrow \infty$, de la solución del problema de Cauchy:

$$\frac{dx}{dt} + f(x) = 0, \quad x > 0, \quad x(0) = u_0, \quad (26)$$

si este límite existe. Denotemos con $x = x(t)$ la solución del problema de Cauchy, y con x_* , la solución de la ecuación $f(x) = 0$. Para su diferencia $z(t) = x(t) - x_*$, tenemos

$$\frac{dz}{dt} + (f(x) - f(x_*)) = \frac{dz}{dt} + f'(\xi) \cdot z, \quad \xi = x_* + \theta z, \quad 0 \leq \theta \leq 1,$$

$$\frac{dz}{dt} + \alpha(t) z = 0, \quad t > 0, \quad z(0) = u_0, \quad \alpha(t) = f'(\xi).$$

De aquí se ve que $|z(t)| \rightarrow 0$ para $t \rightarrow \infty$, si $f'(x) > 0$.

Para resolver la ecuación (26) se debe hacer uso de un método explícito cualquiera. La velocidad de convergencia de $x(t)$ a x_0 depende sólo de la magnitud de la derivada $f'(x)$.

4. Método de las secantes. El cálculo de la derivada $f'(x_n)$, aplicado el método de Newton, puede resultar engorroso. Si sustituimos f'_n por una razón de diferencias $(f_n - f_{n-1})/(x_n - x_{n-1})$, obtendremos el método iterativo de las secantes

$$x_{n+1} = x_n - \frac{(x_n - x_{n-1}) f(x_n)}{f(x_n) - f(x_{n-1})}. \quad (27)$$

El método de las secantes converge con una velocidad menor que el de Newton, sin embargo en (27) se calcula sólo la función, mientras que en (15) es necesario hallar no sólo la función sino también la derivada de ella. Es por esto que el volumen de los cálculos en cada iteración del método de las secantes es, en el caso general, menor.