

Solution Methods:
Iterative Techniques
Lecture 6

1 Motivation

SLIDE 1

Consider a standard second order finite difference discretization of

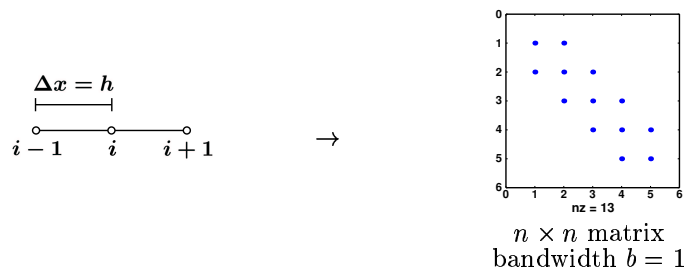
$$-\nabla^2 u = f,$$

on a regular grid, in 1, 2, and 3 dimensions.

$$\Rightarrow \boxed{A u = f}$$

1.1 1D Finite Differences

SLIDE 2

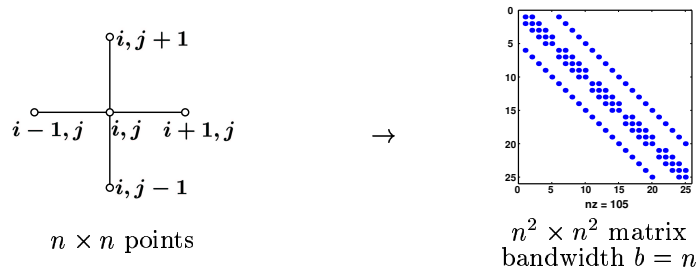


n points

$$\boxed{\text{Cost of Gaussian elimination } O(b^2 n) = O(n)}$$

1.2 2D Finite Differences

SLIDE 3

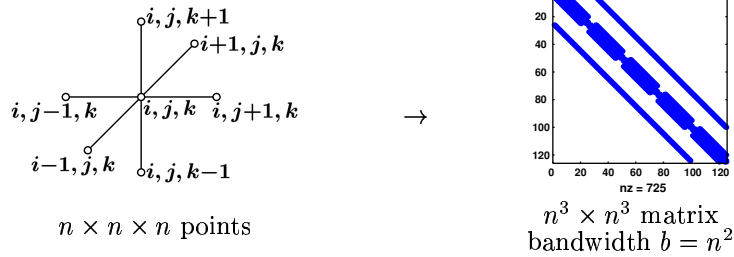


$n \times n$ points

$$\boxed{\text{Cost of Gaussian elimination } O(b^2 n^2) = O(n^4)}$$

1.3 3D Finite Differences

SLIDE 4



Cost of Gaussian elimination $O(b^2 n^3) = O(n^7)!$

This means that we halve the grid spacing, we will have 8 times (2^3) more unknowns and the cost of solving the problem will increase by a factor of 128 (2^7). It is apparent that, at least for practical three dimensional problems, faster methods are needed.

2 Basic Iterative Methods

2.1 Jacobi

2.1.1 Intuitive Interpretation

SLIDE 5

Instead of solving

$$-u_{xx} = f,$$

we solve

$$\frac{\partial u}{\partial t} = u_{xx} + f, \quad \boxed{N1}$$

starting from an arbitrary $u(x, 0)$.

We expect $u(x, t \rightarrow \infty) \rightarrow u(x)$.

That is, we expect the solution of the time dependent problem to converge to the solution of the steady state problem.

Note 1

Solution by relaxation

By adding the time dependent term $\frac{\partial u}{\partial t}$, we have now a parabolic equation. If the boundary conditions are kept fixed, we expect the solution to “converge” to the steady state solution (i.e., $\frac{\partial u}{\partial t} = 0$). Recall that the time dependent heat transfer problem is modeled by this equation. In this case, we expect the temperature to settle to a steady state distribution provided that the heat source, f , and the boundary conditions, do not depend on t .

To solve

$$\frac{\partial u}{\partial t} = u_{xx} + f$$

we use an inexpensive (explicit) method. *Thus avoiding the need to solve a linear system of equations.*

For instance,
$$\frac{\hat{u}_i^{r+1} - \hat{u}_i^r}{\Delta t} = \frac{\hat{u}_{i+1}^r - 2\hat{u}_i^r + \hat{u}_{i-1}^r}{h^2} + \hat{f}_i$$

$$\mathbf{u} = \{\hat{u}_i\}_{i=1}^n, \quad \mathbf{f} = \{\hat{f}_i\}_{i=1}^n$$

Here, we use the super index r to denote iteration (or time level). \mathbf{u} will denote the solution vector. An approximation to \mathbf{u} at iteration r will be denoted by \mathbf{u}^r .

$$\mathbf{u}^{r+1} = \mathbf{u}^r + \Delta t(\mathbf{f} - A\mathbf{u}^r) = (I - \Delta t A) \mathbf{u}^r + \Delta t \mathbf{f} .$$

Stability dictates that

$$\Delta t \leq \frac{h^2}{2}$$

Thus, we take Δt as large as possible, i.e. ($\Delta t = h^2/2$).

$$\mathbf{u}^{r+1} = \left(I - \frac{h^2}{2} A \right) \hat{\mathbf{u}}^r + \frac{h^2}{2} \mathbf{f}$$

$$\Rightarrow \quad \hat{u}_i^{r+1} = \frac{1}{2} (\hat{u}_{i+1}^r + \hat{u}_{i-1}^r + h^2 \hat{f}_i) \quad \text{for } i = 1, \dots, n.$$

2.1.2 Matrix Form

Split A

$$A = D - L - U \quad \begin{cases} D: & \text{Diagonal} \\ L: & \text{Lower triangular} \\ U: & \text{Upper triangular} \end{cases}$$

$$A \mathbf{u} = \mathbf{f} \quad \text{becomes} \quad (D - L - U) \mathbf{u} = \mathbf{f}$$

Iterative method

$$D \mathbf{u}^{r+1} = (L + U) \mathbf{u}^r + \mathbf{f}$$

$$\begin{aligned} \mathbf{u}^{r+1} &= D^{-1}(L + U) \mathbf{u}^r + D^{-1} \mathbf{f} \\ &= D^{-1}(D - A) \mathbf{u}^r + D^{-1} \mathbf{f} \\ &= (I - D^{-1}A) \mathbf{u}^r + D^{-1} \mathbf{f} \end{aligned}$$

$$\mathbf{u}^{r+1} = R_J \mathbf{u}^r + D^{-1} \mathbf{f}$$

$R_J = (I - D^{-1}A)$: Jacobi Iteration Matrix

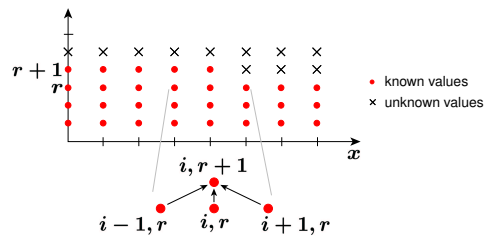
$$D_{ii}^{-1} = h^2/2$$

Note that, in order to implement this method we would typically not form any matrices, but update the unknowns, one component at a time, according to

$$\Rightarrow \hat{u}_i^{r+1} = \frac{1}{2} (\hat{u}_{i+1}^r + \hat{u}_{i-1}^r + h^2 \hat{f}_i) \quad \text{for } i = 1, \dots, n$$

2.1.3 Implementation

SLIDE 10



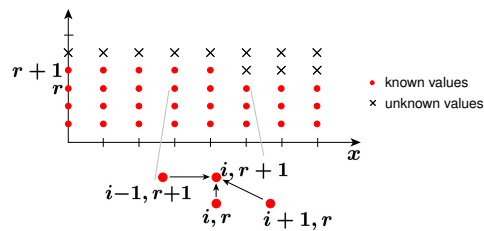
Jacobi iteration

$$u_i^{r+1} = \frac{1}{2} (u_{i+1}^r + u_{i-1}^r + h^2 f_i)$$

2.2 Gauss-Seidel

SLIDE 11

Assuming we are updating the unknowns according to their index i , at the time u_i^{r+1} needs to be calculated, we already know u_{i-1}^{r+1} . The idea of Gauss-Seidel iteration is to always use the most recently computed value.



Gauss-Seidel iteration (consider most recent iterate)

$$u_i^{r+1} = \frac{1}{2} (u_{i+1}^r + u_{i-1}^{r+1} + h^2 f_i)$$

2.2.1 Matrix Form

SLIDE 12

Split A

$$A = D - L - U \quad \begin{cases} D: & \text{Diagonal} \\ L: & \text{Lower triangular} \\ U: & \text{Upper triangular} \end{cases}$$

$$A \mathbf{u} = \mathbf{f} \text{ becomes } (D - L - U) \mathbf{u} = \mathbf{f}$$

Iterative method

$$(D - L) \mathbf{u}^{r+1} = U \mathbf{u}^r + \mathbf{f}$$

The above matrix form assumes that we are updating through our unknowns in ascending order. If we were to update in reverse order, i.e., the last unknown first, the iteration matrix would become $R_{GS} = (D - U)^{-1}L$.

Note 2

Updating order

We see that, unlike in the Jacobi method, the order in which the unknowns are updated in Gauss-Seidel changes the result of the iterative procedure. One could sweep through the unknowns in ascending, descending, or alternate orders. The latter procedure is called symmetric Gauss-Seidel.

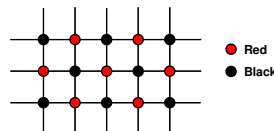
Another effective strategy, known as red-black Gauss-Seidel iteration, is to update the even unknowns first (red)

$$u_{2i}^{r+1} = \frac{1}{2} (u_{2i+1}^r + u_{2i-1}^r + h^2 f_{2i}) ,$$

and then the odd components (black)

$$u_{2i+1}^{r+1} = \frac{1}{2} (u_{2i+1}^{r+1} + u_{2i}^{r+1} + h^2 f_{2i+1}) .$$

The red-black Gauss-Seidel iteration is popular for parallel computation since the red (black) points only require the black (red) points and these can be updated in any order. The method readily extends to multiple dimensions. In 2D, for instance, the red and black points are shown below.



SLIDE 13

$$\mathbf{u}^{r+1} = (D - L)^{-1}U \mathbf{u}^r + (D - L)^{-1} \mathbf{f}$$

$$\boxed{\mathbf{u}^{r+1} = R_{\text{GS}} \mathbf{u}^r + (D - L)^{-1} \mathbf{f}}$$

$R_{\text{GS}} = (D - L)^{-1}U$: **Gauss-Seidel Iteration Matix**

2.3 Error Equation

SLIDE 14

Let \mathbf{u} be the solution of $A \mathbf{u} = \mathbf{f}$.

For an approximate solution \mathbf{u}^r , we define

$$\text{Iteration Error: } \mathbf{e}^r = \mathbf{u} - \mathbf{u}^r$$

$$\text{Residual: } \mathbf{r}^r = \mathbf{f} - A \mathbf{u}^r$$

Subtracting $A \mathbf{u}^r$ from both sides of $A \mathbf{u} = \mathbf{f}$,

$$A \mathbf{u} - A \mathbf{u}^r = \mathbf{f} - A \mathbf{u}^r$$

$$\text{ERROR EQUATION} \rightarrow \boxed{A \mathbf{e}^r = \mathbf{r}^r}$$

N3

We note that, whereas the iteration error is not an accessible quantity, since it requires \mathbf{u} , the residual is easily computable and only requires \mathbf{u}^r .

Note 3 *Relationship between error and residual.*

We have seen that the residual is easily computable and in some sense it measures the amount by which our approximate solution \mathbf{u}^r fails to satisfy the original problem.

It is clear that if $\mathbf{r} = \mathbf{0}$, we have $\mathbf{e} = \mathbf{0}$. However, it is not always the case that when \mathbf{r} is small in norm, \mathbf{e} is also small in norm.

We have that $A \mathbf{u} = \mathbf{f}$ and $A^{-1} \mathbf{r} = \mathbf{e}$. Taking norms, we have

$$\|\mathbf{f}\|_2 \leq \|A\|_2 \|\mathbf{u}\|_2, \quad \|\mathbf{e}\|_2 \leq \|A^{-1}\|_2 \|\mathbf{r}\|_2 .$$

Combining these two expressions, we have

$$\frac{\|\mathbf{e}\|_2}{\|\mathbf{u}\|_2} \leq \underbrace{\|A\|_2 \|A^{-1}\|_2}_{\text{cond}(A)} \frac{\|\mathbf{r}\|_2}{\|\mathbf{f}\|_2} .$$

From here we see that if our matrix is not well conditioned, i.e., $\text{cond}(A)$ is large, then small residuals can correspond to significant errors.

2.3.1 Jacobi

SLIDE 15

$$\begin{aligned} \mathbf{u}^{r+1} &= D^{-1}(L + U) \mathbf{u}^r + D^{-1}f \\ \mathbf{u} &= D^{-1}(L + U) \mathbf{u} + D^{-1}f \end{aligned}$$

subtracting

$$\mathbf{e}^{r+1} = \underbrace{D^{-1}(L + U)}_{R_J} \mathbf{e}^r = R_J \mathbf{e}^r$$

2.3.2 Gauss-Seidel

SLIDE 16

Similarly,

$$\mathbf{e}^{r+1} = \underbrace{(D - L)^{-1}U}_{R_{GS}} \mathbf{e}^r = R_{GS} \mathbf{e}^r$$

It is clear that the above equations satisfied by the error are not useful for practical purposes since $\mathbf{e}^0 = \mathbf{u} - \mathbf{u}^0$, will not be known before the problem is solved. We will see however, that by studying the equation satisfied by the error we can determine useful properties regarding the convergence of our schemes.

2.4 Examples

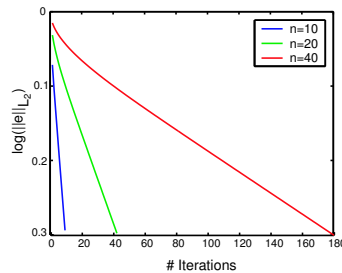
2.4.1 Jacobi

SLIDE 17

$$-u_{xx} = 1$$

$$u(0) = u(1) = 0;$$

$$\mathbf{u}^0 = \mathbf{0}$$



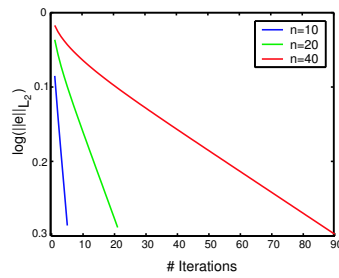
2.4.2 Gauss-Seidel

SLIDE 18

$$-u_{xx} = 1$$

$$u(0) = u(1) = 0;$$

$$\mathbf{u}^0 = \mathbf{0}$$



2.4.3 Observations

SLIDE 19

- The number of iterations required to obtain a certain level of convergence is $O(n^2)$.
- Gauss-Seidel is about twice as fast as Jacobi.

Why?

3 Convergence Analysis

SLIDE 20

$$e^r = R e^{r-1} = R R e^{r-2} = \dots = R^r e^0 .$$

The iterative method will converge if

$$\lim_{r \rightarrow \infty} R^r = 0 \iff \rho(R) = \max |\lambda(R)| < 1 .$$

N4

$\rho(R)$ is the spectral radius.

The spectral radius is the radius of the smallest circle centered at the origin that contains all the eigenvalues. The above condition indicates that the necessary and sufficient condition for the scheme to converge, for arbitrary initial conditions, is that the largest eigenvalue, in modulus, lies within the unit circle.

Note 4

Condition on $\rho(R)$

The condition that

$$\lim_{r \rightarrow \infty} R^r = 0 \quad \text{if and only if} \quad \rho(R) < 1$$

can be easily shown for the case where R is diagonalizable. In this case $R = X^{-1} \Lambda X$, where Λ is the diagonal matrix of eigenvalues. We have that $R^r = X^{-1} \Lambda^r X$ and hence we require that

$$\lim_{r \rightarrow \infty} \Lambda^r = 0$$

which is equivalent to $\rho(R) < 1$.

3.1 Theorems

SLIDE 21

- If the matrix A is *strictly diagonally dominant* then Jacobi and Gauss-Seidel iterations converge starting from an arbitrary initial condition. N5

$A = \{a_{ij}\}$ is strictly diagonally dominant if

$$|a_{ii}| > \sum_{j=1, j \neq i}^n |a_{ij}| \text{ for all } i.$$

- If the matrix A is symmetric positive definite then Gauss-Seidel iteration converges for any initial solution. N6

Note 5 **Convergence of Jacobi iteration for diagonally dominant matrices**

For the Jacobi iteration, $R_J = D^{-1}(L + U)$, it can be easily seen that R has zero diagonal entries. It can be shown (Gershgorin Theorem - see GVL) that a matrix with zeros in the main diagonal, $\rho(R_J) \leq \|R_J\|_\infty$. Hence

$$\rho(R_J) = \|R_J\|_2 \leq \|R_J\|_\infty = \max_{1 \leq j \leq n} \sum_{\substack{j=1 \\ j \neq i}}^n \left| \frac{a_{ij}}{a_{ii}} \right| < 1,$$

the last inequality follows from the diagonal dominance of A .

Note 6 **Convergence of Gauss-Seidel for SPD matrices**

(See Theorem 10.1.2 of GVL)

3.2 Jacobi

SLIDE 22

Unfortunately none of the above convergence theorems helps us to show that Jacobi iteration will converge for our model problem. However, since for our model problem we can calculate the eigenvalues of A explicitly (see previous lectures), convergence can be shown directly.

$$R_J = D^{-1}(L + U) = D^{-1}(D - A) = I - D^{-1}A = I - \frac{h^2}{2}A$$

If $A\mathbf{v}^k = \lambda^k \mathbf{v}^k$,

then $R_J \mathbf{v}^k = \left(I - \frac{h^2}{2}A \right) \mathbf{v}^k = \underbrace{\left(1 - \frac{h^2}{2} \lambda^k(A) \right)}_{\lambda^k(R_J)} \mathbf{v}^k$

$$\lambda^k(R_J) = 1 - \frac{h^2}{2} \lambda^k(A)$$

Eigenvectors of $R_J \equiv$ Eigenvectors of A

SLIDE 23

Recall ...

$$A = \frac{1}{h^2} \begin{pmatrix} 2 & -1 & 0 & \cdots & 0 \\ -1 & 2 & -1 & \ddots & \vdots \\ 0 & \ddots & \ddots & \vdots & 0 \\ \vdots & \ddots & -1 & 2 & -1 \\ 0 & \cdots & 0 & -1 & 2 \end{pmatrix} \begin{matrix} A\mathbf{v}^k = \lambda^k \mathbf{v}^k & k = 1, \dots, n \\ \text{Eigenvectors:} & \boxed{h = \frac{1}{n+1}} \\ v_j^k = \sin(k\pi h j) = \sin\left(\frac{k\pi j}{n+1}\right) \\ \text{Eigenvalues:} & \lambda^k(A) = \frac{2}{h^2} [1 - \cos(k\pi h)] \end{matrix}$$

$n \times n$ *SPD*

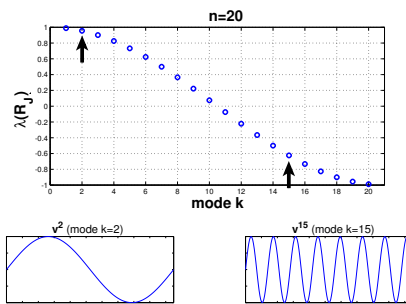
SLIDE 24

$$\begin{aligned} \lambda^k(R_J) &= 1 - \frac{h^2}{2} \lambda^k(A) = 1 - [1 - \cos k\pi h] \\ &= \cos \frac{k\pi}{n+1} < 1, \quad k = 1, \dots, n \end{aligned}$$

Jacobi converges for our model problem.

$$R_J \mathbf{v}^k = \lambda^k(R_J) \mathbf{v}^k$$

SLIDE 25



SLIDE 26

Since A has a complete set of eigenvectors \mathbf{v}^k , $k = 1, \dots, n$,

$$\text{Write } \mathbf{e}^0 = \sum_{k=1}^n c_k \mathbf{v}^k \quad \mathbf{v}^k, k\text{-th eigenvector of } R_J$$

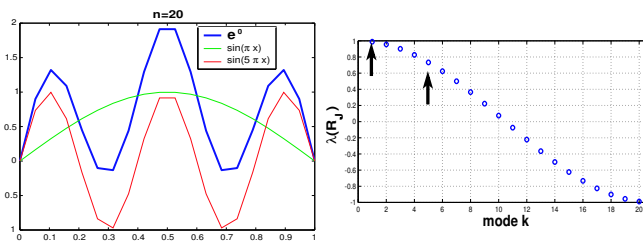
$$\mathbf{e}^1 = R \mathbf{e}^0 = \sum_{k=1}^n c_k \lambda^k(R_J) \mathbf{v}^k$$

$$\mathbf{e}^r = R^r \mathbf{e}^0 = \sum_{k=1}^n c_k (\lambda^k(R_J))^r \mathbf{v}^k$$

3.2.1 Example

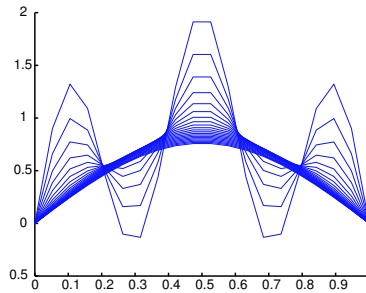
SLIDE 27

We consider the situation in which the initial error only contains modal components corresponding to $k = 1$ and $k = 5$.



SLIDE 28

The figure below show the evolution of the error for each iteration. The faster decay of the $k = 5$ mode can be readily observed.



3.2.2 Convergence rate

$$\lambda^k(R_J) = \cos(k\pi h), \quad k = 1, \dots, n$$

SLIDE 29

Largest $|\lambda^k(R_J)|$ for $k = 1, \dots, n$

Worst case $e^0 = c_1 v^1 \rightarrow e^r = c_1 \rho(R_J)^r v^1$

$$\frac{\|e^r\|}{\|e^0\|} = (\cos(\pi h))^r \simeq \left(1 - \frac{\pi^2 h^2}{2}\right)^r$$

SLIDE 30

For h small (n large) $\pi h \simeq 0$ and $\cos(\pi h)$ can be approximated as $1 - \frac{\pi^2 h^2}{2}$.

To obtain a given level of convergence; e.g., $10^{-\delta}$

$$\frac{\|e^r\|}{\|e^0\|} < 10^{-\delta}$$

$$\Rightarrow \left(1 - \frac{\pi^2 h^2}{2}\right)^r < 10^{-\delta} \rightarrow r = \frac{-\delta}{\log(1 - \frac{\pi^2 h^2}{2})} \simeq \frac{2\delta}{\pi^2 h^2} = \frac{2\delta(n+1)^2}{\pi^2}$$

We have used the fact that for x small $\log(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} \dots$

$$\rightarrow \boxed{r = O(n^2)}$$

The number of iterations required to obtain a given level of convergence is n^2 .

3.2.3 Convergence rate (2D)

SLIDE 31

This analysis can be extended to two and three dimensions in a straightforward manner.

In two dimensions, and for a uniform, $n \times n$, grid, we have

$$\begin{aligned}\lambda^{k\ell}(R_J) &= 1 - \frac{h^2}{4} \lambda(A) = \frac{1}{2} [\cos(k\pi h) + \cos(\ell\pi h)] \\ \rho(R_J) &= \cos(\pi h) = \cos\left(\frac{\pi}{h+1}\right) \quad \boxed{h = \frac{1}{n+1}}\end{aligned}$$

Therefore, $\boxed{r = O(n^2)}$

It is important to note that in 1D, 2D and also in 3D, $r = O(n^2)$. This, combined with the fact that the number of unknowns in 1, 2, and 3D is n , n^2 and n^3 , makes iterative methods potentially attractive for 2, and specially, 3D problems.

3.3 Gauss-Seidel

SLIDE 32

The convergence analysis for the Gauss-Seidel method is analogous to that of the Jacobi method. The expressions for the eigenvectors and eigenvalues, for the model problems, can also be found in closed although the procedure is somewhat more involved.

$$R_{GS} = (D - L)^{-1} U$$

It can be shown (try it !) that

$$\boxed{\lambda^k(R_{GS}) = \cos^2(k\pi h) = [\lambda^k(R_J)]^2 < 1}$$

Gauss-Seidel converges for our problem

The eigenvalues are always positive.

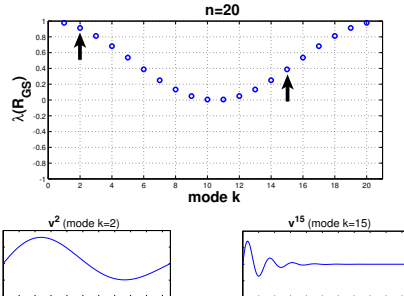
But, Eigenvectors of $R_{GS} \neq$ Eigenvectors of A

The eigenvectors $\mathbf{v}^k = \{v_j^k\}_{j=1}^n$ of R_{GS} are

$$v_j^k = [\sqrt{\lambda^k(R_{GS})}]^j \sin(k\pi h j)$$

These eigenvectors, although not orthogonal, still form a complete basis and can be used to represent the error components.

SLIDE 33



3.3.1 Convergence rate

SLIDE 34

To obtain a given level of convergence; e.g., $10^{-\delta}$

$$\frac{\|e^r\|}{\|e^0\|} < 10^{-\delta}$$

$$\Rightarrow \left(1 - \frac{\pi^2 h^2}{2}\right)^{2r} < 10^{-\delta} \rightarrow r = \frac{-\delta}{2 \log\left(1 - \frac{\pi^2 h^2}{2}\right)} \simeq \frac{\delta}{\pi^2 h^2} = \frac{\delta(n+1)^2}{\pi^2}$$

$$\rightarrow \boxed{r = O(n^2)}$$

Same rate as Jacobi, but requires only half the number of iterations.

4 Comparative cost

SLIDE 35

		Iteration	Gauss Elimination
1D	$n \times n$	$O(n^2 n) = O(n^3)$	$O(n)$
2D	$n^2 \times n^2$	$O(n^2 n^2) = O(n^4)$	$O(n^4)$
3D	$n^3 \times n^3$	$O(n^2 n^3) = O(n^5)$	$O(n^7)$

red # iters
green cost/iter

Only in 3D, give iterative methods some cost advantage over Gaussian elimination.

5 Over/Under Relaxation

5.0.2 Main Idea

SLIDE 36

Typically

$$\mathbf{u}^{r+1} = R \mathbf{u}^r + f^* \quad \begin{array}{ll} f^* = D^{-1} f & \text{Jacobi} \\ f^* = (D - L)^{-1} f & \text{Gauss-Seidel} \end{array}$$

Can we “extrapolate” the changes?

$$\begin{aligned} \mathbf{u}^{r+1} &= \omega(R\mathbf{u}^r + f^*) + (1 - \omega) \mathbf{u}^r \\ &= \underbrace{[\omega R + (1 - \omega) I]}_{R_\omega} \mathbf{u}^r + \omega f^* \quad \boxed{\omega > 0} \end{aligned}$$

5.1 How large can we choose ω ?

SLIDE 37

$$\lambda^k(R_\omega) = \omega \lambda^k(R) + (1 - \omega)$$

Jacobi $\lambda^k(R_J) = \cos k\pi h$

GS $\lambda^k(R_{GS}) = \cos^2 k\pi h$

$$\Rightarrow \rho(R_\omega) < 1 \Rightarrow \begin{array}{ll} 0 \leq \omega_J \leq 1 & \text{can only be} \\ & \text{under-relaxed} \\ 0 \leq \omega_{GS} \leq 2 & \text{can be over-relaxed} \end{array}$$

N7

Note 7

Successive Over-Relaxation (SOR)

Over-relaxed Gauss-Seidel is about a factor of two faster than standard GS but still requires $O(n^2)$ iterations.

The idea of over-relaxation can be used to improve convergence further. If instead of extrapolating the changes for all the unknowns at the same time, we extrapolate the change of each unknown as soon as it has been computed, and use the updated value in subsequent unknown updates we obtain what is called successive over-relaxation (SOR). In a component form SOR is written as

$$u_i^{r+1} = \omega \left[\frac{1}{2} (u_{i+1}^r + u_{i-1}^{r+1}) + \frac{h^2}{2} f_i \right] + (1 - \omega) u_i^r .$$

with $1 < \omega < 2$ and typically $\omega \sim 1.5$.

In matrix form SOR can be written as we

$$\begin{aligned} \mathbf{u}^{r+1} &= (D - \omega L)^{-1} [((1 - \omega) D + \omega U) \mathbf{u} + \omega f] \\ &= R_{\text{SOR}} \mathbf{u}^r + (D - \omega L)^{-1} . \end{aligned}$$

It can be shown that the spectral radius of R_{SOR} , for our model problem, is $1 - O(h)$ and hence the number of iterations scales like n rather than n^2 . SOR was very popular 30 years ago. Nowadays, we can do better.

REFERENCES

[GVL] *G. Golub, C. van Loan, Matrix Computations, Johns Hopkins.*