

LECTURE 28: LINEAR OPERATOR SPLITTING

A first example of problem-specific algorithms is afforded by consideration of the steady-state diffusion problem

$$-\nabla^2 u = -\Delta u = f, \quad (1)$$

where $u(\mathbf{x})$, $\mathbf{x} \in \Omega \subseteq \mathbb{R}^d$ can be understood to denote the temperature in an infinitesimal volume positioned at \mathbf{x} in d -dimensional space, and f is a local rate of heat generation. The above arises from setting the time derivative zero in $\partial_t u - \Delta u = f$. Though often stated in this thermal language, the Poisson equation (1) is generally valid for unresolved transport by Brownian motion. The mathematical concept of an “infinitesimal volume” is interpreted as setting some length scale ℓ much smaller than the length scale L characterizing the size of the domain Ω . As an example, for heat conduction in a column of water of length $L = 1\text{m}$ the length scale of a quasi-infinitesimal volume can be considered as, say, $\ell = 1\text{ }\mu\text{m}$. There is no physical significance to the mathematical limit process $\ell \rightarrow 0$ due to the discrete structure of matter, and for all practical purposes setting $\ell = 1\text{ }\mu\text{m}$ is an acceptable threshold to delimit a phenomenon of interest to behavior that can be neglected. In this case the phenomenon of interest is the average transport of thermal energy in volumes of size greater than $\mathcal{O}(\ell^d)$. The detailed Brownian motion of the water molecules that occurs on length scales $s \approx 1\text{ nm} \ll \ell$ can be neglected and is said to be *unresolved*. The only observable effect of this motion is that temperature gradients lead to a heat flux as described by $\mathbf{f}(u) = -\alpha \nabla u$ (Fourier's law). The same equation (1) arises in epidemiology when $\ell = 10\text{m}$, an average separation between an infected and susceptible individual, $L = 10\text{km}$, the size of a large city, and u is reinterpreted as the percentage of infected individuals in the population. Again, the detailed Brownian steps of size $s \approx 10\text{cm} \ll \ell$ taken by individuals can be neglected.

1. Poisson equation discretization

Understanding the underlying unresolved Brownian motion is useful in constructing numerical solutions. For $f = 0$, (1) becomes $\nabla^2 u = 0$, which states that there is no net heat flux in an infinitesimal volume, $\text{div} \cdot (\alpha \nabla u) = 0$, colloquially: “what flows in on one side goes out on the other”. A function that satisfies Laplace's equation $\nabla^2 u = 0$ is said to be *harmonic*. For $d = 1$, the ordinary differential equation $\partial_x^2 u = 0$ is obtained with solution $u(x) = a + bx$, and boundary conditions u_0 at $x = 0$ and u_1 at $x = 1$ gives $u(x) = u_0 + (u_1 - u_0)x$. A temperature difference $u_0 \neq u_1$ at the boundaries induces diffusive fluxes that lead to the mean value $u(1/2) = (u_0 + u_1)/2$ at the midpoint. An analogous statement is made for $d > 1$ starting from Green's formula on ball B with spherical boundary S of radius R

$$\int_B (u \Delta v - v \Delta u) d\omega = \int_{S=\partial B} \left(u \frac{\partial v}{\partial n} - v \frac{\partial u}{\partial n} \right) d\sigma,$$

for scalar functions u, v . For $d = 3$ and $v = 1/r$ the mean value theorem

$$u(\mathbf{x}) = \frac{1}{4\pi R^2} \int_{S=\partial B} u(\mathbf{y}) d\sigma(\mathbf{y}), \quad (2)$$

is obtained, which states that the value of a harmonic function is the average of the values on a surrounding sphere. For $d = 2$ the analogous statement is

$$u(\mathbf{a}) = \frac{1}{2\pi r} \int_0^{2\pi} u(\mathbf{a} + r \cos \theta \mathbf{e}_x + r \sin \theta \mathbf{e}_y) d\theta. \quad (3)$$

Midpoint quadrature of (3) over four subintervals gives

$$u_0 = \frac{u_1 + u_2 + u_3 + u_4}{4}. \quad (4)$$

Apply the above on a grid covering some arbitrary domain Ω with boundary $\Sigma = \partial\Omega$ to obtain

$$u_{i,j} = \frac{1}{4} (u_{i,j-1} + u_{i-1,j} + u_{i+1,j} + u_{i,j+1}). \quad (5)$$

Complications arise for general boundaries Σ , but for a square domain $\Omega = [0, 1] \times [0, 1]$ grid points $(x_i = ih, y_j = jh)$ align with the boundary, $h = 1/(n+1)$. Instead of two indices, one for each spatial direction, organize the grid values u through a single index $k = (j-1)n + i$, with u_k denoting the value at the interior points (x_i, y_j) . The vector of interior grid values $\mathbf{u} = [u_1 \dots u_m]$ has $m = n^2$ components, and the mean value theorem leads to the linear system

$$\mathbf{A}\mathbf{u} = \mathbf{b}, \quad (6)$$

where A has a regular sparsity pattern induced by the uniform spacing of the grid

$$A = \begin{bmatrix} 4 & -1 & 0 & \dots & 0 & -1 & 0 & \dots & \dots & 0 \\ -1 & 4 & -1 & 0 & \dots & 0 & -1 & 0 & \dots & 0 \\ & \ddots & \ddots & \ddots & & & & \ddots & & \\ & & & & & & & & & \end{bmatrix}.$$

The vector \mathbf{b} contains the boundary values, for example $b_1 = u_{1,0} + u_{0,1}$.

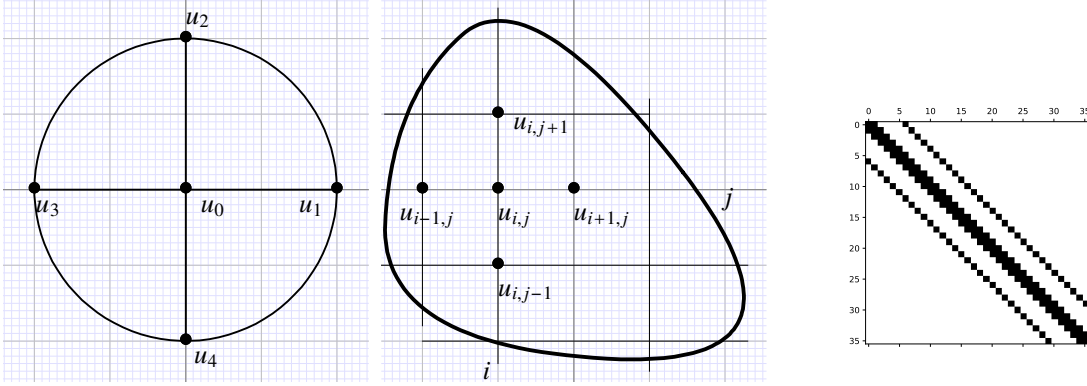


Figure 1. Left: Mean value theorem leads to $u_0 = (u_1 + u_2 + u_3 + u_4) / 4$. Middle: Five-point finite difference stencil for Laplace operator. Right: Structure of A matrix resulting from discretization of Laplace operator.

The system (6) was obtained by discretization of the mean-value integral (3). The same linear system is also obtained by discretization of the differential equation

$$-\nabla^2 u = -u_{xx} - u_{yy} = 0,$$

where indices denote differentiation. The minus sign arises from compatibility with the unsteady form of the heat equation $\partial_t u - \nabla^2 u = f$. A centered finite difference approximation of the derivatives on the uniform grid leads to

$$u_{xx}(x_i, y_j) \cong \frac{u_{i-1,j} - 2u_{i,j} + u_{i+1,j}}{h^2}, u_{yy}(x_i, y_j) \cong \frac{u_{i,j-1} - 2u_{i,j} + u_{i,j+1}}{h^2},$$

and (6) is recovered. For the Poisson equation $f \neq 0$ the right hand side changes to

$$A\mathbf{u} = \mathbf{b} + h^2 \mathbf{f} = \mathbf{c} \quad (7)$$

with $f_k = f(x_i, y_j)$. It is often the case that the same discrete system arises from both the differential and the integral formulation of a conservation law on a uniform grid.

2. Matrix splitting iteration

The above discussion of the underlying physics of the Poisson equation productively guides construction of numerical solution procedures. Solving the linear system (7) by general factorizations such as $A = LU$ or $A = QR$ is costly in terms of memory usage since the sparsity pattern is lost. For the uniform grid and square domain considered above the matrix A need not be explicitly stored at all since $a_{i,i} = 4$, and $a_{i,j} = -1$ when $|i - j| = 1$ or $|i - j| = n$. The discrete mean value theorem (5) suggests that some approximation $\mathbf{u}^{(l)}$ can be improved by the iteration

$$u_{i,j}^{(l+1)} = \frac{1}{4} (u_{i,j-1}^{(l)} + u_{i,j+1}^{(l)} + u_{i-1,j}^{(l)} + u_{i+1,j}^{(l)}). \quad (8)$$

The above is known as Jacobi iteration, and can be stated in matrix form by expressing A as

$$A = L + D + U, \quad (9)$$

with L, U containing non-zero components of A below, above the diagonal, and D containing the diagonal of A . In contrast to the *multiplicative* decompositions considered up to now, the QR, LU, SVD , eigen or Schur decompositions, the decomposition (9) is now *additive*. Note that in (9) L, U are strictly lower, upper diagonal matrices with zeros on the diagonal in contrast to the notation for the standard LU factorization algorithm. Recall that the utility of matrix multiplication was associated with the representation of linear mapping composition. Additive decompositions such as (9) generally are useful when separating different aspects of a physical process, and are a simple example of *operator splitting*. For the discrete Poisson system Jacobi iteration can be expressed as

$$Au = c \Rightarrow (L + D + U)u = c \Rightarrow u^{(l+1)} = D^{-1}(c - Lu^{(l)} - Uu^{(l)}). \quad (10)$$

Several variants of the idea can be pursued. The matrix splitting (9) is useful in theoretical convergence analysis, but implementations directly use (8) within loops over the (i, j) indices. Updated values can be immediately utilized leading to either of the following iterations

$$u^{(l+1)} = D^{-1}(c - Lu^{(l+1)} - Uu^{(l)}), u^{(l+1)} = D^{-1}(c - Lu^{(l)} - Uu^{(l+1)}), \quad (11)$$

depending on loop organization. These are known as *Gauss-Seidel iterations*. Convergence might be accelerated by extrapolation,

$$u^{(l+1)} = u^{(l)} + \omega [D^{-1}(c - Lu^{(l)} - Uu^{(l+1)}) - u^{(l)}] = (1 - \omega)u^{(l)} + \omega D^{-1}(c - Lu^{(l)} - Uu^{(l+1)}), \quad (12)$$

where the new iteration is continued by factor ω in the direction of the Gauss-Seidel update. When $\omega > 1$ this is known as successive over-relaxation (SOR) and goes further in the Gauss-Seidel direction. Choosing $0 < \omega < 1$ leads to successive under-relaxation.

3. Convergence analysis

Turning now from algorithm construction to analysis of its behavior, simplify notation by letting u_k denote the current iterate. The previous notation $u^{(l)}$ was convenient since individual components were referenced as in $u_{i,j}^{(l)}$, but convergence analysis is determined by the properties of the operator splitting and not of the current iterate. Introduce the error δ_k at iteration k as the difference between the exact solution u and the current iterate u_k , $\delta_k = u - u_k$. Also introduce the *residual* $r_k = c - Au_k = A\delta_k$, and the correction to the current iterate $e_k = u_{k+1} - u_k$.

The above methods can be formulated as a residual correction algorithm through the steps:

1. residual computation, $r_k = c - Au_k$
2. correction computation, $e_k = Br_k$
3. approximation update, $u_{k+1} = u_k + e_k$

When $B = A^{-1}$ the exact solution is recovered in one step

$$e_k = A^{-1}(c - Au_k) = u - u_k \Rightarrow u_{k+1} = u_k + e_k = u.$$

Iterative methods use some approximation of the (unknown) inverse, $B \cong A^{-1}$. Jacobi iteration uses $B = D^{-1}$ since

$$u_{k+1} = u_k + D^{-1}[c - (L + D + U)u_k] = D^{-1}[c - (L + U)u_k],$$

recovering (10). Table 1 shows several common choices for B . Two key aspects govern the choice of the inverse approximant:

1. Computational efficiency stated as a requirement that each iteration cost either $\mathcal{O}(m)$ or $\mathcal{O}(m \log m)$ operations;
2. Capturing the essential aspects of A .

Jacobi	D^{-1}	Forward Gauss-Seidel	$(D + L)^{-1}$
Weighted Jacobi	ωD^{-1}	Backward Gauss-Seidel	$(D + U)^{-1}$
SOR	$\omega(D + \omega L)^{-1}$	Symmetric Gauss-Seidel	$(D + U)^{-1}D(D + L)^{-1}$
Symmetric SOR	$\omega(2 - \omega)(D + \omega U)^{-1}D(D + \omega L)$	Richardson	ωI

Table 1. Common iterative methods

The iteration converges to the solution if $\|\delta_k\| \rightarrow 0$ for increasing k . The error at iteration $k + 1$ is expressed as

$$\delta_{k+1} = \mathbf{u} - \mathbf{u}_{k+1} = \mathbf{u} - (\mathbf{u}_k + \mathbf{B}\mathbf{r}_k) = \delta_k - \mathbf{B}\mathbf{A}(\mathbf{u} - \mathbf{u}_k) = (\mathbf{I} - \mathbf{B}\mathbf{A})\delta_k = (\mathbf{I} - \mathbf{B}\mathbf{A})^{k+1}\delta_0. \quad (13)$$

The repeated matrix multiplication indicates that the eigenstructure of the iteration matrix $\mathbf{M} = \mathbf{I} - \mathbf{B}\mathbf{A}$ determines iteration convergence. Indeed the above is simply power iteration for \mathbf{M} and can be expected to converge as

$$\delta_k \rightarrow \mu_1^k c_1 \mathbf{q}_1,$$

with (μ_1, \mathbf{q}_1) the eigenpair that corresponds to the largest eigenvalue in absolute value, $\mathbf{M}\mathbf{q}_1 = \mu_1 \mathbf{q}_1$, known as the spectral radius of \mathbf{M} , denoted by $\rho(\mathbf{M})$. Clearly, the above iterations will exhibit linear order of convergence when $\rho(\mathbf{M}) < 1$. The rate of convergence s_k at iteration k is estimated by the Rayleigh quotient

$$s_k = \frac{\delta_k^T \mathbf{M} \delta_k}{\delta_k^T \delta_k} = \frac{\delta_k^T \delta_{k+1}}{\delta_k^T \delta_k},$$

and is monitored in implementations of iterative methods, and determined by the eigenvalues $\lambda = 1 - \mu$ of $\mathbf{B}\mathbf{A}$.

The eigenstructure of $\mathbf{B}\mathbf{A}$ is difficult for arbitrary matrices \mathbf{A} , but can be carried out when \mathbf{A} has special structure induced by known physical phenomena. The relation between analytical and numerical formulations plays an essential role in convergence analysis. The diffusion equation (5) leads to a symmetric matrix \mathbf{A} , $\mathbf{A} = \mathbf{A}^T$ due to two aspects:

1. the chosen discretization is symmetric using centered finite differences;
2. the operator itself exhibits symmetry.

Insight into the above two aspects is most readily gained from the one-dimensional case $-u_{xx} = f$ with homogeneous Dirichlet boundary conditions $u(0) = u(1) = 0$. The linear system $\mathbf{A}\mathbf{u} = \mathbf{f}$ obtained from the centered finite difference discretization

$$-u_{i-1} + 2u_i - u_{i+1} = h^2 f_i, i = 1, \dots, m, h = 1/(m+1), u_0 = u_{m+1} = 0,$$

has a symmetric tridiagonal system matrix $\mathbf{A} = \text{diag}([-1 \ 2 \ -1])$. The $\mathbf{A}^T = \mathbf{A}$ symmetry can be expressed through scalar products in a way that generalizes to differential operators. Recall that a real-valued scalar product (\mathbf{u}, \mathbf{v}) must satisfy symmetry $(\mathbf{v}, \mathbf{u}) = (\mathbf{u}, \mathbf{v})$. For $\mathbf{u}, \mathbf{v} \in \mathbb{R}^m$ the standard inner product $(\mathbf{u}, \mathbf{v}) = \mathbf{u}^T \mathbf{v}$ has this property. Consider the action of the operator $\mathbf{A} \in \mathbb{R}^{m \times m}$ on the two terms. If $(\mathbf{A}\mathbf{u}, \mathbf{v}) = (\mathbf{u}, \mathbf{A}\mathbf{v})$ the operator \mathbf{A} is said to be symmetric. For the inner product

$$(\mathbf{A}\mathbf{u}, \mathbf{v}) = (\mathbf{A}\mathbf{u})^T \mathbf{v} = \mathbf{u}^T \mathbf{A}^T \mathbf{v}, (\mathbf{u}, \mathbf{A}\mathbf{v}) = \mathbf{u}^T \mathbf{A}\mathbf{v},$$

and the two expressions are equal if $\mathbf{A} = \mathbf{A}^T$. The same approach extends to the $d = 1$ diffusion operator $L = -\partial_{xx}$ using the scalar product

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

Applying integration by parts

$$(\partial_{xx} f, g) = \int_0^1 f''(x) g(x) dx = [f' g]_{x=0}^{x=1} - \int_0^1 f'(x) g'(x) dx = [f' g]_{x=0}^{x=1} - [f g']_{x=0}^{x=1} + \int_0^1 f(x) g''(x) dx = b + (f, \partial_{xx} g).$$

For homogeneous boundary conditions $f(0) = f(1) = g(0) = g(1) = 0$, the symmetry condition is satisfied. Note that symmetry involves both the operator and the boundary conditions of the problem. For $d = 2$ the scalar product

$$(u, v) = \int_{\Omega} u(x, y) v(x, y) dx dy,$$

is defined on the unit square $\Omega = [0, 1] \times [0, 1]$, and for homogeneous Dirichlet boundary conditions two applications of Green's formula leads to $(\nabla^2 u, v) = (u, \nabla^2 v)$, and the Laplace operator is symmetric.

For $d = 1$, $x_j = jh$, $h = 1/(m+1)$, the eigenvalues ν_k of $\mathbf{A} = \text{diag}([1 \ -2 \ 1])$ are inferred from those of the $-\partial_{xx}$ operator with homogeneous boundary conditions at $x = 0$, $x = 1$

$$-\partial_{xx} \sin(\kappa \pi x) = (\kappa \pi)^2 \sin(\kappa \pi x), \kappa \in \mathbb{N}.$$

Positing that an eigenvector \mathbf{q} of \mathbf{A} is the discretization of the continuum eigenfunction leads to

$$q_j = \sin(\kappa \pi x_j) = \sin(\kappa \pi j h) = \sin\left(\frac{j \kappa \pi}{m+1}\right),$$

hypothesis that is verified by the calculation of component j of $\mathbf{A}\mathbf{q}$

$$(\mathbf{A}\mathbf{q})_j = -\sin\left[\frac{(j-1)\kappa\pi}{m+1}\right] + 2\sin\left[\frac{j\kappa\pi}{m+1}\right] - \sin\left[\frac{(j+1)\kappa\pi}{m+1}\right] = 2\left[1 - \cos\left(\frac{\kappa\pi}{m+1}\right)\right] \sin\left(\frac{j\kappa\pi}{m+1}\right),$$

thereby obtaining the eigenvalue

$$\nu_\kappa = 2\left[1 - \cos\left(\frac{\kappa\pi}{m+1}\right)\right] = 4\sin^2\left[\frac{\kappa\pi h}{2}\right],$$

which recovers the analytical eigenvalue in the $h \rightarrow 0$ limit

$$\lim_{h \rightarrow 0} \frac{\nu_\kappa}{h^2} = (\kappa\pi)^2.$$

For Jacobi $\mathbf{B} = \mathbf{D}^{-1}$, so the eigenvalues of \mathbf{BA} are

$$\lambda_\kappa = 2\sin^2\left[\frac{\kappa\pi h}{2}\right], \kappa = 1, 2, \dots, m.$$

The eigenvalues of $\mathbf{M} = \mathbf{I} - \mathbf{BA}$ are therefore

$$\mu_\kappa = 1 - 2\sin^2\left[\frac{\kappa\pi h}{2}\right] = \cos(\kappa\pi h)$$

Replacing $h = 1/(m+1)$ the largest eigenvalue is obtained for $\kappa = 1$

$$\mu_{\max} = \mu_1 = \cos\left(\frac{\pi}{m+1}\right).$$

For large m , $\mu_1 \approx 1$, and slow convergence is expected as verified in the numerical experiment from Fig. 2.

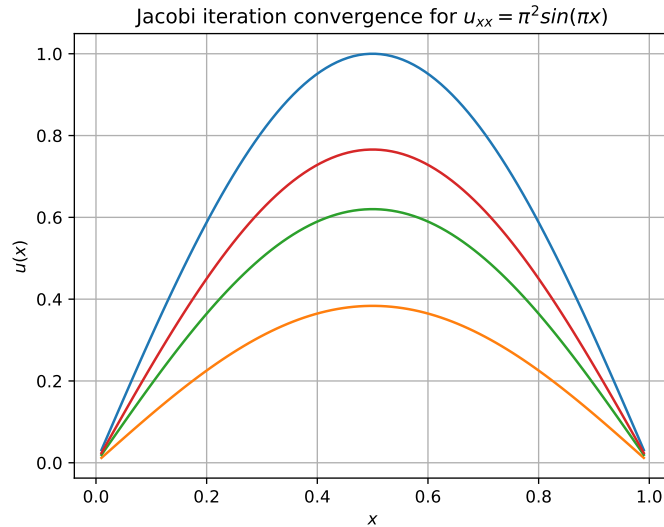


Figure 2. Convergence of Jacobi iteration. Blue: exact solution. Orange, green, red: iterates after 1000, 2000, 3000 iterations.

Jacobi iteration and its variants have limited practical utility by comparison to other iterative procedures due to slow convergence, but the concept of operator splitting has wide applicability. A more consequential example of operator splitting is to consider the advection diffusion equation

$$q_t + \mathbf{v} \cdot \nabla q = \alpha \Delta q,$$

which can be interpreted as stating that the time evolution of q is due to the effect of a diffusion operator $\mathcal{A} = \alpha \Delta$ and an advection operator $\mathcal{B} = -\mathbf{v} \cdot \nabla$

$$q_t = (\mathcal{A} + \mathcal{B})q.$$

Suppose both operators are discretized leading to matrices \mathbf{A}, \mathbf{B} and the discrete system

$$\mathbf{q}_t = (\mathbf{A} + \mathbf{B})\mathbf{q},$$

with initial condition $\mathbf{q}(\mathbf{x}, t=0) = \mathbf{q}_0$. Advancing the solution by a time step Δt can be written as

$$\mathbf{q}(t + \Delta t) = e^{\Delta t(\mathbf{A} + \mathbf{B})} \mathbf{q}(t),$$

and can be separated into two stages

$$\mathbf{q}^{(l+1)} = e^{\Delta t \mathbf{A}} e^{\Delta t \mathbf{B}} \mathbf{q}^{(l)}.$$

The quantity $\tilde{\mathbf{q}} = e^{\Delta t \mathbf{B}} \mathbf{q}^{(l)}$ captures advection effects and $\mathbf{q}^{(l+1)} = e^{\Delta t \mathbf{A}} \tilde{\mathbf{q}}$ is the diffusion correction. Since advection and diffusion are markedly different physical effects describing resolved versus unresolved transport, it can be expected that matrices \mathbf{A}, \mathbf{B} have different properties that require specific solution procedures.

By contrast, the Jacobi iteration splitting $\mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{U}$ does not separate physical effects and simply is suggested by the sparsity of \mathbf{A} and computational efficiency per iteration. For example, the forward Gauss-Seidel iteration with $\mathbf{B} = (\mathbf{D} + \mathbf{L})^{-1}$ leads to

$$\mathbf{u}_{k+1} = \mathbf{u}_k + (\mathbf{D} + \mathbf{L})^{-1} [\mathbf{c} - (\mathbf{L} + \mathbf{D} + \mathbf{U}) \mathbf{u}_k] = (\mathbf{D} + \mathbf{L})^{-1} (\mathbf{c} - \mathbf{U} \mathbf{u}_k) \Rightarrow (\mathbf{D} + \mathbf{L}) \mathbf{u}_{k+1} = \mathbf{c} - \mathbf{U} \mathbf{u}_k. \quad (14)$$

The matrix $\mathbf{D} + \mathbf{L}$ is (non-strictly) lower-triangular and (14) is easily solved by forward substitution. The implementation is very simple to express while preserving two-dimensional indexing.

Algorithm (Componentwise forward Gauss-Seidel)

```

for  $i = 1:m_x$ 
  for  $j = 1:m_y$ 
     $u(i, j) = [c(i, j) + u(i+1, j) + u(i-1, j) + u(i, j+1) + u(i, j-1)] / 4$ 
  end
end

```

In the above implementation a single memory space is used for \mathbf{u} with new values taking place of the old, leading to just four floating point additions and one multiplication per grid point. Since the algorithm essentially takes the current average of neighboring values, it is also known as a *relaxation* method, smoothing out spatial variations in \mathbf{u} .

Though such simple implementation is desirable, the non-physical splitting $\mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{U}$ and associated slow convergence usually outweighs ease of coding effort and suggests looking for alternative approaches. The only scenarios where such simply implemented iterations find practical use is in parallel execution and as a preliminary modification of the system prior to use of some other algorithm.