

Numerical approaches to solving PDE's

1. A general framework for numerical solution of PDE's

Now that we have an idea of the differential equations of interest in applications, let us turn to the problem of finding solutions. Analytical techniques such as separation of variables or Fourier analysis are very useful but for a limited class of problems, generally linear PDE's on domains of simple geometry. Most practical interest arises from non-linear PDE's over domains of complicated shape. The fundamental problem facing us is to determine an approximate function \tilde{q} that approximates the solution of the PDE of interest. A large number of methods have been devised to solve this problem numerically and we shall study individual methods extensively in later chapters. It is instructive to see that basically all methods can be expressed in a general framework that allows comparison of the strengths and weaknesses of individual methods. Say we are faced with the following general problem:

(1)

PROBLEM 1. Find $q : \Omega \rightarrow R^n$, $q \in \mathcal{F}$ that satisfies $Lq = 0$ on Ω and $Bq = 0$ on $\partial\Omega$ where L, B are operators defined on the normed linear space \mathcal{F} .

Let us see how this abstract statement maps onto one of the typical PDE problems, the Neumann problem:

$$(1.1) \quad \begin{cases} \frac{\partial^2 q}{\partial x^2} + \frac{\partial^2 q}{\partial y^2} = \sigma(x, y, q), & (x, y) \in \Omega \\ \frac{\partial q}{\partial n}(x, y) = F(x, y), & (x, y) \in \partial\Omega \end{cases} .$$

We have

$$(1.2) \quad L = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \sigma(x, y, \cdot)$$

$$(1.3) \quad B = \frac{\partial}{\partial n} - F(x, y)$$

with $L : \Omega$, $B : \partial\Omega$ and \mathcal{F} some reasonable space of functions, for instance the space of functions continuous up to second order defined on Ω , $\mathcal{C}^{(2)}(\Omega)$.

The advantage of the abstract formulation is that it allows us to concentrate on the typical steps followed when building a numerical procedure without spending too much time on the individual application at hand. We assume that the exact solution q is difficult to find and therefore concentrate on constructing an approximation $\tilde{q} \cong q$. Applying the operators L, B to the approximation \tilde{q} leads to an error, called a *residual*

$$(1.4) \quad L\tilde{q} = r, \quad B\tilde{q} = s .$$

It is our objective that the residual be as small as possible. In order to quantify the error made, we need a mapping from the function space to which r belongs to a real number. Such a mapping is called a *functional* and the typical example is the norm functional, e.g. $\|r\|$ and $\|s\|$. In numerical approximations we often work in the space of p -integrable functions $\mathcal{F} = \mathcal{L}^{(p)}(\Omega)$ for which the norm is defined by

$$(1.5) \quad \|f\| = \left(\int |f(x)|^p dx \right)^{1/p} .$$

The norm is a good candidate for evaluating the quality of our approximation \tilde{q} since if $\|r\| = 0$ then we know that $r = 0$ and therefore $\tilde{q} = q$. Unfortunately, the presence of the absolute value operation limits the operations that can be carried out on the norm and it is typical to use other functionals in constructing numerical procedures. Instead of using the norm consider the functionals

$$(1.6) \quad I(r) = \int_{\Omega} r w dx, \quad J(s) = \int_{\partial\Omega} s w' dx .$$

Here we have introduced weight functions w, w' that assign differing importance to errors made in various parts of the domain of integration. In adopting $I(r), J(s)$ as measures of the quality of our approximation \tilde{q} we have abandoned the certainty of knowing that $\tilde{q} = q$ when $I(r) = 0$, since a zero value for $I(r)$ could be obtained by cancellation of positive and negative errors throughout the domain. Appropriately chosen weight functions alleviate this concern somewhat and this is tradeoff we accept for now with a view to simplicity of the ensuing algorithm. It is possible to eliminate this drawback using the square of the residual as we shall see later on.

To simplify the presentation let us concentrate on $I(r)$ only. The addition of boundary conditions is usually a straightforward matter. Now that we have established a means by which to quantify the quality of an approximation we have to decide on how to build the approximation itself. Since \mathcal{F} is a normed linear space we can express any element g of \mathcal{F} as a linear combination over a set of basis functions

$$(1.7) \quad g = \sum_{i=1}^{\infty} a_i l_i .$$

For example the set $\{1, \sin x, \cos x, \sin 2x, \cos 2x, \dots\}$ is a basis for the square-integrable functions defined on $[0, 2\pi]$, $\mathcal{L}^{(2)}([0, 2\pi])$, and the set $\{1, x, x^2, \dots\}$ is a basis for the infinitely differentiable functions defined over the reals $C^{(\infty)}(\mathbb{R})$. In practical computations we cannot use an infinite sum, so we construct our approximation \tilde{q} using only N terms

$$(1.8) \quad \tilde{q} = \sum_{i=1}^N c_i l_i .$$

Using (1.8) $I(r)$ becomes

$$(1.9) \quad I(r) = \int_{\Omega} L \left(\sum_{i=1}^N c_i l_i \right) w dx .$$

For arbitrary coefficients c_i we shall have $I(r) \neq 0$. Since we wish \tilde{q} to be a good approximation we can reasonably impose $I(r) = 0$ and thus obtain an equation to

be used in determining the coefficients c_i

$$(1.10) \quad \int_{\Omega} L \left(\sum_{i=1}^N c_i l_i \right) w \, dx = 0 .$$

Now this is just one equation but we have N unknown coefficients. However the weight function is arbitrary so obtaining as many equations as we need is easy: we just choose w from some set of functions $w \in \{w_1, w_2, \dots, w_N\}$ and obtain a system of N equations for N unknowns

$$(1.11) \quad \int_{\Omega} L \left(\sum_{i=1}^N c_i l_i \right) w_j \, dx = 0, \quad j = 1, 2, \dots, N .$$

Up to this point we have concentrated on the approximation of q itself. There still remains the question of how to approximate the presumably complex shape of Ω . Generally this is done by approximating Ω by a set of simple-shaped subdomains ω_k such that the measure ρ of the difference between the two sets goes to zero as we increase the number of subdomains

$$(1.12) \quad \rho \left(\Omega - \bigcup_{m=1}^M \omega_m \right) \rightarrow 0 .$$

Using this approximation of the domain Ω (1.11) becomes

$$(1.13) \quad \sum_{m=1}^M \int_{\omega_m} L \left(\sum_{i=1}^N c_i l_i \right) w_j \, dx = 0, \quad j = 1, 2, \dots, N .$$

This is known as a *weighted residual formulation* and a large number of methods for numerically solving PDE's can be thus expressed. We shall turn to some examples shortly, but let us summarize the basic aspects:

- (1) A function space from which we construct approximations is chosen along with a subset of a basis of this space $\{l_1, l_2, \dots, l_N\}$;
- (2) A set of weight functions $\{w_1, w_2, \dots, w_N\}$ is chosen;
- (3) A discretization of the domain $\{\omega_1, \omega_2, \dots, \omega_M\}$ is chosen.

2. Basic numerical methods

2.1. Finite difference methods.

2.1.1. *Finite difference derivation.* In a finite difference method (FDM) the derivatives appearing in an ODE or PDE are approximated using finite differences. For example the IVP

$$(2.1) \quad \begin{cases} q' = f(t, q) \\ q(t=0) = q_0 \end{cases}$$

can be solved over the domain $[0, T]$ by finite differences using the following procedure. We construct an approximation \tilde{q} by a set of point values $Q^n = \tilde{q}(t^n)$ with $t^n = nk$, $n = 0, 1, \dots, N$ and k a step size $k = T/N$ and assume that \tilde{q} varies linearly between the point values. From the point values we can construct myriad approximations of the value of the derivative of \tilde{q} , for example

$$\tilde{q}'(t^n) \cong \frac{Q^{n+1} - Q^{n-1}}{2k},$$

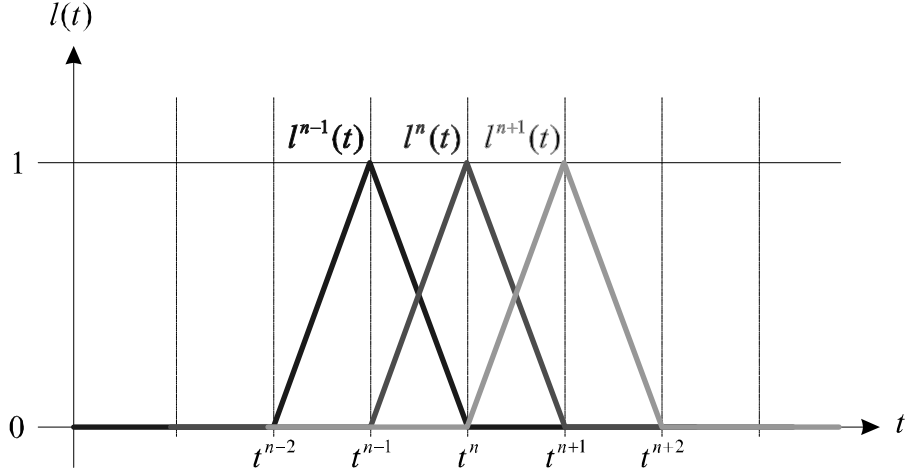


FIGURE 1. The piecewise linear form functions.

and obtain practical algorithms to solve (2.1) such as

$$(2.2) \quad Q^{n+1} = Q^{n-1} + 2k f(t^n, Q^n) ,$$

known as the *midpoint rule*.

2.1.2. *Weighted residual derivation.* We shall analyze (2.1) and similar algorithms extensively later on, but let us now see how the same method can be obtained via the weighted residual formulation and what insights we can thereby gain. We have chosen \tilde{q} as being a piecewise linear approximation defined at the points $t_n = nk$. In the general language of the weighted residual formulation we have $\Omega = [0, T]$, $\omega_m = [t^{m-1}, t^m]$. A basis for the piecewise linear functions defined on this partition of Ω is given by

$$(2.3) \quad l_n(t) = \begin{cases} 0 & t < t^{n-1} \\ \frac{1}{k} (t - t^{n-1}) & t^{n-1} \leq t < t^n \\ \frac{1}{k} (t^{n+1} - t) & t^n \leq t < t^{n+1} \\ 0 & t^{n+1} \leq t \end{cases} .$$

See Fig. (1).

Any piecewise linear function over $\{\omega_m\}$ can be defined as a linear combination of $\{l_n\}$, for instance

$$(2.4) \quad \tilde{q}(t) = \sum_{i=0}^N c_i l_i(t) .$$

It is apparent from (2.3) that

$$(2.5) \quad l_m(t^n) = \delta_{mn} = \begin{cases} 1 & \text{if } m = n \\ 0 & \text{if } m \neq n \end{cases} ,$$

so imposing the conditions $\tilde{q}(t^n) = Q^n$ leads to

$$(2.6) \quad \tilde{q}(t^n) = \sum_{i=0}^N c_i l_i(t^n) = \sum_{i=0}^N c_i \delta_{in} = c_n = Q^n ,$$

i.e. the coefficients of the expansion (2.4) are the nodal values Q^n . A similar expansion is made to approximate the values of the r.h.s. term

$$(2.7) \quad f(t, q) \cong \tilde{f}(t, q) = \sum_{i=0}^N F^i l_i(t)$$

with $F^n = \tilde{f}(t^n, Q^n)$.

We must now choose appropriate weight functions. Since the approximation we are building depends only on nodal values, a reasonable choice would be a set of weight functions that give importance to the residual obtained at the nodes. Such a set is given by

$$(2.8) \quad w_j = \delta(t - t^j)$$

where $\delta(t - t^j)$ is the Dirac delta functional centered on t^j defined by its integral property

$$(2.9) \quad \int_0^T g(t) \delta(t - t^j) dt = g(t^j)$$

for any function $g(t)$.

Having chosen the function space on which to base our approximation \tilde{q} , a basis in this space $\{l_m(t)\}$, the weight functions $\{\delta(t - t^j)\}$ and a discretization of the domain $[0, T]$ we can work through the weighted residual formulation (1.13) to obtain

$$(2.10) \quad \sum_{m=1}^M \int_{t^{m-1}}^{t^m} \left[\frac{d}{dt} \left(\sum_{n=1}^N Q^n l_n(t) \right) - \left(\sum_{n=1}^N F^n l_n(t) \right) \right] \delta(t - t^j) dt = 0 .$$

It is easiest to use the properties of the Dirac- δ function to work through the above expression. We have

$$(2.11) \quad \int_0^T \frac{d}{dt} \left(\sum_{n=1}^N Q^n \right) \delta(t - t^j) dt = \int_0^T \left(\sum_{n=1}^N F^n \right) \delta(t - t^j) dt$$

$$(2.12) \quad \sum_{n=1}^N Q^n \int_0^T l'_n(t) \delta(t - t^j) dt = \sum_{n=1}^N F^n \int_0^T \delta(t - t^j) dt$$

$$(2.13) \quad \sum_{n=1}^N Q^n l'_n(t_j) = \sum_{n=1}^N F^n l_n(t^j) = \sum_{n=1}^N F^n \delta_{nj} = F^j$$

We should be careful in evaluating $l'_n(t_j)$ since $l_n(t)$ is not differentiable at the nodal points $t = t^n$. If we interpret $l'_n(t_j)$ in principal value as the average of the limits to the left and the right we have

$$(2.14) \quad l'_n(t_j) = \begin{cases} -\frac{1}{2k} & n = j - 1 \\ 0 & n \neq j \pm 1 \\ \frac{1}{2k} & n = j + 1 \end{cases} .$$

Eq. (2.13) becomes

$$(2.15) \quad \frac{Q^{j+1} - Q^{j-1}}{2k} = F^j$$

which is exactly the same expression we had obtained previously, Eq. (2.2).

2.1.3. *Comparison of the two derivations.* It is satisfying to see a method derived in two ways, but an immediate question is what is to be gained by more complicated weighted residual procedure in comparison to the straightforward finite difference derivation. Generally the benefit arises in theoretical considerations of the behavior of the method. For instance let us consider the following theorem from approximation theory.

THEOREM 1. *Let \mathcal{V} be a metric linear space with a metric induced by the scalar product (\cdot, \cdot) on \mathcal{V} and let \mathcal{S} be a subspace of \mathcal{V} . Let $v \in \mathcal{V}$ and $u \in \mathcal{S}$. If $v - u$ is orthogonal to any $w \in \mathcal{S}$, then u is the best approximation of v within \mathcal{S} .*

PROOF. Let $d(u, v)$ be the distance induced by the scalar product between u and v . Ask whether any other $w \in \mathcal{S}$ gives a smaller distance to v

$$(2.16) \quad [d(v, w)]^2 = (v - w, v - w) = (v - u + u - w, v - u + u - w)$$

$$(2.17) \quad = (v - u, v - u) + 2(u - w, v - u) + (u - w, u - w)$$

$$(2.18) \quad = \|v - u\|^2 + \|u - w\|^2 + 2(u - w, v - u) .$$

Since $v - u$ is orthogonal to any element in \mathcal{S} it is orthogonal to $v - u$ so $(u - w, v - u) = 0$ and

$$(2.19) \quad [d(v, w)]^2 = \|v - u\|^2 + \|u - w\|^2 \geq \|v - u\|^2 = [d(v, u)]^2$$

and we conclude that u is the best approximation of v within \mathcal{S} . \square

We can apply such theorems to weighted residual derivations to infer the behavior of the numerical approximation. Applied to the above example, \mathcal{V} would be the space of differentiable functions to which q belongs. \mathcal{S} would be the space of piecewise continuous functions where we defined our approximation \tilde{q} . The best approximation we could obtain would be orthogonal to the complement \mathcal{S} of within \mathcal{V} . This subspace would contain functions that are not expressible as an expansion along the set (2.3), for instance functions that vary more rapidly than the time step chosen k . Predictions such as this are typically more difficult to obtain from the simpler finite difference derivation.

2.2. Finite volume methods.

2.3. Finite element methods.

2.4. Spectral methods.