# The Method of Weighted Residuals
# and the *Weak Statement*

Following are the "ingredients" that will comprise our finite element method "recipe."

**Step 1.** State the solution *approximation* as a series of assumed-known spatial functions multiplied by an unknown expansion coefficient set.

**Step 2.** From step 1, recognize that *two decisions* face us, i.e. what are suitable spatial functions to use, and how do we determine the unknown coefficients?

**Step 3.** Answer the latter question by defining a *weak statement* to determine the expansion coefficients, to ensure that the error in the approximation is an absolute minimum for *any choice* of spatial functions.

**Step 4.** Examine candidate spatial functions and insist that they have long-term versatility, hence select piecewise continuous polynomials forming the basis of *Lagrange interpolation*.

**Step 5.** Define a *discretization* of the solution spatial domain into finite elements; hence rearrange the piecewise polynomials into basis functions that yield easy evaluation of the weak statement integrals.

**Step 6.** Complete these integrals for the specified problem data; hence construct the *matrix statement* for determining the expansion coefficients defined in step 1.

**Step 7.** Apply the appropriate boundary data and *solve* the matrix statement, hence insert the previously unknown expansion coefficients into the stated approximation, step 1.

**Step 8.** Determine the accuracy of the approximation, quantize the error, and hence *validate solution quality*.

Let us now walk through the recipe with a problem we already know - one-dimensional steady state conduction in a thick slab with constant thermal conductivity and source. Recall the governing differential equation to be

$$\mathcal{L}(T) = -\frac{d}{dx}\left( k\frac{dT}{dx} \right) - s = 0 \qquad 0 < x < L \tag{1}$$

with associated boundary conditions

$$\ell(T) = k\frac{dT}{dn} - q_{in} = 0 \qquad x = 0 \tag{2a}$$

$$T = T_b \qquad x = L \tag{2b}$$

For this simple problem, we were able to obtain a closed form solution for the exact distribution of $T(x)$

$$T(x) = \frac{sL^2}{2k}\left[ 1 - \left(\frac{x}{L}\right)^2 \right] + \frac{q_{in}L}{k}\left[ 1 - \left(\frac{x}{L}\right) \right] + T_b \tag{3}$$

Heed that for "real" problems, closed form solutions do not exist (else why would we even worry about FEA!). Thus we shall assume the existences of some approximate solution which closely satisfies (1).

$$T(x) \approx T^N(x) = \sum_{\alpha=1}^{N} \Psi_\alpha(x) T_\alpha \tag{4}$$

Approximating the solution to (1) with the series expansion (4) generally will not coincide with the exact solution. Consequently, an error is created which is simply the difference between $T(x)$ and $T^N(x)$. Denoting this error as $e^N(x)$, then

$$e^N(x) = T(x) - T^N(x) \tag{5}$$

clearly state the connection between the approximation and its error. The analyst would certainly seek to make $e^N(x)$ as small as possible, and hence minimize the difference between $T(x)$ and $T^N(x)$. However, the error is a function of x and its minimization requires a general approach. The error must be "measured" in some appropriate manner. One manner is to substitute (5) into (1), which yields a differential equation $\mathcal{L}(e^N)$ for the error

$$\mathcal{L}(e^N) = \mathcal{L}(T) - \mathcal{L}(T^N) = -\mathcal{L}(T^N) \tag{6}$$

since (1) states that $\mathcal{L}(T) = 0$. Equation (6) shows then one connection between error and the approximation (4). Hence

$$\mathcal{L}(T^N) = -\frac{d}{dx}\left( k\frac{dT^N}{dx} \right) - s \neq 0 \tag{7}$$

becomes a precise measure of the distribution of the error in $T^N(x)$ over the domain of (1). It makes sense to require that (7) be as close to zero as possible. However, it is not reasonable to try to force the error to be zero everywhere, since this amounts to determining the exact distribution of $T(x)$! Instead, the pragmatic approach is to require the measure of the approximation error (7) to vanish in an overall integrated sense. The corresponding mathematical statement is that the weighted residual of (7), denoted the *weak statement* WS, must vanish:

$$WS^N = \int_\Omega \Phi_\beta(x) L(T^N)dx \equiv 0 \qquad for \quad 1 \leq \beta \leq N \tag{8}$$

In (8), $\Omega$ denotes the problem statement domain and the weight function set $\Phi_\beta(x)$ is absolutely arbitrary at this point. The decision on specific $\Phi_\beta(x)$ determines the type of weighted residual technique chosen. For the *Galerkin criterion*, each weight function is made identical to the corresponding trial function $\Psi_\alpha$ of the approximation (4), hence

$$\Phi_\beta(x) = \Psi_\alpha(x) \qquad for \quad 1 \leq \alpha, \beta \leq N \tag{9}$$

Selection of the Galerkin criterion definition (9) for (8) exactly recovers the classical Rayleigh-Ritz approximate variational formulation for a linear field problem, e.g. (1). Further, the Galerkin choice is particularly advantageous when solving a (linear or non-linear) problem, since (8) guarantees that the associated *approximation error is minimum* since it is orthogonal to the trial function set $\Psi_\beta$.

The *Galerkin weak statement* (8,9), is thus elegantly simple.

$$GWS^N = \int_0^L \Psi_\beta(x)\left[ -\frac{d}{dx}\left( k\frac{dT^N(x)}{dx} \right) - s \right]dx \equiv 0 \qquad for \quad 1 \leq \beta \leq N \tag{10}$$

The second derivative term in (10) can be reduced one order courtesy integration by parts yielding the Galerkin *symmetric weak statement*

$$GWS^N = \int_0^L \frac{d\Psi_\beta}{dx} k\frac{dT^N}{dx} dx - \int_0^L \Psi_\beta s dx - \Psi_\beta k\frac{dT^N}{dx}\bigg|_0^L = 0 \qquad for \quad 1 \leq \beta \leq N \tag{11}$$

Note that the differentiability requirement on $\Psi_\beta$ has been balanced as it is contained within $T^N$. Substituting (4), the explicit form for (11) is

$$GWS^N = \sum_{\alpha=1}^{N}\left(\int_0^L \frac{d\Psi_\beta}{dx}k\frac{d\Psi_\alpha}{dx}T_\alpha dx\right) - \int_0^L \Psi_\beta s\,dx - \Psi_\beta k\frac{dT^N}{dx}\bigg|_0^L = 0 \qquad for \quad 1 \le \beta \le N \qquad (12)$$

Note that $-k\,dT^N/dx = q_{in}$ at $x=0$, as required by the flux boundary condition (2a), is incorporated *exactly* into (12) for *any* approximation $T^N$, confirming that the finite element method handles all flux boundary conditions in a natural way.

From our investigation of interpolation theory, we found that we could conveniently approximate a function on an global level as the assembly of element level approximations

$$Q^N(x) = \sum_{\alpha=1}^{N}\Psi_\alpha(x)Q_\alpha = \sum_{e=1}^{M}\{N_k(\zeta_i)\}^T\{Q\}_e \qquad (13)$$

Thus, we shall obtain our approximate solution by next discretizing our problem domain into elements and substituting (13) into (12).

$$\begin{aligned}GWS^h &= \sum_{e=1}^{M}GSW_e^h \\ &= \sum_{e=1}^{M}\left(\int_0^{le}\frac{d\{N_k\}}{dx}k\frac{d\{N_k\}^T}{dx}d\bar{x}\{Q\}_e - \int_0^{le}\{N_k\}s\,d\bar{x} - k\frac{dT^N}{dx}\{\delta_k\} = \{0_k\}\right)\end{aligned} \qquad (14)$$

We have added the superscript h to GWS to denote the introduction of the discretization $\Omega^h$. The first two terms in (14) involve integrals only over the generic element domain. The last term is the boundary flux at each end of the domain $\Omega$ which will be replaced directly by the boundary conditions and do not involve any integral over $\Omega_e$

We must now select a finite element basis and perform the required element level derivatives and integrals. We shall complete this exercise with the linear basis and then return to investigate quadratic and cubic. Recall

$$\{N_{k=1}(\zeta_i)\} = \begin{Bmatrix}\zeta_1 \\ \zeta_2\end{Bmatrix}, \qquad \zeta_1(\bar{x}) \equiv 1 - \frac{\bar{x}}{l_e} \quad \zeta_2(\bar{x}) \equiv \frac{\bar{x}}{l_e}, \qquad \bar{x} = x - X_L \qquad (15)$$

Looking at the first integral term of (14), we need to differentiate $\{N_k\} = \{N_1\}$. Employing the chain rule

$$\frac{d\{N_1(\delta_i)\}}{dx} = \frac{1}{l_e}\begin{Bmatrix}-1 \\ 1\end{Bmatrix} \qquad (16)$$

Since the thermal conductivity is constant, the first integral in (14) is then evaluated and written compactly in matrix form as

$$\int_{\Omega_e} \frac{d\{N_1\}}{dx} k \frac{d\{N_1\}^T}{dx} d\overline{x}\{Q\}_e = k\int_0^{l_e} \frac{1}{l_e}\begin{Bmatrix} -1 \\ 1 \end{Bmatrix} \frac{1}{l_e}\{-1 \quad 1\} d\overline{x}\{Q\}_e$$

$$= \frac{k}{l_e^2}\begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}\int_0^{l_e} d\overline{x}\{Q\}_e \qquad (17)$$

$$= \frac{k}{l_e}\begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}\{Q\}_e$$

In the same manner, the constant source term is directly evaluated, yielding

$$\int_{\Omega_e} \{N_1\}s\,d\overline{x} = s\int_0^{l_e} \begin{Bmatrix} 1 - \dfrac{\overline{x}}{l_e} \\ \dfrac{\overline{x}}{l_e} \end{Bmatrix} d\overline{x}$$

$$= \frac{sl_e}{2}\begin{Bmatrix} 1 \\ 1 \end{Bmatrix} \qquad (18)$$

Substituting (17) and (18) into (14) gives the finite element matrix form for the discretized Galerkin weak statement using the linear basis function

$$\boxed{GWS^h = \sum_{e=1}^{M} GSW_e^h = \sum_{e=1}^{M}\left(\frac{k}{l_e}\begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}\{Q\}_e - \frac{sl_e}{2}\begin{Bmatrix} 1 \\ 1 \end{Bmatrix} - k\frac{dT^N}{dx}\begin{Bmatrix} \delta_{e1} \\ \delta_{eM} \end{Bmatrix}\right) = \begin{Bmatrix} 0 \\ 0 \end{Bmatrix}} \qquad (19)$$

**Homework (suggested)**

1.   Complete the integration for (18)