

# Finite-Elements Method <sup>2</sup>

January 29, 2014

---

<sup>2</sup>From *Applied Numerical Analysis* Gerald-Wheatley (2004), Chapter 9.

Finite-element methods (FEM) are based on some mathematical physics techniques and the most fundamental of them is the so-called **Rayleigh-Ritz method** which is used for the solution of boundary value problems.

Two other methods which are more appropriate for the implementation of the FEM will be discussed, these are **the collocation method** and **the Galerkin method**.

# Rayleigh-Ritz method

In the Rayleigh-Ritz (RR) method we solve a boundary-value problem by approximating the solution with a linear approximation of basis functions. The method is based on a part of mathematics called **calculus of variations**. In this method we try to minimize a special class of functions called **functionals**.

The usual form for functional in problems with one variable is

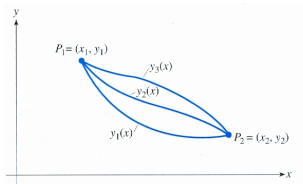
$$I[y] = \int_a^b F(x, y, y') dx \quad (1)$$

The argument of  $I[y]$  is not a simple variable but a function  $y = y(x)$ . The value of  $I[y]$  is varying with  $y(x)$ , but for fixed  $y(x)$  is a scalar quantity.

\*\* We seek the  $y(x)$  that **minimizes**  $I[y]$ .

# Rayleigh-Ritz method: an example I

Let's find the function  $y(x)$  that minimizes the distance between two points. Although, we know the answer, i.e. that it is a straight line, we will pretend that we don't, and that we will choose among the set of curves  $y_i(x)$  as in the figure.



The functional is the integral of the distance along any of these curves:

$$I[y] = \int_{x_1}^{x_2} \sqrt{dx^2 + dy^2} = \int_{x_1}^{x_2} \sqrt{1 + (dy/dx)^2} dx \quad (2)$$

To minimize  $I[y]$ , we set its derivatives to zero.

Each curve must pass through the points  $(x_1, y_1)$  and  $(x_2, y_2)$ .

In addition, for the optimal trajectory, the **Euler-Lagrange equation** must be satisfied:

$$\frac{d}{dx} \left[ \frac{\partial}{\partial y'} F(x, y, y') \right] - \frac{\partial}{\partial y} F(x, y, y') = 0 \quad (3)$$

For the functional of equation (2) we get:

$$\begin{aligned} F(x, y, y') &= (1 + (y')^2)^{1/2} \\ F_y &= 0 \\ F_{y'} &= y' (1 + (y')^2)^{-1/2} \end{aligned} \quad (4)$$

Finally, from the Euler-Lagrange equation (3) we get

$$\frac{d}{dx} \left( \frac{\partial F}{\partial y'} \right) = \frac{d}{dx} \left( \frac{y'}{\sqrt{1 + (y')^2}} \right) = \frac{\partial F}{\partial y} = 0 \quad (5)$$

From which it follows that :

$$\frac{y'}{\sqrt{1 + (y')^2}} = c \quad \rightarrow \quad y' = \sqrt{\frac{c^2}{1 - c^2}} = b \quad \rightarrow \quad y = bx + a \quad (6)$$

# Rayleigh-Ritz method: an example II

Let's try a more complicated case: Apply the RR method to the 2nd order boundary value problem

$$y'' + Q(x)y - F(x) = 0, \quad y(a) = y_0, \quad y(b) = y_n \quad (7)$$

The above boundary condition are known as **Dirichlet conditions**.  
The functional that corresponds to eqn (7) is:

$$I[u] = \int_a^b [(u')^2 - Qu^2 + 2Fu] dx \quad (8)$$

**Note I:** when the above functional is optimized (via Euler-Lagrange equation) leads to the original equation (7).

**Note II:** we now have only 1st order instead of 2nd order derivatives.

- If we know the solution to the ODE, substituting it for  $u$  in eqn (8) will make  $I[u]$  a minimum.
- If the solution is not known, we may try to approximate it by some arbitrary function and see whether we can minimize the functional by a suitable choice of the parameters of the approximations.
- **The Rayleigh-Ritz method is based on this idea.**

We let  $u(x)$ , which is the approximation to  $y(x)$  (the exact solution), be a sum:

$$u(x) = c_0 v_0(x) + c_1 v_1(x) + \cdots + c_n v_n(x) = \sum_{i=0}^n c_i v_i(x) \quad (9)$$

There are two conditions on the trial functions  $v_i(x)$  :

- They must be chosen such that  $u(x)$  meets the boundary conditions
- The individual  $v_i(x)$  should be linearly independent

The  $v_i(x)$  and  $c_i$  are to be chosen to make  $u(x)$  a good approximation to the true solution of eqn (7).

Since we have no prior knowledge of the true function  $y(x)$  we have no chance to guess the  $v_i(x)$  that will provide a solution to closely resemble  $y(x)$ .

Thus we go for the usual choice that is the use of polynomials! and we will try to find a way to get the values of the  $c_i$ .

These requirements can be fulfilled by using the functional of eqn (8). If we substitute  $u(x)$  as defined by equation (9) into the functional of eqn (8) we get:

$$I(c_0, c_1, \dots, c_n) = \int_a^b \left[ \left( \frac{d}{dx} \sum c_i v_i \right)^2 - Q \left( \sum c_i v_i \right)^2 + 2F \sum c_i v_i \right] dx$$

Thus  $I$  is a function of the unknown  $c_i$ . To minimize  $I$  we take the partial derivatives with respect to each unknown  $c_i$  and set zero.

$$\frac{\partial I}{\partial c_i} = 2 \int_a^b u' \frac{\partial u'}{\partial c_i} dx - \int_a^b 2Qu \frac{\partial u}{\partial c_i} dx + \int_a^b 2F \frac{\partial u}{\partial c_i} dx \quad (10)$$

Thus we led to a system of  $n$  equations to solve. This will define the  $u(x)$  of equation (9).



**EXAMPLE:** Solve the ODE  $y'' + y - 3x^2 = 0$  with the boundary conditions  $(0,0)$  and  $(2,3.5)$ .

We will use polynomials up to 3rd order, we can define  $u(x)$  as:

$$u(x) = \frac{7}{4}x + c_2x(x-2) + c_3x^2(x-2) \quad (11)$$

note that the functions  $v$  are linearly independent and that  $u(x)$  satisfies the boundary conditions.

The following terms are needed for the evaluation of equations (10):

$$\begin{aligned} u' &= \frac{7}{4} + 2c_2(x-1) + c_3(3x^2 - 4x) \\ \frac{\partial u'}{\partial c_2} &= 2x - 2 \quad \text{and} \quad \frac{\partial u'}{\partial c_3} = 3x^2 - 4x \\ \frac{\partial u}{\partial c_2} &= x(x-2) \quad \text{and} \quad \frac{\partial u}{\partial c_3} = x^2(x-2) \end{aligned} \quad (12)$$

By substituting the above equations into equation (10) we get two equations for the two constants  $c_2$  and  $c_3$ .

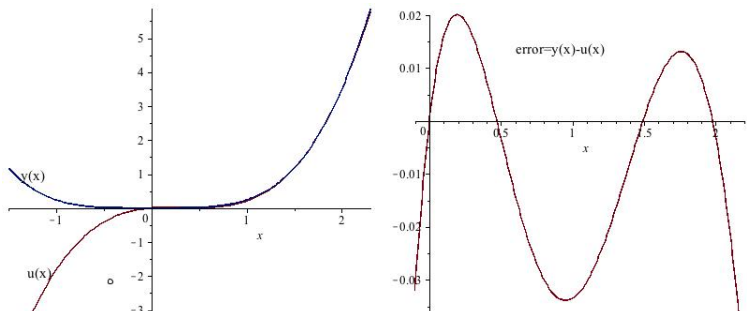
The results which come from trivial integrations are:

$$\frac{\partial I}{\partial c_2} : \frac{16}{5}c_2 + \frac{16}{5}c_3 = \frac{74}{15}$$
$$\frac{\partial I}{\partial c_3} : \frac{16}{5}c_2 + \frac{128}{21}c_3 = \frac{36}{15}$$

From their solution we get the values of  $c_2$  and  $c_3$  and we come to the approximate solution:

$$u(x) = \frac{119}{152}x^3 - \frac{46}{57}x^2 + \frac{53}{228}x \quad (13)$$

**NOTE:** The exact solution is:  $y(x) = 6 \cos x + 3(x^2 - 2)$ .



# The Collocation Method

The **collocation method** uses an alternative method in approximating the solution  $y(x)$  of the BVP defined in (7).

We define the residual function  $R(x)$  as:

$$R(x) = y'' + Qy - F \quad (14)$$

We approximate again  $y(x)$  with  $u(x)$  equal to a sum of trial functions (linearly independent polynomials) as with the RR method, and we try to make  $R(x) = 0$  by suitable choice of the coefficients.

Of course this is not possible to be achieved everywhere in the interval and thus we may choose arbitrarily to make  $R(x) = 0$  at a number of points inside the interval. The number of interior points should be the same as the number of the unknowns.

# The Collocation Method : Example

We will solve the previous example using the collocation method.

We take again

$$u(x) = \frac{7}{4}x + c_2x(x - 2) + c_3x^2(x - 2) \quad (15)$$

The residual is defined as:

$$R(x) = u'' + u - 3x^2 \quad (16)$$

and when we differentiate  $u(x)$  we get:

$$R(x) = 2c_2 + 2c_3(x - 2) + 4c_3x + (7/4)x + c_2x(x - 2) + c_3x^2(x - 2) - 3x^2$$

Since we have 2 unknown constants we will use 2 points in the interval  $[0, 2]$ , e.g  $x = 0.7$  and  $x = 1.3$ .

Then by setting  $R(x) = 0$  for these two choices we get:

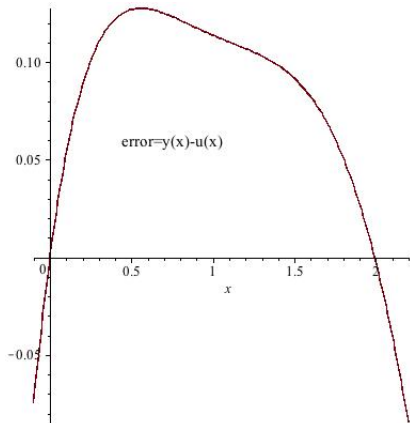
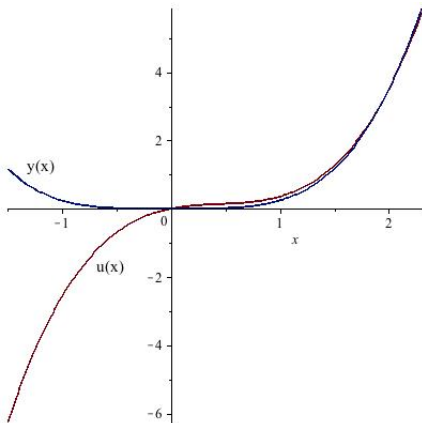
$$x = 0.7 : \quad 1090c_2 - 437c_3 - 245 = 0$$

$$x = 1.3 : \quad 1090c_2 + 2617c_3 - 2795 = 0$$

# The Collocation Method : Example

and by solving for  $c_2$  and  $c_3$  we get

$$u(x) = \frac{425}{509}x^3 - \frac{61607}{55481}x^2 + \frac{140023}{221924}x \quad (17)$$



# The Galerkin Method

- This method can be considered as a variation of the collocation method i.e. is a "residual method" that use the function  $R(x)$  defined in (14). The difference is that here we multiply with weighting functions  $W_i(x)$  which can be chosen in many ways.
- Galerkin showed that the individual trial functions  $v_i(x)$  used in (9) are a good choice.
- Once we have selected the  $v_i(x)$  from eqn (9) we compute the unknown coefficients by setting the integral of the weighted residual to zero:

$$\int_a^b W_i(x)R(x)dx = 0, i = 0, 1, \dots, n \quad \text{where} \quad W_i(x) = v_i(x) \quad (18)$$

- Notice that using the Dirac delta functions for the  $W_i(x)$  we reduce to the collocation method.

# The Galerkin Method: Example

We will solve the previous example using the Galerkin method.

We take again

$$u(x) = \frac{7}{4}x + c_2x(x-2) + c_3x^2(x-2) \quad (19)$$

so that  $v_2 = x(x-2)$  and  $v_3 = x^2(x-2)$ .

The residual is

$$R(x) = u'' + u - 3x^2 \quad (20)$$

and after substituting  $u''$  and  $u$  we get

$$R(x) = 2c_2 + c_3(6x-4) + 7x/4 + c_2x(x-2) + c_3x^2(x-2) - 3x^2 \quad (21)$$

We carry out two integrations:

$$\text{If } v_2 \rightarrow W_i : \int_0^2 [x(x-2)] R(x) dx = 0$$

$$\text{If } v_3 \rightarrow W_i : \int_0^2 [x^2(x-2)] R(x) dx = 0$$

These intergrations give two equations for the  $c_j$ :

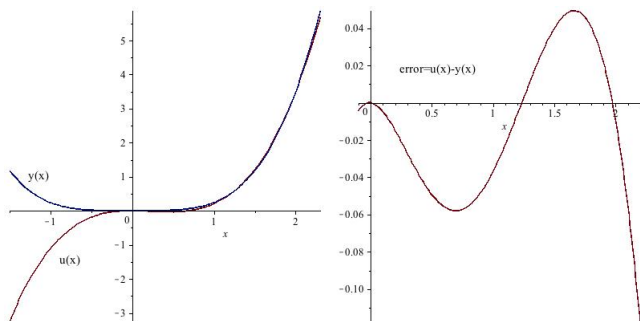
$$24c_2 + 24c_3 - 37 = 0$$

$$21c_2 + 40c_3 - 45 = 0$$

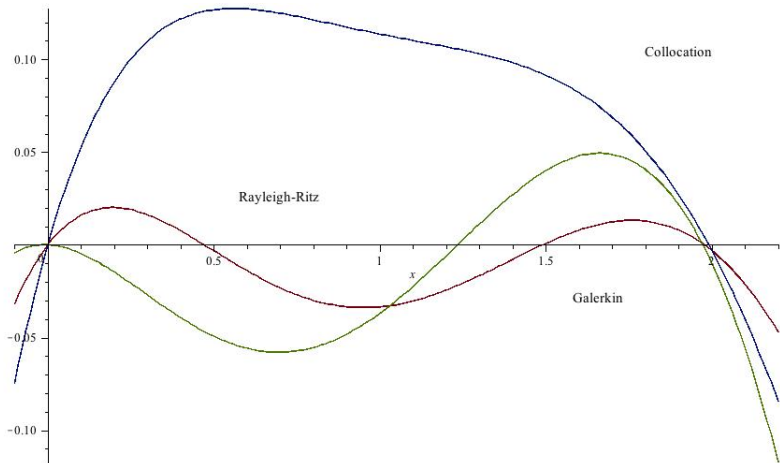
and by solving for  $c_2$  and  $c_3$  we get

$$u(x) = \frac{101}{152}x^3 - \frac{103}{228}x^2 - \frac{1}{128}x \quad (22)$$

**COMMENT:** Although the RR method is more accurate than the others the Galerkin method is much easier to be implemented since we don't need to find the variational form.







# Finite Elements for ODEs

In order to improve the accuracy and also to be able to treat longer intervals  $[a, b]$ , we follow the steps below:

- Subdivide  $[a = x_0, b = x_n]$  into  $n$  subintervals, the **elements**, that join at the **nodes**  $x_1, x_2, \dots, x_{n-1}$
- Apply the Galerkin method to each element separately to interpolate between the end point nodal values  $u(x_{i-1})$  and  $u(x_i)$
- Use a low-degree polynomial for  $u(x)$ , e.g. even a 1st degree can do the work (higher order polynomials are better but too complicated to be implemented)
- When Galerkin's method is applied to element  $(i)$  we get a pair of eqns with unknowns the nodal values at the ends of the element  $(i)$ , the  $c_j$ .
- If we do it for each element we end up with a system of equations involving all the nodal values
- The equations adjusted to take into account the boundary conditions and the solution is an approximation to  $y(x)$  at the nodes; intermediate values can be taken by interpolation.

We will describe step by step the procedure for solving the following boundary value problem:

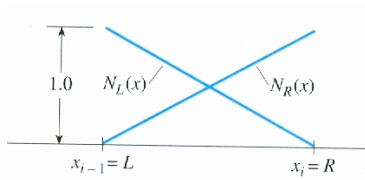
$$y'' + Q(x)y = F(x) \quad \text{with BC at } x = a \text{ and } x = b \quad (23)$$

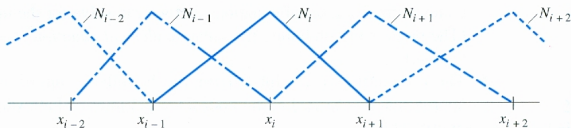
**STEP 1** : Subdivide the interval  $[a, b]$  into  $n$  elements, and for a specific element between  $x_{i-1}$  and  $x_i$  name the left node as  $L$  and the right as  $R$ .

**STEP 2** : Write  $u(x)$  for the element  $(i)$ :

$$u(x) = c_L N_L + c_R N_R = c_L \frac{x - R}{L - R} + c_R \frac{x - L}{R - L} = c_L \frac{x - R}{-h_i} + c_R \frac{x - L}{-h_i} \quad (24)$$

Notice that the  $N$ 's are 1st-degree Lagrange polynomials usually called **hat functions**. The following figure sketches the  $N_L$  and  $N_R$  for the  $(i)$  element.





**STEP 3** : Apply the Galerkin method to element ( $i$ ), the residual is

$$R(x) = y'' + Qy - F = u'' + Qu - F \quad (25)$$

The Galerkin method sets the integral of  $R$  weighted with each of the  $N$

$$\int_L^R N_L R(x) dx = 0$$

$$\int_L^R N_R R(x) dx = 0$$

These lead to the integrals (26)

$$\int_L^R u'' N_L dx + \int_L^R Qu N_L dx - \int_L^R FN_L dx = 0$$

$$\int_L^R u'' N_R dx + \int_L^R Qu N_R dx - \int_L^R FN_R dx = 0$$

**STEP 4** : Transform the previous equations by applying integration by parts (**how?**) to the first integral. In the 2nd integral take  $Q$  out from the integrant as  $Q_{av}$  (an average value within the element). We do the same also for the 3rd integral.

$$\int_L^R u' N_L' dx - Q_{av} \int_L^R u N_L dx + F_{av} \int_L^R N_L dx - N_L u'|_{x=R} + N_L u'|_{x=L} = 0 \quad (27)$$

Remember that  $N_L = 1$  at  $L$  and  $0$  at  $R$ . Thus the equation can be written as:

$$\int_L^R u' N_L' dx - Q_{av} \int_L^R u N_L dx + F_{av} \int_L^R N_L dx - u'|_{x=L} = 0 \quad (28)$$

and similarly the other equation:

$$\int_L^R u' N_R' dx - Q_{av} \int_L^R u N_R dx + F_{av} \int_L^R N_R dx + u'|_{x=R} = 0 \quad (29)$$

**STEP 5** : Substitute in the previous integrals  $u$ ,  $u'$ ,  $N'_L$  and  $N'_R$  and carry out the integrations.

$$\int_L^R u' N'_L dx = \dots = (c_L - c_R) / h_i \quad (30)$$

$$-Q_{av} \int_L^R u N_L dx = \dots = -Q_{av} h_i (2c_L - c_R) / 6 \quad (31)$$

$$F_{av} \int_L^R N_L dx = \dots = F_{av} h_i / 2 \quad (32)$$

and we do the same for equation (29)

**STEP 6** : Substitute the results of the previous steps into equations (28) and (29) and by rearrangement we get 2 equations of the unknowns  $c_L$  and  $c_R$ :

$$\left(\frac{1}{h_i} - \frac{Q_{av} h_i}{3}\right) c_L - \left(\frac{1}{h_i} + \frac{Q_{av} h_i}{6}\right) c_R = -\frac{F_{av} h_i}{2} - u'|_{x=L} \quad (33)$$

$$\left(\frac{-1}{h_i} - \frac{Q_{av} h_i}{3}\right) c_L + \left(\frac{1}{h_i} - \frac{Q_{av} h_i}{6}\right) c_R = -\frac{F_{av} h_i}{2} + u'|_{x=L} \quad (34)$$

These are the so called **element equations**.

\*\* We do the same for each element to get  $n$  such pairs. \*\*

## STEP 7 :

- We assemble all the element equations to form a system of linear equations for the problem.
- Notice that the point  $R$  in the element  $(i)$  is the same as the point  $L$  in the element  $(i + 1)$ .
- Renumber the  $c$  as  $c_0, c_1, \dots, c_n$
- Notice that the gradient  $u'$  must be the same on either side of the elements i.e.  $u'_{x=R}$  in the element  $(i)$  equals  $u'_{x=L}$  in element  $(i + 1)$ . Thus these terms will cancel when we assemble except in the first and last equations.
- The results is a system of  $n + 1$  equations of the form

$$\mathbf{K}\vec{c} = \vec{b} \quad (35)$$

The matrix  $\mathbf{K}$  contains combination of the quantities  $Q_{av,i}$  and  $h_i$ .

The matrix  $\mathbf{K}$  is tridiagonal

The vector  $\vec{c}$  contains the  $c_0, c_1, \dots, c_n$

The vector  $\vec{b}$  contains combinations of  $F_{av,i}$  and  $h_i$ .

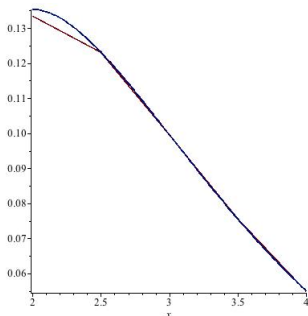


**STEP 8** : Adjust the set of the equations for the boundary conditions which can be either **Dirichlet** or **Neumann**.

**STEP 9** : Solve the system and find the  $c_0, c_1, \dots, c_n$ .

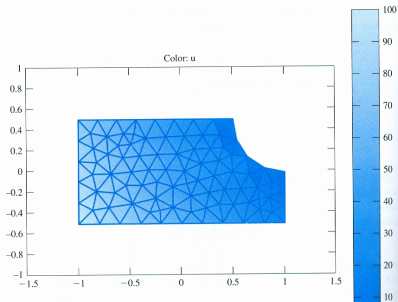
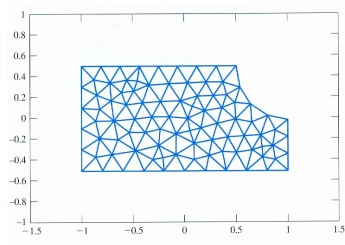
\*\*\* Then relax :-)

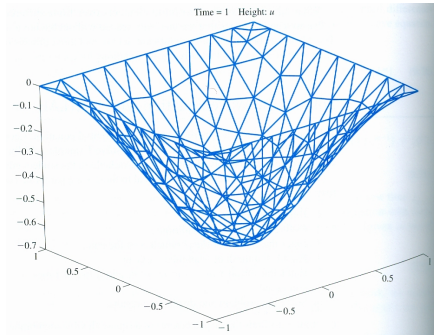
**EXAMPLE:** Solve  $y'' - (x + 1)y = e^{-x}(x^2 - x + 2)$  with the Neumann conditions  $y'(2) = 0$  and  $y'(4) = -0.036631$ . (The exact solution is  $y(x) = e^{-x}(x - 2)$ ).



# Finite Elements for PDEs

Quite complicated.





Solving the wave equation

# Finite Elements for PDEs

