# 數值分析 MA-3021

### Chapter 6. Numerical Ordinary Differential Equations

- §6.1 Introduction
- §6.2 Taylor-Series Method
- §6.3 Euler's Method
- §6.4 The Runge-Kutta Method
- §6.5 Collocation Method
- §6.6 Finite Difference Method
- §6.7 Finite Element Method

#### Definition

A differential equation is a mathematical equation that relates some unknown function with its derivatives. A differential equation is called an *ordinary differential equation* (ODE) if it contains an unknown function of one independent variable and its derivatives. A differential equation is called a *partial differential equation* (PDE) if it contains unknown multi-variable functions and their partial derivatives.

#### Definition

The *order* of a differential equation is the order of the highest-order derivatives presented in the equation. A differential equation of order 1 is called first order, order 2 second order, etc.

#### Definition

A differential equation is a mathematical equation that relates some unknown function with its derivatives. A differential equation is called an *ordinary differential equation* (ODE) if it contains an unknown function of one independent variable and its derivatives. A differential equation is called a *partial differential equation* (PDE) if it contains unknown multi-variable functions and their partial derivatives.

#### Definition

The **order** of a differential equation is the order of the highest-order derivatives presented in the equation. A differential equation of order 1 is called first order, order 2 second order, etc.

#### Definition

A differential equation is a mathematical equation that relates some unknown function with its derivatives. A differential equation is called an *ordinary differential equation* (ODE) if it contains an unknown function of one independent variable and its derivatives. A differential equation is called a *partial differential equation* (PDE) if it contains unknown multi-variable functions and their partial derivatives.

#### Definition

The *order* of a differential equation is the order of the highest-order derivatives presented in the equation. A differential equation of order 1 is called first order, order 2 second order, etc.

#### Remark:

1 In general, an *n*-th order (scalar) ODE has the form

$$F(t, y, y', \cdots, y^{(n)}) = 0,$$

where  $\frac{\partial F}{\partial y^{(n)}} \neq 0$ . By the implicit function theorem,

$$y^{(n)} = \varphi(t, y, y', \dots, y^{(n-1)}) \tag{*}$$

for some function  $\varphi$  (locally).

② Assume that y satisfies  $(\star)$ . Let  $\mathbf{x} = (y, y', y'', \dots, y^{(n-1)})^{\top}$ . Then  $\mathbf{x}$  satisfies

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x})$$

where the vector-valued function f is given by

$$f(t, \mathbf{x}) = f(t, x_1, \cdots, x_n) = \begin{bmatrix} x_2 \\ \vdots \\ x_n \\ \varphi(t, \mathbf{x}) \end{bmatrix}.$$

#### Remark:

• In general, an *n*-th order (scalar) ODE has the form

$$F(t, y, y', \cdots, y^{(n)}) = 0,$$

where  $\frac{\partial F}{\partial y^{(n)}} \neq 0$ . By the **implicit function theorem**,  $y^{(n)} = \varphi\left(t, y, y', \cdots, y^{(n-1)}\right)$ 

$$y^{(n)} = \varphi(t, y, y', \dots, y^{(n-1)}) \tag{*}$$

for some function  $\varphi$  (locally).

**3** Assume that y satisfies  $(\star)$ . Let  $\mathbf{x} = (y, y', y'', \dots, y^{(n-1)})^{\top}$ . Then x satisfies

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x}),$$

where the vector-valued function  $\mathbf{f}$  is given by

$$\mathbf{f}(t,\mathbf{x}) = \mathbf{f}(t,x_1,\cdots,x_n) = \begin{bmatrix} x_2 \\ \vdots \\ x_n \\ \varphi(t,\mathbf{x}) \end{bmatrix}.$$

### Remark (cont'd):

① It is also possible to consider n-th order **vector-valued** ODE. For example, let r(t) denote the position of a planet (of mass m) moving around the sun (of mass M) which locates at the origin. Then Newton's second law of motion implies that

$$m\mathbf{r}''(t) = -\frac{GMm}{\|\mathbf{r}\|^3}\mathbf{r}(t)$$
.

In general, when considering this kind of equation, we assume that it can be written as

$$\mathbf{y}^{(n)} = \boldsymbol{\varphi}(t, \mathbf{y}, \mathbf{y}', \cdots, \mathbf{y}^{(n-1)}).$$

Then the same procedure (by letting  $\mathbf{x} = (\mathbf{y}, \mathbf{y}', \dots, \mathbf{y}^{(n-1)})^{\top}$  and find an equation that  $\mathbf{x}$  satisfies) shows that the equation above reduces to an first order ODE

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x})$$
.



### Remark (cont'd):

3 It is also possible to consider n-th order **vector-valued** ODE. For example, let r(t) denote the position of a planet (of mass m) moving around the sun (of mass M) which locates at the origin. Then Newton's second law of motion implies that

$$m\mathbf{r}''(t) = -\frac{GMm}{\|\mathbf{r}\|^3}\mathbf{r}(t)$$
.

In general, when considering this kind of equation, we assume that it can be written as

$$\mathbf{y}^{(n)} = \varphi(t, \mathbf{y}, \mathbf{y}', \cdots, \mathbf{y}^{(n-1)}).$$

Then the same procedure (by letting  $\mathbf{x} = (\mathbf{y}, \mathbf{y}', \cdots, \mathbf{y}^{(n-1)})^{\top}$  and find an equation that  $\mathbf{x}$  satisfies) shows that the equation above reduces to an first order ODE

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x})$$
.



Initial-value problem (IVP): find x(t) such that

$$\begin{cases} & \textbf{\textit{x}}'(t) = \textbf{\textit{f}}\big(t,\textbf{\textit{x}}(t)\big), & \text{(ordinary differential equations)} \\ & \textbf{\textit{x}}(t_0) = \textbf{\textit{x}}_0, & \text{(initial condition)} \end{cases}$$

where  $\mathbf{f}(t, \mathbf{x})$ ,  $t_0 \in \mathbb{R}^1$ ,  $\mathbf{x}_0 \in \mathbb{R}^n$  are given.

#### Example

$$x'(t) = x(t)\tan(t+3), \qquad x(-3) = 1.$$

The analytic solution of this IVP is  $x(t) = \sec(t+3)$ . The solution is valid only for  $-\frac{\pi}{2} < t+3 < \frac{\pi}{2}$ .

#### Example

$$x'(t) = x, \qquad x(0) = 1.$$

Try  $x(t) = ce^{rt} \Rightarrow cre^{rt} = ce^{rt} \Rightarrow r = 1$ ,  $x = ce^t$  general solution Use  $x(0) = 1 \Rightarrow x = e^t$  particular solution



**Existence and Uniqueness of Solution**: Not all IVPs have a solution. Even if there exists a solution, the solution may not be unique.

#### Example

There is no solution to the initial value problem

$$\exp\left(x'(t)\right) = 0, \qquad x(0) = 0$$

even if complex-valued solutions are allowed.

#### Example

The initial value problem

$$x'(t) = 3x(t)^{\frac{2}{3}}, \qquad x(0) = 0$$

has infinitely many solutions. In fact, the function

$$x_c(t) = \begin{cases} 0 & \text{if } t \le c \\ (t-c)^3 & \text{if } t > c \end{cases}$$

is a solution to the given ODE for all  $c \ge 0$ .

#### Theorem (Existence and Uniqueness of Solution)

Consider the initial value problem

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x}), \qquad \mathbf{x}(t_0) = \mathbf{x}_0.$$

If **f** and the first partial derivatives of **f** with respect to all its variables, possibly except t, are continuous functions in some rectangular domain R that contains the point  $(t_0, \mathbf{x}_0)$  in the interior, then the initial value problem has a unique solution  $\varphi(t)$  in some interval  $I = (t_0 - h, t_0 + h)$  for some positive number h.

#### Remark

- ① If f is k-times continuously differentiable, then x is (k+1)-times continuously differentiable.
- ② The length of the time interval of existence "usually" is inverse proportional to the maximum of  $\|\mathbf{f}\| + \|\mathbf{f}_{\mathbf{y}}\|$ .

### Theorem (Fundamental Theorem of ODE)

Consider the initial value problem

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x}), \qquad \mathbf{x}(t_0) = \mathbf{x}_0.$$

If  $\mathbf{f}$  and the first partial derivatives of  $\mathbf{f}$  with respect to all its variables, possibly except t, are continuous functions in some rectangular domain R that contains the point  $(t_0, \mathbf{x}_0)$  in the interior, then the initial value problem has a unique solution  $\varphi(t)$  in some interval  $I = (t_0 - h, t_0 + h)$  for some positive number h.

#### Remark

- ① If f is k-times continuously differentiable, then x is (k+1)-times continuously differentiable.
- ② The length of the time interval of existence "usually" is inverse proportional to the maximum of  $\|\mathbf{f}\| + \|\mathbf{f}_{\mathbf{x}}\|$ .

### Theorem (Fundamental Theorem of ODE)

Consider the initial value problem

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x}), \qquad \mathbf{x}(t_0) = \mathbf{x}_0.$$

If  $\mathbf{f}$  and the first partial derivatives of  $\mathbf{f}$  with respect to all its variables, possibly except t, are continuous functions in some rectangular domain R that contains the point  $(t_0, \mathbf{x}_0)$  in the interior, then the initial value problem has a unique solution  $\varphi(t)$  in some interval  $I = (t_0 - h, t_0 + h)$  for some positive number h.

#### Remark:

- If f is k-times continuously differentiable, then x is (k+1)-times continuously differentiable.
- ② The length of the time interval of existence "usually" is inverse proportional to the maximum of  $\|\mathbf{f}\| + \|\mathbf{f}_{\mathbf{x}}\|$ .



#### Remark:

③ Suppose that for a constant  $k \in (0,1)$  the "cube" Q centered at  $(t_0, \mathbf{x}_0)$  with width 2k is a subset of the rectangular domain R, and for some  $M \ge 1$ ,  $|\mathbf{f}(t, \mathbf{y})| + |\mathbf{f}_{\mathbf{x}}(t, \mathbf{x})| \le M$  for all  $(t, \mathbf{x}) \in Q$ . Then the solution  $\mathbf{x}$  to the initial value problem

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x}), \qquad \mathbf{x}(t_0) = \mathbf{x}_0.$$

exists at least in the interval  $\left[t_0 - \frac{k}{M}, t_0 + \frac{k}{M}\right]$ .

### Definition (Informal definition)

A numerical method of solving the ODE  $\mathbf{x}' = \mathbf{f}(t, \mathbf{x})$  with  $\mathbf{x}(t_0) = \mathbf{x}_0$ is an **iterative scheme** which, when the **step size** h > 0 is given, generates a unique sequence of vectors  $\{x_1, \dots, x_N\}$  (for some N which in general depends on h) such that some function  $\varphi$  which interpolates the data  $(t_0, \mathbf{x}_0), (t_1, \mathbf{x}_1), \cdots, (t_n, \mathbf{x}_N)$  where,  $t_n = t_0 +$ *nh*, **resembles** the solution to x' = f(t, x) with initial condition  $\mathbf{x}(t_0) = \mathbf{x}_0$  in the time interval  $[t_0, t_N]$ . The function  $\varphi$  is called the **numerical solution** generated by this numerical method with step size h.

### Definition (Informal definition (cont'd))

A numerical method of solving the ODE  $\mathbf{x}' = \mathbf{f}(t, \mathbf{x})$  is called a kstep method if it requires  $\mathbf{x}_n, \mathbf{x}_{n+1}, \cdots, \mathbf{x}_{n+k-1}$  to determine  $\mathbf{x}_{n+k}$ for all  $n \in \{0, \cdots, N-k\}$ .

A numerical method of solving the ODE x' = f(t, x) is said to be **explicit** if it does not require "nonlinear procedures" to obtain some  $x_n$ 's, and is said to be **implicit** if it is not explicit.

**Remark**: A one-step explicit method is often (but not always) given in the form  $\mathbf{x}_{n+1} = \mathbf{x}_n + h\Phi(t_n, \mathbf{x}_n)$  for some function  $\Phi$ , while a k-step explicit method is often (but not always) given in the form

$$\mathbf{x}_{n+1} = \alpha_1 \mathbf{x}_n + \alpha_2 \mathbf{x}_{n-1} + \dots + \alpha_k \mathbf{x}_{n-k+1} + h \Big[ \beta_1 \mathbf{f}(t_n, \mathbf{x}_n) + \dots + \beta_k \mathbf{f}(t_{n-k+1}, \mathbf{x}_{n-k+1}) \Big].$$

#### Example

The forward Euler method of solving the ordinary differential equations  $\mathbf{y}' = \mathbf{f}(t, \mathbf{y})$  is an explicit one-step method given by

$$\mathbf{x}_n = \mathbf{x}_{n-1} + h\mathbf{f}(t_{n-1}, \mathbf{x}_{n-1}) \quad \forall n \in \{1, 2, \cdots, N\},$$

while the backward Euler method is an implicit one-step method given by

$$\mathbf{x}_n = \mathbf{x}_{n-1} + h\mathbf{f}(t_n, \mathbf{x}_n) \qquad \forall n \in \{1, 2, \cdots, N\}.$$

#### Example

The Runge-Kutta method is an explicit one-step method given by

$$\mathbf{x}_{n+1} = \mathbf{x}_n + h\left(\frac{\mathbf{k}_{n1} + 2\mathbf{k}_{n2} + 2\mathbf{k}_{n3} + \mathbf{k}_{n4}}{6}\right),$$

where

$$\mathbf{k}_{n1} = \mathbf{f}(t_n, \mathbf{x}_n), \qquad \mathbf{k}_{n2} = \mathbf{f}(t_n + \frac{1}{2}h, \mathbf{x}_n + \frac{1}{2}h\mathbf{k}_{n1}), \mathbf{k}_{n3} = \mathbf{f}(t_n + \frac{1}{2}h, \mathbf{x}_n + \frac{1}{2}h\mathbf{k}_{n2}), \quad \mathbf{k}_{n4} = \mathbf{f}(t_n + h, \mathbf{x}_n + h\mathbf{k}_{n3}).$$



#### Example

The forward Euler method of solving the ordinary differential equations  $\mathbf{y}' = \mathbf{f}(t, \mathbf{y})$  is an explicit one-step method given by

$$\mathbf{x}_n = \mathbf{x}_{n-1} + h\mathbf{f}(t_{n-1}, \mathbf{x}_{n-1}) \qquad \forall n \in \{1, 2, \cdots, N\},$$

while the backward Euler method is an implicit one-step method given by

$$\mathbf{x}_n = \mathbf{x}_{n-1} + h\mathbf{f}(t_n, \mathbf{x}_n) \qquad \forall n \in \{1, 2, \cdots, N\}.$$

#### Example

The Runge-Kutta method is an explicit one-step method given by

$$\mathbf{x}_{n+1} = \mathbf{x}_n + h\left(\frac{\mathbf{k}_{n1} + 2\mathbf{k}_{n2} + 2\mathbf{k}_{n3} + \mathbf{k}_{n4}}{6}\right),\,$$

where

$$\mathbf{k}_{n1} = \mathbf{f}(t_n, \mathbf{x}_n), \qquad \mathbf{k}_{n2} = \mathbf{f}(t_n + \frac{1}{2}h, \mathbf{x}_n + \frac{1}{2}h\mathbf{k}_{n1}), \mathbf{k}_{n3} = \mathbf{f}(t_n + \frac{1}{2}h, \mathbf{x}_n + \frac{1}{2}h\mathbf{k}_{n2}), \quad \mathbf{k}_{n4} = \mathbf{f}(t_n + h, \mathbf{x}_n + h\mathbf{k}_{n3}).$$



#### Example

The forward Euler method of solving the ordinary differential equations  $\mathbf{y}' = \mathbf{f}(t, \mathbf{y})$  is an explicit one-step method given by

$$\mathbf{x}_n = \mathbf{x}_{n-1} + h\mathbf{f}(t_{n-1}, \mathbf{x}_{n-1}) \quad \forall n \in \{1, 2, \cdots, N\},$$

while the backward Euler method is an implicit one-step method given by

$$\mathbf{x}_n = \mathbf{x}_{n-1} + h\mathbf{f}(t_n, \mathbf{x}_n) \qquad \forall n \in \{1, 2, \cdots, N\}.$$

#### Example

The Runge-Kutta method is an explicit one-step method given by

$$\mathbf{x}_{n+1} = \mathbf{x}_n + h\left(\frac{\mathbf{k}_{n1} + 2\mathbf{k}_{n2} + 2\mathbf{k}_{n3} + \mathbf{k}_{n4}}{6}\right),\,$$

where

$$\mathbf{k}_{n1} = \mathbf{f}(t_n, \mathbf{x}_n),$$
  $\mathbf{k}_{n2} = \mathbf{f}(t_n + \frac{1}{2}h, \mathbf{x}_n + \frac{1}{2}h\mathbf{k}_{n1}),$   $\mathbf{k}_{n3} = \mathbf{f}(t_n + \frac{1}{2}h, \mathbf{x}_n + \frac{1}{2}h\mathbf{k}_{n2}),$   $\mathbf{k}_{n4} = \mathbf{f}(t_n + h, \mathbf{x}_n + h\mathbf{k}_{n3}).$ 



There are three fundamental sources of error of a numerical solution:

- **1** The iterative scheme used to produce the sequence  $\{x_1, \dots, x_N\}$  is an approximate one. In other words, at each step the numerical method does not produce the correct value of the solution at the next time step. This relates to the local/global truncation error.
- ② The input data used in the iterative scheme are only approximations to the actual values of the solution at each  $t_k$ . For example, one should use  $\mathbf{x}(t_k)$  to generate  $\mathbf{x}_{k+1}$  but we are forced to start with  $\mathbf{x}_k$ . This relates to the global truncation error.
- The precision of calculations of the computer is finite. In other words, at each step only a finite number of digits can be retained. This relates to the round-off error (or machine error).

There are three fundamental sources of error of a numerical solution:

- The iterative scheme used to produce the sequence  $\{x_1, \cdots, x_N\}$  is an approximate one. In other words, at each step the numerical method does not produce the correct value of the solution at the next time step. This relates to the local/global truncation error.
- ② The input data used in the iterative scheme are only approximations to the actual values of the solution at each  $t_k$ . For example, one should use  $\mathbf{x}(t_k)$  to generate  $\mathbf{x}_{k+1}$  but we are forced to start with  $\mathbf{x}_k$ . This relates to the global truncation error.
- The precision of calculations of the computer is finite. In other words, at each step only a finite number of digits can be retained. This relates to the round-off error (or machine error).

There are three fundamental sources of error of a numerical solution:

- **1** The iterative scheme used to produce the sequence  $\{x_1, \dots, x_N\}$  is an approximate one. In other words, at each step the numerical method does not produce the correct value of the solution at the next time step. This relates to the local/global truncation error.
- ② The input data used in the iterative scheme are only approximations to the actual values of the solution at each  $t_k$ . For example, one should use  $\mathbf{x}(t_k)$  to generate  $\mathbf{x}_{k+1}$  but we are forced to start with  $\mathbf{x}_k$ . This relates to the global truncation error.
- The precision of calculations of the computer is finite. In other words, at each step only a finite number of digits can be retained. This relates to the round-off error (or machine error).

#### Definition

Let  $\varphi$  be a numerical solution obtained by a specific numerical method (with step size h>0 fixed) of solving ODE  $\mathbf{x}'=\mathbf{f}(t,\mathbf{x})$  with initial data  $\mathbf{x}(t_0)=\mathbf{x}_0$ . At each time step  $t_n$ ,

- the **global truncation error** (associated with this numerical method) is the number  $\mathbf{E}_n(h) = \mathbf{x}(t_n) \varphi(t_n)$ ;
- 2 the **local truncation error** (associated with this numerical method) is the number  $\tau_n(h) = x(t_{n+1}) x_{n+1}$ , where  $x(\cdot)$  is the exact solution and  $x_{n+1}$  is obtained according to the iterative scheme with  $x_i = x(t_i)$  for all  $j \in \{0, 1, \dots, n\}$ .
- **1** the **round-off error** or **machine error** (associated with this numerical method) is the number  $\mathbf{R}_n = \varphi(t_n) \mathbf{X}_n$ , where  $\mathbf{X}_n$  is the actual value computed from the numerical method.

#### Definition

Let  $\varphi$  be a numerical solution obtained by a specific numerical method (with step size h>0 fixed) of solving ODE  $\mathbf{x}'=\mathbf{f}(t,\mathbf{x})$  with initial data  $\mathbf{x}(t_0)=\mathbf{x}_0$ . At each time step  $t_n$ ,

- **1** the **global truncation error** (associated with this numerical method) is the number  $\mathbf{E}_n(h) = \mathbf{x}(t_n) \varphi(t_n)$ ;
- ② the **local truncation error** (associated with this numerical method) is the number  $\tau_n(h) = \mathbf{x}(t_{n+1}) \mathbf{x}_{n+1}$ , where  $\mathbf{x}(\cdot)$  is the exact solution and  $\mathbf{x}_{n+1}$  is obtained according to the iterative scheme with  $\mathbf{x}_i = \mathbf{x}(t_i)$  for all  $j \in \{0, 1, \dots, n\}$ .
- 3 the **round-off error** or **machine error** (associated with this numerical method) is the number  $\mathbf{R}_n = \boldsymbol{\varphi}(t_n) \mathbf{X}_n$ , where  $\mathbf{X}_n$  is the actual value computed from the numerical method.

#### Definition

Let  $\varphi$  be a numerical solution obtained by a specific numerical method (with step size h>0 fixed) of solving ODE  $\mathbf{x}'=\mathbf{f}(t,\mathbf{x})$  with initial data  $\mathbf{x}(t_0)=\mathbf{x}_0$ . At each time step  $t_n$ ,

- the **global truncation error** (associated with this numerical method) is the number  $\mathbf{E}_n(h) = \mathbf{x}(t_n) \varphi(t_n)$ ;
- ② the **local truncation error** (associated with this numerical method) is the number  $\tau_n(h) = \mathbf{x}(t_{n+1}) \mathbf{x}_{n+1}$ , where  $\mathbf{x}(\cdot)$  is the exact solution and  $\mathbf{x}_{n+1}$  is obtained according to the iterative scheme with  $\mathbf{x}_j = \mathbf{x}(t_j)$  for all  $j \in \{0, 1, \dots, n\}$ .
- **1** the **round-off error** or **machine error** (associated with this numerical method) is the number  $\mathbf{R}_n = \varphi(t_n) \mathbf{X}_n$ , where  $\mathbf{X}_n$  is the actual value computed from the numerical method.

In other words, the local truncation error measures the accuracy of the numerical method for each time step, while the global truncation error measure the errors accumulated from the beginning of this iterative scheme.

#### Definition

A numerical method is said to be *consistent* if

$$\lim_{h\to 0} \max_{0 \le n \le N-1} \left| \frac{\tau_n(h)}{h} \right| = 0,$$

where  $\tau_n(h)$  is the local truncation error associated with the numerical method with step size h.

A numerical method is said to have order k if for any sufficiently smooth solution of the initial value problem, the local truncation error is  $\mathcal{O}(h^{k+1})$ .

In other words, the local truncation error measures the accuracy of the numerical method for each time step, while the global truncation error measure the errors accumulated from the beginning of this iterative scheme.

#### Definition

A numerical method is said to be *consistent* if

$$\lim_{h\to 0} \max_{0\leqslant n\leqslant N-1} \left| \frac{\tau_n(h)}{h} \right| = 0,$$

where  $\tau_n(h)$  is the local truncation error associated with the numerical method with step size h.

A numerical method is said to **have order** k if for any sufficiently smooth solution of the initial value problem, the local truncation error is  $\mathcal{O}(h^{k+1})$ .

In other words, the local truncation error measures the accuracy of the numerical method for each time step, while the global truncation error measure the errors accumulated from the beginning of this iterative scheme.

#### Definition

A numerical method is said to be *consistent* if

$$\lim_{h\to 0} \max_{0\leqslant n\leqslant N-1} \left| \frac{\tau_n(h)}{h} \right| = 0,$$

where  $\tau_n(h)$  is the local truncation error associated with the numerical method with step size h.

A numerical method is said to **have order** k if for any sufficiently smooth solution of the initial value problem, the local truncation error is  $\mathcal{O}(h^{k+1})$ .

For a second order ODE

$$y'' + p(t)y' + q(t)y = g(t),$$

instead of imposing the initial condition  $y(t_0)=y_0$  and  $y'(t_0)=y_1$  sometimes the boundary condition  $y(\alpha)=y_0$  and  $y(\beta)=y_1$  are imposed and consider the two-point boundary value problem (BVP)

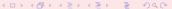
$$y'' + p(x)y' + q(x)y = g(x)$$
  $\forall x \in (\alpha, \beta), y(\alpha) = y_0, y(\beta) = y_1.$ 

Let 
$$z(x)=y(x)-rac{x-lpha}{eta-lpha}y_1-rac{x-eta}{lpha-eta}y_0$$
. Then  $z$  satisfies

$$z'' + p(x)z' + q(x)z = G(x) \quad \forall x \in (\alpha, \beta), \ z(\alpha) = z(\beta) = 0, \quad (\star)$$

where 
$$G(x) = g(x) - p(x)\frac{y_0 - y_1}{\alpha - \beta} - q(x)(\frac{x - \alpha}{\beta - \alpha}y_1 + \frac{x - \beta}{\alpha - \beta}y_0)$$
. (\*) is

called a homogeneous two-point boundary value problem



For a second order ODE

$$y'' + p(t)y' + q(t)y = g(t),$$

instead of imposing the initial condition  $y(t_0)=y_0$  and  $y'(t_0)=y_1$  sometimes the boundary condition  $y(\alpha)=y_0$  and  $y(\beta)=y_1$  are imposed and consider the two-point boundary value problem (BVP)

$$y'' + p(x)y' + q(x)y = g(x)$$
  $\forall x \in (\alpha, \beta), y(\alpha) = y_0, y(\beta) = y_1.$ 

Let 
$$z(x) = y(x) - \frac{x - \alpha}{\beta - \alpha} y_1 - \frac{x - \beta}{\alpha - \beta} y_0$$
. Then z satisfies

$$z'' + p(x)z' + q(x)z = G(x) \quad \forall x \in (\alpha, \beta), \ z(\alpha) = z(\beta) = 0, \quad (\star)$$

where 
$$G(x) = g(x) - p(x) \frac{y_0 - y_1}{\alpha - \beta} - q(x) \left( \frac{x - \alpha}{\beta - \alpha} y_1 + \frac{x - \beta}{\alpha - \beta} y_0 \right)$$
. (\*) is

called a homogeneous two-point boundary value problem.



Remark: Even though the IVP

$$y'' + p(t)y' + q(t)y = g(t),$$
  $y(t_0) = y_0,$   $y'(t_0) = y_1$ 

looks quite similar to the boundary value problem, they actually differ in some very important ways. For example, if p, q, g are continuous, the IVP above always has a unique solution, while the BVP might have no solution or infinitely many solutions:

- y'' + y = 0 with boundary condition  $y(0) = y(\pi) = 0$  has infinite many solutions  $y_c(x) = c \sin x$ .
- ②  $y'' + y = \sin x$  with boundary condition  $y(0) = y(\pi) = 0$  has no solution.

On the other hand, some BVPs has a unique solution. For example,

$$y'' + 2y = 0$$
,  $\forall x \in (0, \pi)$ ,  $y(0) = 1$ ,  $y(\pi) = 0$ 

has a unique solution  $y(x) = \cos \sqrt{2}x - \cot \sqrt{2}\pi \sin \sqrt{2}x$ .

←□▶←□▶←≣▶←≣▶
■

Remark: Even though the IVP

$$y'' + p(t)y' + q(t)y = g(t),$$
  $y(t_0) = y_0,$   $y'(t_0) = y_1$ 

looks quite similar to the boundary value problem, they actually differ in some very important ways. For example, if p, q, g are continuous, the IVP above always has a unique solution, while the BVP might have no solution or infinitely many solutions:

- y'' + y = 0 with boundary condition  $y(0) = y(\pi) = 0$  has infinite many solutions  $y_c(x) = c \sin x$ .
- ②  $y'' + y = \sin x$  with boundary condition  $y(0) = y(\pi) = 0$  has no solution.

On the other hand, some BVPs has a unique solution. For example,

$$y'' + 2y = 0$$
,  $\forall x \in (0, \pi)$ ,  $y(0) = 1$ ,  $y(\pi) = 0$ 

has a unique solution  $y(x) = \cos \sqrt{2}x - \cot \sqrt{2}\pi \sin \sqrt{2}x$ .

#### Theorem

Let  $\alpha, \beta$  be real numbers and  $\alpha < \beta$ . Suppose that the function f = f(t, y, p) is continuous on the set

$$D = \{(x, y, p) \mid x \in [\alpha, \beta], y, p \in \mathbb{R}\}$$

and the partial derivatives  $f_y$  and  $f_p$  are also continuous on D. If

- $f_y(t, y, p) > 0$  for all  $(t, y, p) \in D$ , and
- 2 there exists a constant M > 0 such that

$$|f_p(t,y,p)| \leqslant M \quad \forall (t,y,p) \in D,$$

then the boundary value problem

$$y'' = f(t, y, y') \quad \forall x \in (\alpha, \beta), \ y(\alpha) = y(\beta) = 0$$

has a unique solution.



#### Theorem

Let  $\alpha, \beta$  be real numbers and  $\alpha < \beta$ . Suppose that  $p : [\alpha, \beta] \to \mathbb{R}$  is continuously differentiable, and  $q : [\alpha, \beta] \to \mathbb{R}$  is continuous. Then

$$y'' + p(x)y' + q(x)y = g(x) \quad \forall x \in (\alpha, \beta), \ y(\alpha) = y(\beta) = 0$$

has a solution if and only if  $g: [\alpha, \beta] \to \mathbb{R}$  is integrable and

$$\int_{\alpha}^{\beta} g(x)\varphi(x) \, dx = 0$$

for all function  $\varphi$  satisfying

$$\varphi'' - p(x)\varphi' + (q(x) - p'(x))\varphi = 0 \quad \forall x \in (\alpha, \beta), \ \varphi(\alpha) = \varphi(\beta) = 0.$$

The solution is unique if the ODE y'' + p(x)y' + q(x)y = 0 with  $y(\alpha) = y(\beta) = 0$  has only trivial solution  $y \equiv 0$ .

## §6.2 Taylor-Series Method

Consider the ODE x' = f(t, x). Then the chain rule can be used to compute  $x^{(k)}(t)$ :

$$x''(t) = f_{t}(t,x) + f_{x}(t,x)x',$$

$$x'''(t) = f_{tt}(t,x) + 2f_{tx}(t,x)x' + f_{xx}(t,x)(x')^{2} + f_{x}(t,x)x'',$$

$$x^{(4)}(t) = f_{ttt}(t,x) + 3f_{ttx}(t,x)x' + 3f_{txx}(t,x)(x')^{2} + f_{xxx}(t,x)(x')^{3} + f_{tx}(t,x)x'' + f_{xx}(t,x)x''x'' + f_{x}(t,x)x'''$$

$$\vdots$$

thus Taylor's Theorem implies that

$$x(t+h) = x(t) + hf(t,x) + \frac{h^2}{2} \left[ f_t(t,x) + f_x(t,x)f(t,x) \right]$$

$$+ \frac{h^3}{3} \left[ f_{tt}(t,x) + 2f_{tx}(t,x)f(t,x) + f_{xx}(t,x)f(t,x)^2 + f_x(t,x) \left( f_t(t,x) + f_x(t,x)f(t,x) \right) \right] + \cdots$$

Consider the ODE x' = f(t, x). Then the chain rule can be used to compute  $x^{(k)}(t)$ :

$$x''(t) = f_{t}(t, x) + f_{x}(t, x)x',$$

$$x'''(t) = f_{tt}(t, x) + 2f_{tx}(t, x)x' + f_{xx}(t, x)(x')^{2} + f_{x}(t, x)x'',$$

$$x^{(4)}(t) = f_{ttt}(t, x) + 3f_{ttx}(t, x)x' + 3f_{txx}(t, x)(x')^{2} + f_{xxx}(t, x)(x')^{3} + f_{tx}(t, x)x'' + f_{xx}(t, x)x''x'' + f_{x}(t, x)x'''$$

$$\vdots$$

thus Taylor's Theorem implies that

$$x(t+h) = x(t) + hf(t,x) + \frac{h^2}{2} \left[ f_t(t,x) + f_x(t,x)f(t,x) \right]$$

$$+ \frac{h^3}{3} \left[ f_{tt}(t,x) + 2f_{tx}(t,x)f(t,x) + f_{xx}(t,x)f(t,x)^2 + f_x(t,x) \left( f_t(t,x) + f_x(t,x)f(t,x) \right) \right] + \cdots$$

Depending on the smoothness of f and how many terms one wants to keep (deleting high order terms), we can obtain various Taylor's Method:

Taylor's method of order one:

$$x_{n+1} = x_n + hf(t_n, x_n).$$

2 Taylor's method of order two:

$$x_{n+1} = x_n + hf(t_n, x_n) + \frac{h^2}{2} [f_t(t_n, x_n) + f_x(t_n, x_n)f(t_n, x_n)].$$

3 Taylor's method of order three:

$$x_{n+1} = x_n + hf(t_n, x_n) + \frac{h^2}{2} \left[ f_t(t_n, x_n) + f_x(t_n, x_n) f(t_n, x_n) \right]$$

$$+ \frac{h^3}{3} \left[ f_{tt}(t_n, x_n) + 2f_{tx}(t_n, x_n) f(t_n, x_n) + f_{xx}(t_n, x_n) f(t_n, x_n)^2 + f_x(t_n, x_n) \left( f_t(t_n, x_n) + f_x(t_n, x_n) f(t_n, x_n) \right) \right].$$

#### Example

We use a concrete example to illustrate the method. Consider the following IVP

$$x'(t) = \cos t - \sin x + t^2$$
,  $x(-1) = 3$ .

By the Fundamental Theorem of ODE, the solution x is infinitely many times differentiable. By the Taylor series for x, we have

$$x(t+h) = x(t) + hx'(t) + \frac{h^2}{2!}x''(t) + \frac{h^3}{3!}x'''(t) + \frac{h^4}{4!}x^{(4)}(t) + \mathcal{O}(h^5).$$

Since  $x'(t) = \cos t - \sin x + t^2$ ,

$$\begin{cases} x''(t) = -\sin t - (\cos x)x' + 2t, \\ x'''(t) = -\cos t + \sin x(x')^2 - (\cos x)x'' + 2, \\ x^{(4)}(t) = \sin t + (\cos x)(x'')^3 + 3(\sin x)x'x'' - (\cos x)x'''. \end{cases}$$

thus by truncating at  $h^4$ , the **local truncation error** for obtaining x(t+h) is  $\mathcal{O}(h^5)$ . Such a method is of order 4.

#### Example (cont'd)

Starting t=-1 with h=0.01, we can compute the solution in  $\left[-1,1\right]$  with 200 steps:

input 
$$M \leftarrow 200$$
,  $h \leftarrow 0.01$ ,  $t \leftarrow -1$ ,  $x \leftarrow 3$  output  $0, t, x$ 

$$\quad \text{for } \textit{k} = 1 \text{ to } \textit{M} \text{ do}$$

$$x' \leftarrow \cos t - \sin x + t^{2}$$

$$x'' \leftarrow -\sin t - (\cos x)x' + 2t$$

$$x''' \leftarrow -\cos t + \sin x(x')^{2} - (\cos x)x'' + 2$$

$$x^{(4)} \leftarrow \sin t + (\cos x)(x')^{3} + 3(\sin x)x'x'' - (\cos x)x'''$$

$$x \leftarrow x + h(x' + \frac{h}{2}(x'' + \frac{h}{3}(x''' + \frac{h}{4}x^{(4)})))$$

output k, t, x

 $t \leftarrow t + h$ 

end do

#### Example (cont'd)

1 The local truncation error can be estimated by looking at

$$E_n = \frac{1}{(n+1)!} h^{n+1} x^{(n+1)} (t + \theta h)$$
 for some  $\theta \in (0,1)$ .

Hence

$$E_4 = \frac{1}{5!} h^5 x^{(5)} (t + \theta h) \quad \theta \in (0, 1).$$

**2** Replace  $x^{(5)}(t+\theta h)$  by a simple finite-difference approximation

$$E_4 \approx \frac{1}{5!} h^5 \left( \frac{x^{(4)}(t+h) - x^{(4)}(t)}{h} \right) = \frac{h^4}{120} \left( x^{(4)}(t+h) - x^{(4)}(t) \right).$$

**3** Suppose that the local truncation error (LTE) is  $\mathcal{O}(h^{n+1})$ . The accumulation of all many LTEs gives rise the global truncation error (GTE):

$$GTE \approx \frac{T - t_0}{h} \mathcal{O}(h^{n+1}) = \mathcal{O}(h^n).$$

And we say the numerical method is of  $\mathcal{O}(h^n)$ .

#### Taylor-Series method for systems:

For each  $1 \le i \le n$ , by the Taylor Theorem

$$x_i(t+h) = x_i(t) + hx_i'(t) + \frac{h^2}{2}x_i''(t) + \dots + \frac{h^n}{n!}x_i^{(n)}(t) + \mathcal{O}(h^{n+1})$$

which implies the vector form

$$\mathbf{x}(t+h) = \mathbf{x}(t) + h\mathbf{x}'(t) + \frac{h^2}{2}\mathbf{x}''(t) + \dots + \frac{h^n}{n!}\mathbf{x}^{(n)}(t) + \mathcal{O}(h^{n+1})$$

Using  $\mathbf{x}' = \mathbf{f}(t, \mathbf{x})$ , we find that

$$\mathbf{x}(t+h) = \mathbf{x}(t) + h\mathbf{f}(t,\mathbf{x}(t)) + \frac{h^2}{2}\mathbf{x}''(t) + \dots + \frac{h^n}{n!}\mathbf{x}^{(n)}(t) + \mathcal{O}(h^{n+1}).$$

Since 
$$\mathbf{x}''(t) = \mathbf{f}_t(t, \mathbf{x}) + \sum_{k=1}^n \frac{\partial \mathbf{f}}{\partial x_k}(t, \mathbf{x}) x_k'(t) \equiv [\mathbf{f}_t + (\mathbf{f} \cdot \nabla) \mathbf{f}](t, \mathbf{x})$$
, we have

$$\mathbf{x}(t+h) = \mathbf{x}(t) + h\mathbf{f}(t,\mathbf{x}(t)) + \frac{h^2}{2} \left[ \mathbf{f}_t + (\mathbf{f} \cdot \nabla)\mathbf{f} \right](t,\mathbf{x}(t)) + \mathcal{O}(h^3)$$

#### Taylor-Series method for systems:

For each  $1 \le i \le n$ , by the Taylor Theorem

$$x_i(t+h) = x_i(t) + hx_i'(t) + \frac{h^2}{2}x_i''(t) + \dots + \frac{h^n}{n!}x_i^{(n)}(t) + \mathcal{O}(h^{n+1})$$

which implies the vector form

$$\mathbf{x}(t+h) = \mathbf{x}(t) + h\mathbf{x}'(t) + \frac{h^2}{2}\mathbf{x}''(t) + \dots + \frac{h^n}{n!}\mathbf{x}^{(n)}(t) + \mathcal{O}(h^{n+1})$$

Using  $\mathbf{x}' = \mathbf{f}(t, \mathbf{x})$ , we find that

$$\mathbf{x}(t+h) = \mathbf{x}(t) + h\mathbf{f}(t,\mathbf{x}(t)) + \frac{h^2}{2}\mathbf{x}''(t) + \dots + \frac{h^n}{n!}\mathbf{x}^{(n)}(t) + \mathcal{O}(h^{n+1}).$$

Since 
$$\mathbf{x}''(t) = \mathbf{f}_t(t, \mathbf{x}) + \sum_{k=1}^n \frac{\partial \mathbf{f}}{\partial x_k}(t, \mathbf{x}) x_k'(t) \equiv [\mathbf{f}_t + (\mathbf{f} \cdot \nabla) \mathbf{f}](t, \mathbf{x})$$
, we have

$$\mathbf{x}(t+h) = \mathbf{x}(t) + h\mathbf{f}(t,\mathbf{x}(t)) + \frac{h^2}{2} \left[ \mathbf{f}_t + (\mathbf{f} \cdot \nabla)\mathbf{f} \right](t,\mathbf{x}(t)) + \mathcal{O}(h^3)$$

#### Taylor-Series method for systems:

For each  $1 \le i \le n$ , by the Taylor Theorem

$$x_i(t+h) = x_i(t) + hx_i'(t) + \frac{h^2}{2}x_i''(t) + \dots + \frac{h^n}{n!}x_i^{(n)}(t) + \mathcal{O}(h^{n+1})$$

which implies the vector form

$$\mathbf{x}(t+h) = \mathbf{x}(t) + h\mathbf{x}'(t) + \frac{h^2}{2}\mathbf{x}''(t) + \dots + \frac{h^n}{n!}\mathbf{x}^{(n)}(t) + \mathcal{O}(h^{n+1})$$

Using  $\mathbf{x}' = \mathbf{f}(t, \mathbf{x})$ , we find that

$$\mathbf{x}(t+h) = \mathbf{x}(t) + h\mathbf{f}(t,\mathbf{x}(t)) + \frac{h^2}{2}\mathbf{x}''(t) + \dots + \frac{h^n}{n!}\mathbf{x}^{(n)}(t) + \mathcal{O}(h^{n+1}).$$

Since 
$$\mathbf{x}''(t) = \mathbf{f}_t(t, \mathbf{x}) + \sum_{k=1}^n \frac{\partial \mathbf{f}}{\partial x_k}(t, \mathbf{x}) x_k'(t) \equiv \left[\mathbf{f}_t + (\mathbf{f} \cdot \nabla)\mathbf{f}\right](t, \mathbf{x})$$
, we have

$$\mathbf{x}(t+h) = \mathbf{x}(t) + h\mathbf{f}(t,\mathbf{x}(t)) + \frac{h^2}{2} \big[ \mathbf{f}_t + (\mathbf{f} \cdot \nabla)\mathbf{f} \big](t,\mathbf{x}(t)) + \mathcal{O}(h^3)$$

#### **Disadvantages:**

- The method depends on repeated differentiation of the differential equation, unless we intend to use only the method of order
  - 1. For high order methods, f(t, x) must have partial derivatives of sufficient high order in the region where are solving the problem. Such an assumption is not necessary for the existence of a solution.
- The various derivatives formula need to be programmed.

- Conceptually simple.
- Potential for high precision. For example, if we get 20 derivatives of x(t), then the method is order 20 (that is, terms up to and including the one involving  $h^{20}$ ).

#### **Disadvantages:**

- The method depends on repeated differentiation of the differential equation, unless we intend to use only the method of order
   1. For high order methods, f(t, x) must have partial derivatives of sufficient high order in the region where are solving the problem. Such an assumption is not necessary for the existence of a solution.
- The various derivatives formula need to be programmed.

- Conceptually simple.
- Potential for high precision. For example, if we get 20 derivatives of x(t), then the method is order 20 (that is, terms up to and including the one involving  $h^{20}$ ).

#### **Disadvantages:**

- The method depends on repeated differentiation of the differential equation, unless we intend to use only the method of order
   For high order methods, f(t, x) must have partial derivatives of sufficient high order in the region where are solving the problem. Such an assumption is not necessary for the existence of a solution.
- The various derivatives formula need to be programmed.

- Conceptually simple.
- Potential for high precision. For example, if we get 20 derivatives of x(t), then the method is order 20 (that is, terms up to and including the one involving  $h^{20}$ ).

#### **Disadvantages:**

- The method depends on repeated differentiation of the differential equation, unless we intend to use only the method of order
   For high order methods, f(t, x) must have partial derivatives of sufficient high order in the region where are solving the problem. Such an assumption is not necessary for the existence of a solution.
- The various derivatives formula need to be programmed.

- Conceptually simple.
- Potential for high precision. For example, if we get 20 derivatives of x(t), then the method is order 20 (that is, terms up to and including the one involving  $h^{20}$ ).

### §6.3 Euler's Method

• If n = 1, the Taylor series method reduces to (forward/explicit) Euler's method:

$$x_{n+1}=x_n+hf(t_n,x_n).$$

- Advantage of the method is not to require any differentiation of f.
- Oisadvantage of the method is that the necessity of taking small value for h to gain acceptable precision.
- One can also consider the backward/implicit Euler's method:

$$x_{n+1} = x_n + hf(t_n + h, x_{n+1}).$$

To obtain  $x_{n+1}$ , it is required to solve a nonlinear equation.



The concept of the Runge-Kutta method is to provide a general procedure of deriving higher order methods, which does **NOT** involve the derivatives of **f**, of solving the IVP

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x}), \qquad \mathbf{x}(t_0) = \mathbf{x}_0.$$

To see that this is possible, we note that Taylor's method of order 2 provides that

$$x(t+h) = x(t) + hf(t,x) + \frac{h^2}{2} [f_t(t,x) + f_x(t,x)f(t,x)] + \mathcal{O}(h^3).$$

On the other hand, Taylor's Theorem implies that

$$f(t+h,x+hf(t,x)) = f(t,x) + f_t(t,x)h + f_x(t,x)hf(t,x) + \mathcal{O}(h^2);$$

thus

$$x(t+h) = x(t) + hf(t,x) + \frac{h}{2} [f(t+h,x+hf(t,x)) - f(t,x)] + \mathcal{O}(h^3)$$

which provides "a" second order Runge-Kutta method

$$x_{n+1} = x_n + hf(t_n, x_n) + \frac{h}{2} [f(t_{n+1}, x_n + hf(t_n, x_n)) - f(t_n, x_n)].$$

The concept of the Runge-Kutta method is to provide a general procedure of deriving higher order methods, which does **NOT** involve the derivatives of **f**, of solving the IVP

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x}), \qquad \mathbf{x}(t_0) = \mathbf{x}_0.$$

To see that this is possible, we note that Taylor's method of order 2 provides that

$$x(t+h) = x(t) + hf(t,x) + \frac{h^2}{2} [f_t(t,x) + f_x(t,x)f(t,x)] + \mathcal{O}(h^3).$$

On the other hand, Taylor's Theorem implies that

$$f(t + h, x + hf(t, x)) = f(t, x) + f_t(t, x)h + f_x(t, x)hf(t, x) + O(h^2);$$

thus

$$x(t+h) = x(t) + hf(t,x) + \frac{h}{2} \left[ f(t+h,x+hf(t,x)) - f(t,x) \right] + \mathcal{O}(h^3)$$
 which provides "a" second order Runge-Kutta method

$$x_{n+1} = x_n + hf(t_n, x_n) + \frac{h}{2} [f(t_{n+1}, x_n + hf(t_n, x_n)) - f(t_n, x_n)].$$

The concept of the Runge-Kutta method is to provide a general procedure of deriving higher order methods, which does **NOT** involve the derivatives of **f**, of solving the IVP

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x}), \qquad \mathbf{x}(t_0) = \mathbf{x}_0.$$

To see that this is possible, we note that Taylor's method of order 2 provides that

$$x(t+h) = x(t) + hf(t,x) + \frac{h^2}{2} [f_t(t,x) + f_x(t,x)f(t,x)] + \mathcal{O}(h^3).$$

On the other hand, Taylor's Theorem implies that

$$f(t+h,x+hf(t,x)) = f(t,x) + f_t(t,x)h + f_x(t,x)hf(t,x) + \mathcal{O}(h^2);$$
thus

$$x(t+h) = x(t) + hf(t,x) + \frac{h}{2} \left[ f(t+h, x+hf(t,x)) - f(t,x) \right] + \mathcal{O}(h^3)$$

which provides "a" second order Runge-Kutta method

$$x_{n+1} = x_n + hf(t_n, x_n) + \frac{h}{2} [f(t_{n+1}, x_n + hf(t_n, x_n)) - f(t_n, x_n)].$$

The concept of the Runge-Kutta method is to provide a general procedure of deriving higher order methods, which does **NOT** involve the derivatives of **f**, of solving the IVP

$$\mathbf{x}' = \mathbf{f}(t, \mathbf{x}), \qquad \mathbf{x}(t_0) = \mathbf{x}_0.$$

To see that this is possible, we note that Taylor's method of order 2 provides that

$$x(t+h) = x(t) + hf(t,x) + \frac{h^2}{2} [f_t(t,x) + f_x(t,x)f(t,x)] + \mathcal{O}(h^3).$$

On the other hand, Taylor's Theorem implies that

$$f(t+h,x+hf(t,x)) = f(t,x) + f_t(t,x)h + f_x(t,x)hf(t,x) + \mathcal{O}(h^2);$$
thus

$$x(t+h) = x(t) + hf(t,x) + \frac{h}{2} [f(t+h,x+hf(t,x)) - f(t,x)] + \mathcal{O}(h^3)$$

which provides "a" second order Runge-Kutta method

$$x_{n+1} = x_n + + \frac{h}{2} [f(t_{n+1}, x_n + hf(t_n, x_n)) + f(t_n, x_n)].$$

The 2nd-order Runge-Kutta (RK) method derived previously is often written as the following alternative form

$$x_{n+1} = x_n + \frac{h}{2}(k_1 + k_2),$$

where

$$k_1 = f(t_n, x_n), \quad k_2 = f(t_n + h, x_n + hk_1).$$

This is also known as Heun's method.

**Remark**: Heun's method can be viewed as evaluating the number  $\int_{t_n}^{t_{n+1}} f(t, x(t)) dt$  using the trapezoidal rule, while the exact value of  $x(t_{n+1})$  is unknown and is replaced by the forward Euler formula

$$x(t_{n+1}) - x(t_n) = \int_{t_n}^{t_{n+1}} x'(t) dt = \int_{t_n}^{t_{n+1}} f(t, x(t)) dt$$

$$\approx \frac{t_{n+1} - t_n}{2} \left[ f(t_n, x(t_n)) + f(t_{n+1}, x(t_{n+1})) \right]$$

$$\approx \frac{t_{n+1} - t_n}{2} \left[ f(t_n, x(t_n)) + f(t_{n+1}, x(t_n) + hf(t_n, x_n)) \right]$$

The 2nd-order Runge-Kutta (RK) method derived previously is often written as the following alternative form

$$x_{n+1} = x_n + \frac{h}{2}(k_1 + k_2),$$

where

$$k_1 = f(t_n, x_n), \quad k_2 = f(t_n + h, x_n + hk_1).$$

This is also known as Heun's method.

**Remark**: Heun's method can be viewed as evaluating the number  $\int_{t_n}^{t_{n+1}} f(t,x(t)) dt$  using the trapezoidal rule, while the exact value of  $x(t_{n+1})$  is unknown and is replaced by the forward Euler formula  $x_n + hf(t_n, x_n)$ :

$$x(t_{n+1}) - x(t_n) = \int_{t_n}^{t_{n+1}} x'(t) dt = \int_{t_n}^{t_{n+1}} f(t, x(t)) dt$$

$$\approx \frac{t_{n+1} - t_n}{2} \left[ f(t_n, x(t_n)) + f(t_{n+1}, x(t_{n+1})) \right]$$

$$\approx \frac{t_{n+1} - t_n}{2} \left[ f(t_n, x(t_n)) + f(t_{n+1}, x(t_n)) + hf(t_n, x_n) \right]$$

The 2nd-order Runge-Kutta (RK) method derived previously is often written as the following alternative form

$$x_{n+1} = x_n + \frac{h}{2}(k_1 + k_2),$$

where

$$k_1 = f(t_n, x_n), \quad k_2 = f(t_n + h, x_n + hk_1).$$

This is also known as Heun's method.

**Remark**: Heun's method can be viewed as evaluating the number  $\int_{t_n}^{t_{n+1}} f(t, x(t)) dt$  using the trapezoidal rule, while the exact value of  $x(t_{n+1})$  is unknown and is replaced by the forward Euler formula  $x_n + hf(t_n, x_n)$ :

$$x(t_{n+1}) - x(t_n) = \int_{t_n}^{t_{n+1}} x'(t) dt = \int_{t_n}^{t_{n+1}} f(t, x(t)) dt$$

$$\approx \frac{t_{n+1} - t_n}{2} \left[ f(t_n, x(t_n)) + f(t_{n+1}, x(t_{n+1})) \right]$$

$$\approx \frac{t_{n+1} - t_n}{2} \left[ f(t_n, x(t_n)) + f(t_{n+1}, x(t_n) + hf(t_n, x_n)) \right]$$

The 2nd-order Runge-Kutta (RK) method derived previously is often written as the following alternative form

$$x_{n+1} = x_n + \frac{h}{2}(k_1 + k_2),$$

where

$$k_1 = f(t_n, x_n), \quad k_2 = f(t_n + h, x_n + hk_1).$$

This is also known as Heun's method.

**Remark**: Heun's method can be viewed as evaluating the number  $\int_{t_n}^{t_{n+1}} f(t, x(t)) dt$  using the trapezoidal rule, while the exact value of  $x(t_{n+1})$  is unknown and is replaced by the forward Euler formula  $x_n + hf(t_n, x_n)$ :

$$x(t_{n+1}) - x(t_n) = \int_{t_n}^{t_{n+1}} x'(t) dt = \int_{t_n}^{t_{n+1}} f(t, x(t)) dt$$

$$\approx \frac{t_{n+1} - t_n}{2} \left[ f(t_n, x(t_n)) + f(t_{n+1}, x(t_{n+1})) \right]$$

$$\approx \frac{t_{n+1} - t_n}{2} \left[ f(t_n, x(t_n)) + f(t_{n+1}, x(t_n) + hf(t_n, x_n)) \right].$$

In general, a 2nd order Runge-Kutta method is of to choose real numbers  $\omega_1,\omega_2$  and  $\alpha,\beta$  so that

$$\begin{aligned} x(t+h) &= x(t) + \omega_1 h f(t,x) + \omega_2 h \underline{f(t+\alpha h,x+\beta h f(t,x))} + \mathcal{O}(h^3), \\ &= x(t) + \omega_1 h f(t,x) + \omega_2 h \underline{\left[\underline{f(t,x) + \alpha h f_t(t,x) + \beta h f(t,x) f_x(t,x)\right]} \right.} \\ &+ \mathcal{O}(h^3) \\ &= x(t) + (\omega_1 + \omega_2) h f(t,x) + h^2 \underline{\left[\omega_2 \alpha f_t(t,x) + \omega_2 \beta f(t,x) f_x(t,x)\right]} \\ &+ \mathcal{O}(h^3). \end{aligned}$$

Comparing with

$$x(t+h) = x(t) + hf + \frac{h^2}{2} [f_t(t,x) + f_x(t,x)f(t,x)] + \mathcal{O}(h^3),$$

$$\omega_1 + \omega_2 = 1,$$
  
$$\omega_2 \alpha = \omega_2 \beta = 1/2.$$



In general, a 2nd order Runge-Kutta method is of to choose real numbers  $\omega_1,\omega_2$  and  $\alpha,\beta$  so that

$$\begin{aligned} x(t+h) &= x(t) + \omega_1 h f(t,x) + \omega_2 h \underline{f(t+\alpha h,x+\beta h f(t,x))} + \mathcal{O}(h^3), \\ &= x(t) + \omega_1 h f(t,x) + \omega_2 h \overline{\left[\underline{f(t,x) + \alpha h f_t(t,x) + \beta h f(t,x) f_x(t,x)}\right]} \\ &+ \mathcal{O}(h^3) \\ &= x(t) + (\omega_1 + \omega_2) h f(t,x) + h^2 \left[\omega_2 \alpha f_t(t,x) + \omega_2 \beta f(t,x) f_x(t,x)\right] \\ &+ \mathcal{O}(h^3). \end{aligned}$$

Comparing with

$$x(t+h) = x(t) + hf + \frac{h^2}{2} [f_t(t,x) + f_x(t,x)f(t,x)] + \mathcal{O}(h^3),$$

$$\omega_1 + \omega_2 = 1,$$
  
$$\omega_2 \alpha = \omega_2 \beta = 1/2.$$



In general, a 2nd order Runge-Kutta method is of to choose real numbers  $\omega_1,\omega_2$  and  $\alpha,\beta$  so that

$$\begin{aligned} x(t+h) &= x(t) + \omega_1 h f(t,x) + \omega_2 h \underline{f(t+\alpha h,x+\beta h f(t,x))} + \mathcal{O}(h^3), \\ &= x(t) + \omega_1 h f(t,x) + \omega_2 h \overline{\left[\underline{f(t,x) + \alpha h f_t(t,x) + \beta h f(t,x) f_x(t,x)}\right]} \\ &+ \mathcal{O}(h^3) \\ &= x(t) + (\omega_1 + \omega_2) h f(t,x) + h^2 \left[\omega_2 \alpha f_t(t,x) + \omega_2 \beta f(t,x) f_x(t,x)\right] \\ &+ \mathcal{O}(h^3). \end{aligned}$$

Comparing with

$$x(t+h) = x(t) + hf + \frac{h^2}{2} [f_t(t,x) + f_x(t,x)f(t,x)] + \mathcal{O}(h^3),$$

$$\omega_1 + \omega_2 = 1,$$
  
$$\omega_2 \alpha = \omega_2 \beta = 1/2$$



In general, a 2nd order Runge-Kutta method is of to choose real numbers  $\omega_1,\omega_2$  and  $\alpha,\beta$  so that

$$\begin{split} x(t+h) &= x(t) + \omega_1 h f(t,x) + \omega_2 h \underline{f(t+\alpha h,x+\beta h f(t,x))} + \mathcal{O}(h^3), \\ &= x(t) + \omega_1 h f(t,x) + \omega_2 h \overline{\left[\underline{f(t,x) + \alpha h f_t(t,x) + \beta h f(t,x) f_x(t,x)}\right]} \\ &+ \mathcal{O}(h^3) \\ &= x(t) + (\omega_1 + \omega_2) h f(t,x) + h^2 \Big[\omega_2 \alpha f_t(t,x) + \omega_2 \beta f(t,x) f_x(t,x)\Big] \\ &+ \mathcal{O}(h^3). \end{split}$$

Comparing with

$$x(t+h) = x(t) + hf + \frac{h^2}{2} [f_t(t,x) + f_x(t,x)f(t,x)] + \mathcal{O}(h^3),$$

$$\omega_1 + \omega_2 = 1,$$
  
$$\omega_2 \alpha = \omega_2 \beta = 1/2.$$



• The previous method (Heun's method) is obtained by setting

$$\begin{cases} \omega_1 = \omega_2 = 1/2, \\ \alpha = \beta = 1. \end{cases}$$

The modified Euler method is obtained by setting

$$\begin{cases} \omega_1 = 0, \\ \omega_2 = 1, \\ \alpha = \beta = 1/2, \end{cases}$$

we obtain the following:

$$x(t+h) \approx x(t) + hk_2$$

where

$$k_1 = f(t, x), \quad k_2 = f(t + \frac{h}{2}, x + \frac{h}{2}k_1).$$



- The derivations of higher order RK methods are tedious. However, the formulas are rather elegant and easily programmed once they have been derived.
- The most popular 4th order Runge-Kutta method is based on

$$x(t+h) = x(t) + \frac{h}{6}(k_1 + 2k_2 + 2k_3 + k_4) + \mathcal{O}(h^5),$$

where

$$\begin{cases} k_1 = f(t, x), \\ k_2 = f(t + \frac{h}{2}, x + \frac{h}{2}k_1), \\ k_3 = f(t + \frac{h}{2}, x + \frac{h}{2}k_2), \\ k_4 = f(t + h, x + hk_3). \end{cases}$$

We note that  $k_1, k_2, k_3, k_4$  are all approximated value of x'(t).

- The derivations of higher order RK methods are tedious. However, the formulas are rather elegant and easily programmed once they have been derived.
- The most popular 4th order Runge-Kutta method is given by

$$x_{n+1} = x_n + \frac{h}{6}(k_{n1} + 2k_{n2} + 2k_{n3} + k_{n4}),$$

where

$$\begin{cases} k_{n1} = f(t_n, x_n), \\ k_{n2} = f(t_n + \frac{h}{2}, x_n + \frac{h}{2}k_{n1}), \\ k_{n3} = f(t_n + \frac{h}{2}, x_n + \frac{h}{2}k_{n2}), \\ k_{n4} = f(t_n + h, x_n + hk_{n3}). \end{cases}$$

We note that  $k_{n1}$ ,  $k_{n2}$ ,  $k_{n3}$ ,  $k_{n4}$  are all approximated value of  $x'(t_n)$ .

To see that the method provided in the previous page is indeed fourth order, we apply Taylor's Theorem and obtain that

$$\begin{split} k_2 &= f(t + \frac{h}{2}, x + \frac{h}{2}k_1) \\ &= f(t, x) + f_t(t, x) \cdot \frac{h}{2} + f_x(t, x) \cdot \frac{h}{2}k_1 \\ &+ \frac{1}{2} \Big[ f_{tt}(t, x) \cdot \frac{h^2}{4} + 2f_{tx}(t, x) \cdot \frac{h}{2} \cdot \frac{h}{2}k_1 + f_{xx}(t, x) \cdot \frac{h^2k_1^2}{4} \Big] \\ &+ \frac{1}{6} \Big[ f_{ttt}(t, x) \cdot \frac{h^3}{8} + 3f_{ttx}(t, x) \cdot \frac{h^2}{4} \cdot \frac{h}{2}k_1 + 3f_{txx}(t, x) \cdot \frac{h}{2} \cdot \frac{h^2k_1^2}{4} \\ &+ f_{xxx}(t, x) \cdot \frac{h^3k_1^3}{8} \Big] + \mathcal{O}(h^4) \,. \end{split}$$

The reason why we truncate at the order  $\mathcal{O}(\mathit{h}^4)$  is that we are going to prove that

$$x(t+h) = x(t) + \frac{h}{6}(k_1 + 2k_2 + 2k_3 + k_4) + \mathcal{O}(h^5),$$

is an order four numerical method for solving ODE.

To see that the method provided in the previous page is indeed fourth order, we apply Taylor's Theorem and obtain that

$$\begin{split} k_2 &= f(t + \frac{h}{2}, x + \frac{h}{2}k_1) \\ &= f(t, x) + f_t(t, x) \cdot \frac{h}{2} + f_x(t, x) \cdot \frac{h}{2}k_1 \\ &+ \frac{1}{2} \Big[ f_{tt}(t, x) \cdot \frac{h^2}{4} + 2f_{tx}(t, x) \cdot \frac{h}{2} \cdot \frac{h}{2}k_1 + f_{xx}(t, x) \cdot \frac{h^2 k_1^2}{4} \Big] \\ &+ \frac{1}{6} \Big[ f_{ttt}(t, x) \cdot \frac{h^3}{8} + 3f_{ttx}(t, x) \cdot \frac{h^2}{4} \cdot \frac{h}{2}k_1 + 3f_{txx}(t, x) \cdot \frac{h}{2} \cdot \frac{h^2 k_1^2}{4} \\ &+ f_{xxx}(t, x) \cdot \frac{h^3 k_1^3}{8} \Big] + \mathcal{O}(h^4) \,. \end{split}$$

The reason why we truncate at the order  $\mathcal{O}(\mathit{h}^4)$  is that we are going to prove that

$$x(t+h) = x(t) + \frac{h}{6}(k_1 + 2k_2 + 2k_3 + k_4) + \mathcal{O}(h^5),$$

is an order four numerical method for solving ODE.

$$k_{2} = f(t,x) + \frac{h}{2} f_{t}(t,x) + \frac{h}{2} k_{1} f_{x}(t,x) + \frac{h^{2}}{8} \left[ f_{tt}(t,x) + 2k_{1} f_{tx}(t,x) + k_{1}^{2} f_{xx}(t,x) \right] + \frac{h^{3}}{48} \left[ f_{ttt}(t,x) + 3k_{1} f_{ttx}(t,x) + 3k_{1}^{2} f_{txx}(t,x) + k_{1}^{3} f_{xxx}(t,x) \right] + \mathcal{O}(h^{4}), k_{3} = f(t,x) + \frac{h}{2} f_{t}(t,x) + \frac{h}{2} k_{2} f_{x}(t,x) + \frac{h^{2}}{8} \left[ f_{tt}(t,x) + 2k_{2} f_{tx}(t,x) + k_{2}^{2} f_{xx}(t,x) \right] + \frac{h^{3}}{48} \left[ f_{ttt}(t,x) + 3k_{2} f_{ttx}(t,x) + 3k_{2}^{2} f_{txx}(t,x) + k_{2}^{3} f_{xxx}(t,x) \right] + \mathcal{O}(h^{4}), k_{4} = f(t,x) + h f_{t}(t,x) + h k_{3} f_{x}(t,x) + \frac{h^{2}}{2} \left[ f_{tt}(t,x) + 2k_{3} f_{tx}(t,x) + k_{3}^{2} f_{txx}(t,x) + k_{3}^{3} f_{xxx}(t,x) \right] + \frac{h^{3}}{6} \left[ f_{ttt}(t,x) + 3k_{3} f_{ttx}(t,x) + 3k_{3}^{2} f_{txx}(t,x) + k_{3}^{3} f_{xxx}(t,x) \right] + \mathcal{O}(h^{4}).$$

$$\begin{aligned} k_2 &= f(t,x) + \frac{h}{2} f_t(t,x) + \frac{h}{2} k_1 f_x(t,x) \\ &+ \frac{h^2}{8} \left[ f_{tt}(t,x) + 2k_1 f_{tx}(t,x) + k_1^2 f_{xx}(t,x) \right] \\ &+ \frac{h^3}{48} \left[ f_{ttt}(t,x) + 3k_1 f_{ttx}(t,x) + 3k_1^2 f_{txx}(t,x) + k_1^3 f_{xxx}(t,x) \right] + \mathcal{O}(h^4), \\ k_3 &= f(t,x) + \frac{h}{2} f_t(t,x) + \frac{h}{2} k_2 f_x(t,x) \\ &+ \frac{h^2}{8} \left[ f_{tt}(t,x) + 2k_2 f_{tx}(t,x) + k_2^2 f_{xx}(t,x) \right] \\ &+ \frac{h^3}{48} \left[ f_{ttt}(t,x) + 3k_2 f_{ttx}(t,x) + 3k_2^2 f_{txx}(t,x) + k_2^3 f_{xxx}(t,x) \right] + \mathcal{O}(h^4), \\ k_4 &= f(t,x) + h f_t(t,x) + h k_3 f_x(t,x) \\ &+ \frac{h^2}{2} \left[ f_{tt}(t,x) + 2k_3 f_{tx}(t,x) + k_3^2 f_{xx}(t,x) + k_3^3 f_{xxx}(t,x) \right] + \mathcal{O}(h^4). \end{aligned}$$

$$k_{2} = f(t,x) + \frac{h}{2}f_{t}(t,x) + \frac{h}{2}k_{1}f_{x}(t,x) + \frac{h^{2}}{8} \Big[ f_{tt}(t,x) + 2k_{1}f_{tx}(t,x) + k_{1}^{2}f_{xx}(t,x) \Big] + \frac{h^{3}}{48} \Big[ f_{ttt}(t,x) + 3k_{1}f_{ttx}(t,x) + 3k_{1}^{2}f_{txx}(t,x) + k_{1}^{3}f_{xxx}(t,x) \Big] + \mathcal{O}(h^{4}), k_{3} = f(t,x) + \frac{h}{2}f_{t}(t,x) + \frac{h}{2}k_{2}f_{x}(t,x) + \frac{h^{2}}{8} \Big[ f_{tt}(t,x) + 2k_{2}f_{tx}(t,x) + k_{2}^{2}f_{xx}(t,x) \Big] + \frac{h^{3}}{48} \Big[ f_{ttt}(t,x) + 3k_{2}f_{ttx}(t,x) + 3k_{2}^{2}f_{txx}(t,x) + k_{2}^{3}f_{xxx}(t,x) \Big] + \mathcal{O}(h^{4}), k_{4} = f(t,x) + hf_{t}(t,x) + hk_{3}f_{x}(t,x) + \frac{h^{2}}{2} \Big[ f_{tt}(t,x) + 2k_{3}f_{tx}(t,x) + k_{3}^{2}f_{xx}(t,x) + k_{3}^{3}f_{xxx}(t,x) \Big] + \frac{h^{3}}{6} \Big[ f_{ttt}(t,x) + 3k_{3}f_{ttx}(t,x) + 3k_{3}^{2}f_{txx}(t,x) + k_{3}^{3}f_{xxx}(t,x) \Big] + \mathcal{O}(h^{4}).$$

Therefore,

$$\begin{split} k_1 + 2k_2 + 2k_3 + k_4 \\ &= 6f(t,x) + 3hf_t(t,x) + h(k_1 + k_2 + k_3)f_x(t,x) \\ &+ h^2f_{tt}(t,x) + \frac{h^2}{2}(k_1 + k_2 + 2k_3)f_{tx}(t,x) \\ &+ \frac{h^2}{4}(k_1^2 + k_2^2 + 2k_3^2)f_{xx}(t,x) + \frac{h^3}{4}f_{ttt}(t,x) \\ &+ \frac{h^3}{8}(k_1 + k_2 + 4k_3)f_{ttx}(t,x) + \frac{h^3}{8}(k_1^2 + k_2^2 + 4k_3^2)f_{txx}(t,x) \\ &+ \frac{h^3}{24}(k_1^3 + k_2^3 + 4k_3^3)f_{xxx}(t,x) + \mathcal{O}(h^4) \,. \end{split}$$

Next, we need to compute

- ①  $k_1 + k_2 + k_3$  accurate at least of order  $\mathcal{O}(h^2)$ .
- ②  $k_1+k_2+2k_3$  and  $k_1^2+k_2^2+2k_3^2$  accurate at least of order  $\mathcal{O}(h)$ .
- ①  $k_1 + k_2 + 4k_3$ ,  $k_1^2 + k_2^2 + 4k_3^2$  and  $k_1^3 + k_2^3 + 4k_3^3$  accurate at least of order  $\mathcal{O}(1)$ .

Therefore,

$$\begin{split} k_1 + 2k_2 + 2k_3 + k_4 \\ &= 6f(t, x) + 3hf_t(t, x) + h(k_1 + k_2 + k_3)f_x(t, x) \\ &+ h^2 f_{tt}(t, x) + \frac{h^2}{2}(k_1 + k_2 + 2k_3)f_{tx}(t, x) \\ &+ \frac{h^2}{4}(k_1^2 + k_2^2 + 2k_3^2)f_{xx}(t, x) + \frac{h^3}{4}f_{ttt}(t, x) \\ &+ \frac{h^3}{8}(k_1 + k_2 + 4k_3)f_{ttx}(t, x) + \frac{h^3}{8}(k_1^2 + k_2^2 + 4k_3^2)f_{txx}(t, x) \\ &+ \frac{h^3}{24}(k_1^3 + k_2^3 + 4k_3^3)f_{xxx}(t, x) + \mathcal{O}(h^4) \,. \end{split}$$

Next, we need to compute

- **1**  $k_1 + k_2 + k_3$  accurate at least of order  $\mathcal{O}(h^2)$ .
- $k_1 + k_2 + 2k_3$  and  $k_1^2 + k_2^2 + 2k_3^2$  accurate at least of order  $\mathcal{O}(h)$ .
- **3**  $k_1 + k_2 + 4k_3$ ,  $k_1^2 + k_2^2 + 4k_3^2$  and  $k_1^3 + k_2^3 + 4k_3^3$  accurate at least of order  $\mathcal{O}(1)$ .

Using the formula for  $k_2$ ,  $k_3$ , and  $k_4$ , we find that

$$\begin{aligned} k_1 + k_2 + k_3 &= 3f(t,x) + hf_t(t,x) + \frac{h}{2}(k_1 + k_2)f_x(t,x) + \frac{h^2}{4}f_{tt}(t,x) \\ &\quad + \frac{h^2}{4}(k_1 + k_2)f_{tx}(t,x) + \frac{h^2}{8}(k_1^2 + k_2^2)f_{xx}(t,x) + \mathcal{O}(h^3), \\ k_1 + k_2 + 2k_3 &= 4f(t,x) + \frac{3h}{2}f_t(t,x) + \frac{h}{2}(k_1 + 2k_2)f_x(t,x) + \mathcal{O}(h^2), \\ k_1^2 + k_2^2 + 2k_3^2 &= 4f(t,x)^2 + hf(t,x) \left[ 3f_t(t,x) + (k_1 + 2k_2)f_x(t,x) \right] \\ &\quad + \mathcal{O}(h^2), \\ k_1 + k_2 + 4k_3 &= 6f(t,x) + \mathcal{O}(h), \\ k_1^2 + k_2^2 + 4k_3^2 &= 6f(t,x)^2 + \mathcal{O}(h), \\ k_1^3 + k_2^3 + 4k_3^3 &= 6f(t,x)^3 + \mathcal{O}(h). \end{aligned}$$

Since  $k_1 + 2k_2 = 3f(t, x) + \mathcal{O}(h)$ , to continue we compute

- ①  $k_1 + k_2$  accurate at least of order  $\mathcal{O}(h^2)$ .
- 2  $k_1^2 + k_2^2$  accurate at least of order  $\mathcal{O}(h)$ .



Using the formula for  $k_2$ ,  $k_3$ , and  $k_4$ , we find that

$$k_{1}+k_{2}+k_{3} = 3f(t,x)+hf_{t}(t,x)+\frac{h}{2}(k_{1}+k_{2})f_{x}(t,x)+\frac{h^{2}}{4}f_{tt}(t,x)$$

$$+\frac{h^{2}}{4}(k_{1}+k_{2})f_{tx}(t,x)+\frac{h^{2}}{8}(k_{1}^{2}+k_{2}^{2})f_{xx}(t,x)+\mathcal{O}(h^{3}),$$

$$k_{1}+k_{2}+2k_{3} = 4f(t,x)+\frac{3h}{2}f_{t}(t,x)+\frac{h}{2}(k_{1}+2k_{2})f_{x}(t,x)+\mathcal{O}(h^{2}),$$

$$k_{1}^{2}+k_{2}^{2}+2k_{3}^{2} = 4f(t,x)^{2}+hf(t,x)\left[3f_{t}(t,x)+(k_{1}+2k_{2})f_{x}(t,x)\right]$$

$$+\mathcal{O}(h^{2}),$$

$$k_{1}+k_{2}+4k_{3} = 6f(t,x)+\mathcal{O}(h),$$

$$k_{1}^{2}+k_{2}^{2}+4k_{3}^{2} = 6f(t,x)^{2}+\mathcal{O}(h),$$

$$k_{1}^{3}+k_{2}^{3}+4k_{3}^{3} = 6f(t,x)^{3}+\mathcal{O}(h).$$

Since  $k_1 + 2k_2 = 3f(t, x) + \mathcal{O}(h)$ , to continue we compute

- $k_1 + k_2$  accurate at least of order  $\mathcal{O}(h^2)$ .
- 2  $k_1^2 + k_2^2$  accurate at least of order  $\mathcal{O}(h)$ .



Since

$$k_1 + k_2 = 2f(t, x) + \frac{h}{2} \left[ f_t(t, x) + f(t, x) f_x(t, x) \right] + \mathcal{O}(h^2)$$
  

$$k_1^2 + k_2^2 = 2f(t, x)^2 + hf(t, x) \left[ f(t, x) + f(t, x) f_t(t, x) \right] + \mathcal{O}(h^2)$$

we obtain that

$$\begin{aligned} k_1 + k_2 + k_3 &= 3f(t,x) + h \big[ f_t(t,x) + f(t,x) f_x(t,x) \big] \\ &+ \frac{h^2}{4} \big[ f_{tt}(t,x) + f_t(t,x) f_x(t,x) + f(t,x) f_x(t,x)^2 \\ &+ 2f(t,x) f_{tx}(t,x) + f(t,x)^2 f_{xx}(t,x) \big] + \mathcal{O}(h^3) \\ k_1 + k_2 + 2k_3 &= 4f(t,x) + \frac{3h}{2} \big[ f_t(t,x) + f(t,x) f_x(t,x) \big] + \mathcal{O}(h^2), \\ k_1^2 + k_2^2 + 2k_3^2 &= 4f(t,x)^2 + 3hf(t,x) \big[ f_t(t,x) + f(t,x) f_x(t,x) \big] + \mathcal{O}(h^2), \\ k_1 + k_2 + 4k_3 &= 6f(t,x) + \mathcal{O}(h), \\ k_1^2 + k_2^2 + 4k_3^2 &= 6f(t,x)^2 + \mathcal{O}(h), \\ k_1^3 + k_2^3 + 4k_3^3 &= 6f(t,x)^3 + \mathcal{O}(h). \end{aligned}$$

Therefore,

$$\begin{split} k_1 + 2k_2 + 2k_3 + k_4 \\ &= 6f(t,x) + 3h \big[ f_t(t,x) + f(t,x) f_x(t,x) \big] \\ &+ h^2 \big[ f_t(t,x) + f(t,x) f_x(t,x) \big] f_x(t,x) \\ &+ \frac{h^3}{4} \big[ f_{tt}(t,x) + f_t(t,x) f_x(t,x) + f(t,x) f_x(t,x)^2 \\ &\quad + 2f(t,x) f_{tx}(t,x) + f(t,x)^2 f_{xx}(t,x) \big] f_x(t,x) + h^2 f_{tt}(t,x) \\ &+ 2h^2 f(t,x) f_{tx}(t,x) + \frac{3h^3}{4} \big[ f_t(t,x) + f(t,x) f_x(t,x) \big] f_{tx}(t,x) \\ &+ h^2 f(t,x)^2 f_{xx}(t,x) + \frac{3h^3}{4} f(t,x) \big[ f_t(t,x) + f(t,x) f_x(t,x) \big] f_{xx}(t,x) \\ &+ \frac{h^3}{4} f_{ttt}(t,x) + \frac{3h^3}{4} f(t,x) f_{ttx}(t,x) + \frac{3h^3}{4} f(t,x)^2 f_{txx}(t,x) \\ &+ \frac{h^3}{4} f(t,x)^3 f_{xxx}(t,x) + \mathcal{O}(h^4) \; . \end{split}$$

Therefore,

$$k_{1} + 2k_{2} + 2k_{3} + k_{4}$$

$$= 6f(t, x) + 3h [f_{t}(t, x) + f(t, x)f_{x}(t, x)]$$

$$+ h^{2} [f_{tt}(t, x) + f_{t}(t, x)f_{x}(t, x) + f(t, x)f_{x}(t, x)^{2} + 2f(t, x)f_{tx}(t, x)$$

$$+ f(t, x)^{2} f_{xx}(t, x)]$$

$$+ \frac{h^{3}}{4} [f_{ttt}(t, x) + f_{tt}(t, x)f_{x}(t, x) + f_{t}(t, x)f_{x}(t, x)^{2} + f(t, x)f_{x}(t, x)^{3}$$

$$+ f(t, x)^{3} f_{xxx}(t, x) + 2f(t, x)f_{x}(t, x)f_{tx}(t, x)$$

$$+ 3f(t, x)f_{ttx}(t, x) + 3f_{t}(t, x)f_{tx}(t, x)$$

$$+ 3f(t, x)^{2} f_{txx}(t, x) + 4f(t, x)^{2} f_{x}(t, x)f_{xx}(t, x)] + \mathcal{O}(h^{4})$$

which implies that

$$x(t) + \frac{h}{6} \left[ k_1 + 2k_2 + 2k_3 + k_4 \right]$$

$$= x(t) + hf(t, x) + \frac{h^2}{2} \left[ f_t(t, x) + f(t, x) f_x(t, x) \right]$$

$$+ \frac{h^3}{6} \left[ f_{tt}(t, x) + f_t(t, x) f_x(t, x) + f(t, x) f_x(t, x)^2 + 2f(t, x) f_{tx}(t, x) \right]$$

$$+ f(t, x)^2 f_{xx}(t, x)$$

$$+ \frac{h^4}{24} \left[ f_{ttt}(t, x) + f_{tt}(t, x) f_x(t, x) + f_t(t, x) f_x(t, x)^2 + f(t, x) f_x(t, x)^3 \right]$$

$$+ f(t, x)^3 f_{xxx}(t, x) + 2f(t, x) f_x(t, x) f_{tx}(t, x)$$

$$+ 3f(t, x) f_{ttx}(t, x) + 3f_t(t, x) f_{tx}(t, x)$$

$$+ 3f(t, x) f_x(t, x) f_{tx}(t, x) + 3f(t, x) f_t(t, x) f_{xx}(t, x)$$

$$+ 3f(t, x)^2 f_{txx}(t, x) + 4f(t, x)^2 f_x(t, x) f_{xx}(t, x) \right] + \mathcal{O}(h^5)$$

$$= x(t) + hx'(t) + \frac{h^2}{2} x''(t) + \frac{h^3}{6} x'''(t) + \frac{h^4}{24} x^{(4)}(t) + \mathcal{O}(h^5).$$

which implies that

$$x(t) + \frac{h}{6} \left[ k_1 + 2k_2 + 2k_3 + k_4 \right]$$

$$= x(t) + hf(t, x) + \frac{h^2}{2} \left[ f_t(t, x) + f(t, x) f_x(t, x) \right]$$

$$+ \frac{h^3}{6} \left[ f_{tt}(t, x) + f_t(t, x) f_x(t, x) + f(t, x) f_x(t, x)^2 + 2f(t, x) f_{tx}(t, x) \right]$$

$$+ f(t, x)^2 f_{xx}(t, x)$$

$$+ \frac{h^4}{24} \left[ f_{ttt}(t, x) + f_{tt}(t, x) f_x(t, x) + f_t(t, x) f_x(t, x)^2 + f(t, x) f_x(t, x)^3 \right]$$

$$+ f(t, x)^3 f_{xxx}(t, x) + 2f(t, x) f_x(t, x) f_{tx}(t, x)$$

$$+ 3f(t, x) f_{ttx}(t, x) + 3f_t(t, x) f_t(t, x) f_t(t, x) f_x(t, x)$$

$$+ 3f(t, x)^2 f_{txx}(t, x) + 4f(t, x)^2 f_x(t, x) f_{xx}(t, x) \right] + \mathcal{O}(h^5)$$

$$= x(t) + hx'(t) + \frac{h^2}{2} x''(t) + \frac{h^3}{6} x'''(t) + \frac{h^4}{24} x^{(4)}(t) + \mathcal{O}(h^5).$$

which implies that

$$x(t) + \frac{h}{6} [k_1 + 2k_2 + 2k_3 + k_4]$$

$$= x(t) + h f(t, x) + \frac{h^2}{2} [f_t(t, x) + f(t, x) f_x(t, x)]$$

$$+ \frac{h^3}{6} [f_{tt}(t, x) + f_t(t, x) f_x(t, x) + f(t, x) f_x(t, x)^2 + 2f(t, x) f_{tx}(t, x)$$

$$+ f(t, x)^2 f_{xx}(t, x)]$$

$$+ \frac{h^4}{24} [f_{ttt}(t, x) + f_{tt}(t, x) f_x(t, x) + f_t(t, x) f_x(t, x)^2 + f(t, x) f_x(t, x)^3$$

$$+ f(t, x)^3 f_{xxx}(t, x) + 2f(t, x) f_x(t, x) f_{tx}(t, x)$$

$$+ 3f(t, x) f_{ttx}(t, x) + 3f_t(t, x) f_t(t, x) f_t(t, x) f_x(t, x)$$

$$+ 3f(t, x)^2 f_{txx}(t, x) + 4f(t, x)^2 f_x(t, x) f_{xx}(t, x)] + \mathcal{O}(h^5)$$

$$= x(t) + h x'(t) + \frac{h^2}{2} x''(t) + \frac{h^3}{6} x'''(t) + \frac{h^4}{24} x^{(4)}(t) + \mathcal{O}(h^5).$$

#### **Computer Project:**

• Use the most popular 4th order Runge-Kutta with h=1/128, h=1/256 and h=1/512 to solve the following IVP for  $t\in [1,3]$  and then plot the piecewise linear approximate solution:

$$\begin{cases} x'(t) = t^{-2}(tx - x^2), \\ x(1) = 2. \end{cases}$$

2 Also plot the exact solution:

$$x(t) = (1/2 + \ln t)^{-1}t.$$

**3** Find the global truncation error E(h) and "verify" numerically that the 4th order Runge-Kutta is indeed a fourth order numerical method of solving IVPs.

end do

# §6.4 The Runge-Kutta Method

$$\begin{array}{l} \text{input } M \leftarrow 256, \ t \leftarrow 1.0, \ h \leftarrow 0.0078125, \ x \leftarrow 2.0 \\ \text{define } f(t,x) = (tx-x^2)/t^2, \ u(t) = t/(1/2 + \ln t) \\ e \leftarrow |u(t)-x| \\ \text{output } 0,t,x,e \\ \text{for } k = 1 \ \text{to } M \ \text{do} \\ k_1 \leftarrow f(t,x) \\ k_2 \leftarrow f(t+\frac{h}{2},x+\frac{h}{2}k_1) \\ k_3 \leftarrow f(t+\frac{h}{2},x+\frac{h}{2}k_2) \\ k_4 \leftarrow f(t+h,x+hk_3) \\ x \leftarrow x+\frac{h}{6}(k_1+2k_2+2k_3+k_4) \\ t \leftarrow t+h \\ e \leftarrow |u(t)-x| \\ \text{output } k,t,x,e \end{array}$$

For a system of equations  $\mathbf{x}'(t) = \mathbf{f}(t, \mathbf{x})$ , we have 4th-order Runge-Kutta method:

Other methods, they are all similar to the single equation case.

The idea of **collocation methods** is to choose a finite-dimensional space of candidate solutions (usually polynomials up to a certain degree) and a number of points in the domain (called collocation points), and to select that solution which satisfies the given equation at the collocation points.

Suppose that the ordinary differential equation

$$y'(t) = f(t, y(t)),$$
  $y(t_0) = y_0$ 

is to be solved over the interval  $[t_0, t_0 + c_k h]$ . Choose  $c_k$  from  $0 \le c_1 < c_2 < ... < c_n \le 1$ . The corresponding (polynomial) collocation method approximates the solution p by the polynomial p of degree p which satisfies the initial condition  $p(t_0) = y_0$  and the differential equation  $p'(t_k) = f(t_k, p(t_k))$ .

The idea of **collocation methods** is to choose a finite-dimensional space of candidate solutions (usually polynomials up to a certain degree) and a number of points in the domain (called collocation points), and to select that solution which satisfies the given equation at the collocation points.

Suppose that the ordinary differential equation

$$y'(t) = f(t, y(t)),$$
  $y(t_0) = y_0$ 

is to be solved over the interval  $[t_0, t_0 + c_k h]$ . Choose  $c_k$  from  $0 \le c_1 < c_2 < ... < c_n \le 1$ . The corresponding (polynomial) collocation method approximates the solution y by the polynomial p of degree n which satisfies the initial condition  $p(t_0) = y_0$  and the differential equation  $p'(t_k) = f(t_k, p(t_k))$ .

#### Example

Pick to collocation points  $c_1=0$  and  $c_2=1$ . The collocation method is then looking for a polynomial p of degree 2 satisfying the collocation conditions

$$p(t_0) = y_0,$$
  
 $p'(t_0) = f(t_0, p(t_0)),$   
 $p'(t_1) = f(t_1, p(t_1)).$ 

If  $p(t) = \alpha(t - t_0)^2 + \beta(t - t_0) + \gamma$ , then the collocation conditions above become

$$\begin{split} \gamma &= y_0, \\ \beta &= f(t_0, \gamma), \\ \alpha &= \frac{1}{2h} \Big[ f(t_0 + h, \alpha h^2 + \beta h + \gamma) - f(t_0, \gamma) \Big]. \end{split}$$

It remains to solve for lpha from a nonlinear equation.



#### Example

Pick to collocation points  $c_1=0$  and  $c_2=1$ . The collocation method is then looking for a polynomial p of degree 2 satisfying the collocation conditions

$$p(t_0) = y_0,$$
  
 $p'(t_0) = f(t_0, p(t_0)),$   
 $p'(t_1) = f(t_1, p(t_1)).$ 

If  $p(t) = \alpha(t - t_0)^2 + \beta(t - t_0) + \gamma$ , then the collocation conditions above become

$$\gamma = y_0, 
\beta = f(t_0, \gamma), 
\alpha = \frac{1}{2h} \Big[ f(t_0 + h, \alpha h^2 + \beta h + \gamma) - f(t_0, \gamma) \Big].$$

It remains to solve for  $\alpha$  from a nonlinear equation.

Suppose that we have a linear differential operator L and we wish to solve the equation: given a function f, find the function u satisfying

$$Lu(t) = f(t), \quad a < t < b.$$

- Let  $\{v_1, v_2, \dots, v_n\}$  be a set of functions that are linearly independent. Suppose that  $u(t) \approx c_1 v_1(t) + c_2 v_2(t) + \cdots + c_n v_n(t)$ .
- 2 Then solve  $L(\sum_{i=1}^{n} c_i v_i(t)) = f(t)$ . How to determine  $c_i$ ?
- **1** Let  $t_1, t_2, \dots, t_n$  be *n* prescribed points (collocation points) in the domain of u and f. Then we require that

$$\sum_{i=1}^{n} c_{j}(Lv_{j})(t_{i}) = f(t_{i}), \quad i = 1, 2, \cdots, n.$$

• This is a system of n linear equations in n unknowns  $c_i$ . The functions  $v_i$  and the points  $t_i$  should be chosen so that the matrix with entries  $(Lv_j)(t_i)$  is non-singular.

#### Example (Collocation method for Sturm-Liouville BVPs)

Consider a Sturm-Liouville two-point BVP:

$$\begin{cases} u''(t) + p(t)u'(t) + q(t)u(t) = f(t) & \forall t \in (0,1), \\ u(0) = u(1) = 0, \end{cases}$$

where p, q, f are given continuous functions on [0, 1]

• Let Lu := u'' + pu' + qu. Define the vector space

$$\mathcal{V} = \Big\{ u \in C^2((0,1)) \cap C([0,1]) \, \Big| \, u(0) = u(1) = 0 \Big\}.$$

If u is an exact solution of  $(\Box)$ , then  $u \in \mathcal{V}$ .

One set of functions is given by

$$v_{ik}(t) = t^{j}(1-t)^{k} \in C^{2}([0,1]), \quad 1 \le j \le m, 1 \le k \le n.$$

The finite-difference methods (FDM) are discretizations used for solving differential equations by approximating them with difference equations that finite differences approximate the derivatives.

#### Example

Again we consider the Sturm-Liouville two-point BVP:

$$\begin{cases} u''(t) + p(t)u'(t) + q(t)u(t) = f(t) & \forall t \in (0,1), \\ u(0) = u(1) = 0, \end{cases}$$

Let  $t_k = \frac{k}{n+1}$  for  $k = 0, 1, \dots, n+1$ . Using the central difference

approximation, (a) implies that

$$\frac{u(t_{k-1}) - 2u(t_k) + u(t_{k+1})}{h^2} + p(t_k) \frac{u(t_{k+1}) - u(t_{k-1})}{2h} + q(t_k) u(t_k)$$

$$= f(t_k) + \mathcal{O}(h^2) \quad \text{for } k = 1, 2, \dots, n,$$

$$u(t_0) = u(t_{n+1}) = 0.$$

The finite-difference methods (FDM) are discretizations used for solving differential equations by approximating them with difference equations that finite differences approximate the derivatives.

#### Example

Again we consider the Sturm-Liouville two-point BVP:

$$\begin{cases} u''(t) + p(t)u'(t) + q(t)u(t) = f(t) & \forall t \in (0,1), \\ u(0) = u(1) = 0, \end{cases}$$

Let  $t_k = \frac{k}{n+1}$  for  $k = 0, 1, \dots, n+1$ . Using the central difference approximation, ( $\square$ ) implies that

$$\begin{split} \frac{\textit{u}(t_{k-1}) - 2\textit{u}(t_k) + \textit{u}(t_{k+1})}{\textit{h}^2} + \textit{p}(t_k) \frac{\textit{u}(t_{k+1}) - \textit{u}(t_{k-1})}{2\textit{h}} + \textit{q}(t_k) \textit{u}(t_k) \\ &= \textit{f}(t_k) + \mathcal{O}(\textit{h}^2) \qquad \text{for } \textit{k} = 1, 2, \cdots, \textit{n}, \\ \textit{u}(t_0) &= \textit{u}(t_{n+1}) = 0. \end{split}$$

The finite-difference methods (FDM) are discretizations used for solving differential equations by approximating them with difference equations that finite differences approximate the derivatives.

#### Example

Again we consider the Sturm-Liouville two-point BVP:

$$\begin{cases} u''(t) + p(t)u'(t) + q(t)u(t) = f(t) & \forall t \in (0,1), \\ u(0) = u(1) = 0, \end{cases}$$

Let  $t_k = \frac{k}{n+1}$  for  $k = 0, 1, \dots, n+1$ . Using the central difference approximation, ( $\square$ ) implies that

$$\begin{split} \frac{\textit{u}(t_{k-1}) - 2\textit{u}(t_k) + \textit{u}(t_{k+1})}{\textit{h}^2} + \frac{\textit{p}(t_k)\textit{u}(t_{k+1}) + 2\textit{h}\textit{q}(t_k)\textit{u}(t_k) - \textit{p}(t_k)\textit{u}(t_{k-1})}{2\textit{h}} \\ &= \textit{f}(t_k) + \mathcal{O}(\textit{h}^2) \qquad \text{for } \textit{k} = 1, 2, \cdots, \textit{n}, \\ \textit{u}(t_0) &= \textit{u}(t_{n+1}) = 0. \end{split}$$

#### Example (Cont'd)

Define A(h) as the  $n \times n$  matrix

$$\frac{1}{h^2} \begin{bmatrix} -2 & 1 & 0 & \cdots & 0 \\ 1 & -2 & 1 & \ddots & 0 \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & 1 & -2 & 1 \\ 0 & \cdots & 0 & 1 & -2 \end{bmatrix} + \frac{1}{2h} \begin{bmatrix} 2hq(t_1) & p(t_1) & 0 & \cdots & 0 \\ -p(t_2) & 2hq(t_2) & p(t_2) & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & -p(t_{n-1}) & 2hq(t_{n-1}) & p(t_{n-1}) \\ 0 & \cdots & 0 & -p(t_n) & 2hq(t_n) \end{bmatrix}.$$

Then

$$A(h) \begin{bmatrix} u(t_1) \\ u(t_2) \\ \vdots \\ u(t_n) \end{bmatrix} = \begin{bmatrix} f(t_1) \\ f(t_2) \\ \vdots \\ f(t_n) \end{bmatrix} + \mathcal{O}(h^2).$$

#### Example (Cont'd)

In particular, if  $p \equiv 0$  and q < 0 on (0,1), then A(h) is diagonal dominant:

$$\left|\frac{-2}{h^2} + q(t_k)\right| = \frac{2}{h^2} - q(t_k) > \frac{1}{h^2} + \frac{1}{h^2}.$$

Therefore, A(h) is invertible if  $p \equiv 0$  and q < 0.

**Remark**: Let f(t, y, y') = f(t) - q(t)y(t) - p(t)y'(t). Then the condition q < 0 on (0, 1) corresponds to the condition

$$f_y > 0$$

in the statement of the existence theorem. This is an indication why we need this condition in the existence theorem.

We begin with considering the following two-point boundary value problem:

$$\begin{cases} -u''(x) = f(x) & \forall x \in (0,1), \\ u(0) = u(1) = 0, \end{cases}$$
 (D)

where f is a given function in C([0,1]).

Remark: (D) has a unique classical solution.

Let  $v \in C([0,1])$ , piecewise differentiable, and v(0) = v(1) = 0.

Then integration by parts implies that

$$\int_{0}^{1} f(x)v(x) dx = -\int_{0}^{1} u''(x)v(x) dx$$

$$= -u'(x)v(x)\Big|_{x=0}^{x=1} + \int_{0}^{1} u'(x)v'(x) dx$$

$$= \int_{0}^{1} u'(x)v'(x) dx.$$

We begin with considering the following two-point boundary value problem:

$$\begin{cases} -u''(x) = f(x) & \forall x \in (0,1), \\ u(0) = u(1) = 0, \end{cases}$$
 (D)

where f is a given function in C([0,1]).

Remark: (D) has a unique classical solution.

Let  $v \in C([0,1])$ , piecewise differentiable, and v(0) = v(1) = 0.

Then integration by parts implies that

$$\int_0^1 f(x)v(x) dx = -\int_0^1 u''(x)v(x) dx$$

$$= -u'(x)v(x)\Big|_{x=0}^{x=1} + \int_0^1 u'(x)v'(x) dx$$

$$= \int_0^1 u'(x)v'(x) dx.$$

Note that if u is a solution to (D), then u also satisfy the identity

$$\int_0^1 u'(x)v'(x) dx = \int_0^1 f(x)v(x) dx \qquad \forall \ v \in \mathcal{V},$$
 (V)

where  $\mathcal V$  is the collection of all continuous, piecewise differentiable functions on [0,1] that vanish at x=0 and x=1. Therefore, instead of solving for (D), we first look for a function  $\underline{u} \in \mathcal V$  satisfies (V). (V) is called the **variational form** of (D).

We also note that if (V) has a solution u, the function u might not satisfy (D) because u may not be twice differentiable. However, if  $u \in C^2((0,1)) \cap \mathcal{V}$  satisfies (V), then for all  $v \in \mathcal{V}$ ,

$$\int_0^1 [f(x) + u''(x)] v(x) dx = \int_0^1 [f(x)v(x) - u'(x)v'(x)] dx = 0;$$

Note that if u is a solution to (D), then u also satisfy the identity

$$\int_0^1 u'(x)v'(x) dx = \int_0^1 f(x)v(x) dx \qquad \forall \ v \in \mathcal{V},$$
 (V)

where  $\mathcal{V}$  is the collection of all continuous, piecewise differentiable functions on [0,1] that vanish at x=0 and x=1. Therefore, instead of solving for (D), we first look for a function  $\underline{u} \in \mathcal{V}$  satisfies (V). (V) is called the **variational form** of (D).

We also note that if (V) has a solution u, the function u might not satisfy (D) because u may not be twice differentiable. However, if  $u \in C^2((0,1)) \cap \mathcal{V}$  satisfies (V), then for all  $v \in \mathcal{V}$ ,

$$\int_0^1 [f(x) + u''(x)] v(x) dx = \int_0^1 [f(x)v(x) - u'(x)v'(x)] dx = 0;$$

Note that if u is a solution to (D), then u also satisfy the identity

$$\int_0^1 u'(x)v'(x) dx = \int_0^1 f(x)v(x) dx \qquad \forall \ v \in \mathcal{V},$$
 (V)

where  $\mathcal V$  is the collection of all continuous, piecewise differentiable functions on [0,1] that vanish at x=0 and x=1. Therefore, instead of solving for (D), we first look for a function  $\underline{u\in\mathcal V}$  satisfies (V). (V) is called the **variational form** of (D).

We also note that if (V) has a solution u, the function u might not satisfy (D) because u may not be twice differentiable. However, if  $u \in C^2((0,1)) \cap \mathcal{V}$  satisfies (V), then for all  $v \in \mathcal{V}$ ,

$$\int_0^1 [f(x) + u''(x)] v(x) dx = \int_0^1 [f(x)v(x) - u'(x)v'(x)] dx = 0;$$

Note that if u is a solution to (D), then u also satisfy the identity

$$\int_0^1 u'(x)v'(x) dx = \int_0^1 f(x)v(x) dx \qquad \forall \ v \in \mathcal{V},$$
 (V)

where  $\mathcal{V}$  is the collection of all continuous, piecewise differentiable functions on [0,1] that vanish at x=0 and x=1. Therefore, instead of solving for (D), we first look for a function  $\underline{u} \in \mathcal{V}$  satisfies (V). (V) is called the **variational form** of (D).

We also note that if (V) has a solution u, the function u might not satisfy (D) because u may not be twice differentiable. However, if  $u \in C^2((0,1)) \cap \mathcal{V}$  satisfies (V), then for all  $v \in \mathcal{V}$ ,

$$\int_0^1 [f(x) + u''(x)] v(x) dx = \int_0^1 [f(x)v(x) - u'(x)v'(x)] dx = 0;$$

Let  $F: \mathcal{V} \to \mathbb{R}$  be defined by

$$F(v) = \frac{1}{2} \int_0^1 |v'(x)|^2 dx - \int_0^1 f(x)v(x) dx$$

and consider the problem of finding  $u \in \mathcal{V}$  such that

$$F(u) \leqslant F(v) \qquad \forall \ v \in \mathcal{V} \,. \tag{M}$$

In the following, we prove that "if  $u \in V$ , then u satisfies (V) if and only u satisfies (M). Moreover, (V) has at most one solution in V."

Before proceeding, for the purpose of simplifying the notation, we define an "inner product" on V:

$$\langle f, g \rangle \equiv \int_0^1 f(x)g(x) dx.$$

Using this notation, (V) can be rewritten as

$$\langle u', v' \rangle - \langle f, v \rangle = 0 \quad \forall v \in \mathcal{V}.$$

Let  $F: \mathcal{V} \to \mathbb{R}$  be defined by

$$F(v) = \frac{1}{2} \int_0^1 |v'(x)|^2 dx - \int_0^1 f(x)v(x) dx$$

and consider the problem of finding  $u \in \mathcal{V}$  such that

$$F(u) \leqslant F(v) \qquad \forall \ v \in \mathcal{V} \,. \tag{M}$$

In the following, we prove that "if  $u \in \mathcal{V}$ , then u satisfies (V) if and only u satisfies (M). Moreover, (V) has at most one solution in  $\mathcal{V}$ ."

Before proceeding, for the purpose of simplifying the notation, we define an "inner product" on  $\mathcal{V}$ :

$$\langle f, g \rangle \equiv \int_0^1 f(x)g(x) dx.$$

Using this notation, (V) can be rewritten as

$$\langle u', v' \rangle - \langle f, v \rangle = 0 \quad \forall v \in \mathcal{V}.$$



Let  $F: \mathcal{V} \to \mathbb{R}$  be defined by

$$F(v) = \frac{1}{2} \int_0^1 |v'(x)|^2 dx - \int_0^1 f(x)v(x) dx$$

and consider the problem of finding  $u \in \mathcal{V}$  such that

$$F(u) \leqslant F(v) \qquad \forall \ v \in \mathcal{V} \,. \tag{M}$$

In the following, we prove that "if  $u \in \mathcal{V}$ , then u satisfies (V) if and only u satisfies (M). Moreover, (V) has at most one solution in  $\mathcal{V}$ ."

Before proceeding, for the purpose of simplifying the notation, we define an "inner product" on V:

$$\langle f, g \rangle \equiv \int_0^1 f(x)g(x) dx.$$

Using this notation, (V) can be rewritten as

$$\langle u', v' \rangle - \langle f, v \rangle = 0 \quad \forall v \in \mathcal{V}.$$



**1** (V)  $\Rightarrow$  (M): Let u be a solution of problem (V). Let  $v \in \mathcal{V}$  and  $w = v - u \in \mathcal{V}$ . Then v = u + w and

$$F(v) = F(u+w) = \frac{1}{2} \langle (u+w)', (u+w)' \rangle - \langle f, u+w \rangle$$

$$= \frac{1}{2} \langle u', u' \rangle + \langle u', w' \rangle + \frac{1}{2} \langle w', w' \rangle - \langle f, u \rangle - \langle f, w \rangle$$

$$= \frac{1}{2} \langle u', u' \rangle + \frac{1}{2} \langle w', w' \rangle - \langle f, u \rangle \geqslant \frac{1}{2} \langle u', u' \rangle - \langle f, u \rangle = F(u).$$

② (M)  $\Rightarrow$  (V): Let u be a solution of problem (M). Note that if  $v \in \mathcal{V}$ , then  $u + \varepsilon v \in \mathcal{V}$  for all  $\varepsilon > 0$ . Therefore,

$$F(u) \leqslant F(u + \varepsilon v) \qquad \forall \ v \in \mathcal{V} \text{ and } \varepsilon \in \mathbb{R}.$$
  $(\diamond)$ 

Define  $g(\varepsilon) = F(u + \varepsilon v)$ . Then  $(\diamond)$  implies that g attains its global minimum at  $\varepsilon = 0$ . Therefore,

$$0 = g'(0) = \langle u', v' \rangle - \langle f, v \rangle$$



① (V)  $\Rightarrow$  (M): Let u be a solution of problem (V). Let  $v \in \mathcal{V}$  and  $w = v - u \in \mathcal{V}$ . Then v = u + w and

$$F(v) = F(u+w) = \frac{1}{2} \langle (u+w)', (u+w)' \rangle - \langle f, u+w \rangle$$

$$= \frac{1}{2} \langle u', u' \rangle + \langle u', w' \rangle + \frac{1}{2} \langle w', w' \rangle - \langle f, u \rangle - \langle f, w \rangle$$

$$= \frac{1}{2} \langle u', u' \rangle + \frac{1}{2} \langle w', w' \rangle - \langle f, u \rangle \geqslant \frac{1}{2} \langle u', u' \rangle - \langle f, u \rangle = F(u).$$

② (M)  $\Rightarrow$  (V): Let u be a solution of problem (M). Note that if  $v \in \mathcal{V}$ , then  $u + \varepsilon v \in \mathcal{V}$  for all  $\varepsilon > 0$ . Therefore,

$$F(u) \leqslant F(u + \varepsilon v) \qquad \forall \ v \in \mathcal{V} \ \text{and} \ \varepsilon \in \mathbb{R}.$$
 ( $\diamond$ )

Define  $g(\varepsilon)=F(u+\varepsilon v)$ . Then ( $\diamond$ ) implies that g attains its global minimum at  $\varepsilon=0$ . Therefore,

$$0 = \mathbf{g}'(0) = \langle \mathbf{u}', \mathbf{v}' \rangle - \langle \mathbf{f}, \mathbf{v} \rangle$$



(V) has at most one solution in  $\mathcal{V}$ : Suppose that  $u_1, u_2 \in \mathcal{V}$  both satisfy (V). Then

$$\langle u_1', v' \rangle - \langle f, v \rangle = \langle u_2', v' \rangle - \langle f, v \rangle = 0 \qquad \forall v \in \mathcal{V}.$$

Note that  $u_1 - u_2 \in \mathcal{V}$ ; thus letting  $v = u_1 - u_2$  in the equality above we find that

$$\langle u_1', (u_1 - u_2)' \rangle - \langle f, u_1 - u_2 \rangle = 0 , \langle u_2', (u_1 - u_2)' \rangle - \langle f, u_1 - u_2 \rangle = 0 ,$$

Therefore, forming the difference of the two equations above shows that

$$\int_0^1 \left| \left[ u_1(x) - u_2(x) \right]' \right|^2 dx = \left\langle (u_1 - u_2)', (u_1 - u_2)' \right\rangle = 0.$$

The identity above implies that  $\left[u_1(x) - u_2(x)\right]' = 0$ ; thus  $u_1(x) - u_2(x) = 0$  for all  $x \in [0, 1]$  which concludes that  $u_1 = u_2$ .

#### FEM for the model problem with piecewise linear functions:

Construct a finite-dimensional space  $V_h$  (finite element space) as follows: let  $0 = x_0 < x_2 < \cdots < x_M < x_{M+1} = 1$  be a partition of [0,1], and set

- $I_j := [x_{j-1}, x_j], \quad j = 1, 2, \cdots, M+1.$
- $h_j := x_j x_{j-1}, \quad j = 1, 2, \cdots, M+1.$
- $h:=\max_{j=1,2,\cdots,M+1}h_j$ , a measure of how fine the partition is.

Define

$$\mathcal{V}_h := \Big\{ v_h \in \mathcal{V} \ \Big| \ v_h \text{ is linear on each subinterval } I_j, \\ v_h(0) = v_h(1) = 0 \Big\}.$$

Notice that  $V_h \subseteq V$ .



For  $j = 1, 2, \dots, M$ , we define  $\varphi_j \in \mathcal{V}_h$  by

$$\varphi_j(x_i) = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{if } i \neq j. \end{cases}$$

Then every  $f \in \mathcal{V}_h$  can be expressed as

$$f(x) = \sum_{j=1}^{M} f(x_j) \varphi_j(x);$$

thus  $\{\varphi_j\}_{j=1}^M$  is a basis of the finite-dimensional vector space  $\mathcal{V}_h$ .

Two numerical methods for approximating the solution of (D):

- Q Ritz method:
  - Find  $u_h \in \mathcal{V}_h$  such that  $F(u_h) \leqslant F(v_h)$  for all  $v_h \in \mathcal{V}_h$ . (M<sub>h</sub>)
- @ Galerkin method (finite element method):

Find 
$$u_h \in \mathcal{V}_h$$
 such that  $\langle u_h', v_h' \rangle = \langle f, v_h \rangle$  for all  $v_h \in \mathcal{V}_h$ . (V<sub>h</sub>)

Similar to the proof of (M)  $\Leftrightarrow$  (V), one can show that  $(M_h) \Leftrightarrow (V_h)$ .

For  $j = 1, 2, \dots, M$ , we define  $\varphi_j \in \mathcal{V}_h$  by

$$\varphi_j(x_i) = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{if } i \neq j. \end{cases}$$

Then every  $f \in \mathcal{V}_h$  can be expressed as

$$f(x) = \sum_{j=1}^{M} f(x_j) \varphi_j(x);$$

thus  $\{\varphi_j\}_{j=1}^M$  is a basis of the finite-dimensional vector space  $\mathcal{V}_h$ .

Two numerical methods for approximating the solution of (D):

• Ritz method:

Find 
$$u_h \in \mathcal{V}_h$$
 such that  $F(u_h) \leqslant F(v_h)$  for all  $v_h \in \mathcal{V}_h$ .  $(M_h)$ 

② Galerkin method (finite element method):

Find 
$$u_h \in \mathcal{V}_h$$
 such that  $\langle u_h', v_h' \rangle = \langle f, v_h \rangle$  for all  $v_h \in \mathcal{V}_h$ .  $(V_h)$ 

Similar to the proof of (M)  $\Leftrightarrow$  (V), one can show that  $(M_h) \Leftrightarrow (V_h)$ .

Next we focus on solving  $(V_h)$ . We first claim that

Find  $u_h \in \mathcal{V}_h$  such that  $\langle u_h', v_h' \rangle = \langle f, v_h \rangle$  for all  $v_h \in \mathcal{V}_h$ .

 $\Leftrightarrow$  Find  $u_h \in \mathcal{V}_h$  such that  $\langle u_h', \varphi_j' \rangle = \langle f, \varphi_j \rangle$  for all  $1 \leqslant j \leqslant M$ .

#### Proof.

- (⇒) Trivial! (Choosing test function  $v = \varphi_j$ ).
- ( $\Leftarrow$ ) For any  $v_h \in \mathcal{V}_h$ ,  $v_h = \sum\limits_{j=1}^M \eta_j \varphi_j$  for some  $(\eta_1, \eta_2 \cdots, \eta_M) \in \mathbb{R}^M$ . Therefore,

$$\langle u_h', v_h' \rangle = \left\langle u_h', \sum_{j=1}^M \eta_j \varphi_j' \right\rangle = \sum_{j=1}^M \eta_j \langle u_h', \varphi_j' \rangle$$
$$= \sum_{j=1}^M \eta_j \langle f, \varphi_j \rangle = \left\langle f, \sum_{j=1}^M \eta_j \varphi_j \right\rangle = \langle f, v_h \rangle.$$

Note that a solution  $u_h \in \mathcal{V}_h$ , if exists, can be written as  $u_h(x) = \sum_{j=1}^M \xi_j \varphi_j(x)$  for some  $(\xi_1, \xi_2, \cdots, \xi_M) \in \mathbb{R}^M$ . Therefore,

Find 
$$u_h \in \mathcal{V}_h$$
 such that  $\langle u_h', \varphi_j' \rangle = \langle f, \varphi_j \rangle$  for all  $1 \leqslant j \leqslant M$ 

$$\Leftrightarrow \left\langle \sum_{k=1}^{M} \xi_{k} \varphi_{k}', \varphi_{j}' \right\rangle = \left\langle f, \varphi_{j} \right\rangle \text{ for all } 1 \leqslant j \leqslant M$$

$$\Leftrightarrow \sum_{k=1}^{M} \langle \varphi_k', \varphi_j' \rangle \xi_k = \langle f, \varphi_j \rangle \text{ for all } 1 \leqslant j \leqslant M$$

 $\Leftrightarrow A\xi = \mathbf{b}$  for some matrix  $M \times M$  matrix A.

Here, the matrix  $A = [a_{ij}]_{M \times M}$  is defined by  $a_{ij} = \langle \varphi_i', \varphi_j' \rangle$  and is called the **stiffness matrix**, while the vector  $\mathbf{b} = [b_i]_{M \times 1}$  is defined by  $b_i = \langle f, \varphi_i \rangle$  and is called the **load vector**.

$$A\boldsymbol{\xi} = \boldsymbol{b}$$

$$\Leftrightarrow \begin{bmatrix} \langle \varphi_{1}', \varphi_{1}' \rangle & \langle \varphi_{2}', \varphi_{1}' \rangle & \cdots & \langle \varphi_{M}', \varphi_{1}' \rangle \\ \langle \varphi_{1}', \varphi_{2}' \rangle & \langle \varphi_{2}', \varphi_{2}' \rangle & \cdots & \langle \varphi_{M}', \varphi_{2}' \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle \varphi_{1}', \varphi_{M}' \rangle & \langle \varphi_{2}', \varphi_{M}' \rangle & \cdots & \langle \varphi_{M}', \varphi_{M}' \rangle \end{bmatrix} \begin{bmatrix} \xi_{1} \\ \xi_{2} \\ \vdots \\ \xi_{M} \end{bmatrix} = \begin{bmatrix} \langle f, \varphi_{1} \rangle \\ \langle f, \varphi_{2} \rangle \\ \vdots \\ \langle f, \varphi_{M} \rangle \end{bmatrix}.$$

#### Remark:

- A is symmetric.
- **2** A is tri-diagonal:  $\langle \varphi_i', \varphi_i' \rangle = 0$  if |i j| > 1.
- **3** A is positive-definitie: if  $\eta \neq 0$ , then

$$\boldsymbol{\eta}^{\mathrm{T}} A \boldsymbol{\eta} = \left\langle \sum_{j=1}^{M} \eta_{i} \varphi_{i}', \sum_{j=1}^{M} \eta_{j} \varphi_{j}' \right\rangle = \int_{0}^{1} \left| \left( \sum_{j=1}^{M} \eta_{j} \varphi_{j}(\mathbf{x}) \right)' \right|^{2} d\mathbf{x} > 0.$$

**4** Since A is SPD,  $A\boldsymbol{\xi} = \boldsymbol{b}$  has a unique solution.



For  $j = 1, 2, \dots, M$ , we have

$$\begin{split} \langle \varphi'_{j}, \varphi'_{j} \rangle &= \int_{x_{j-1}}^{x_{j}} \varphi'_{j}(x)^{2} dx + \int_{x_{j}}^{x_{j+1}} \varphi'_{j}(x)^{2} dx \\ &= \int_{x_{j-1}}^{x_{j}} \frac{1}{h_{j}^{2}} dx + \int_{x_{j}}^{x_{j+1}} \frac{1}{h_{j+1}^{2}} dx = \frac{1}{h_{j}} + \frac{1}{h_{j+1}}, \\ \langle \varphi'_{j}, \varphi'_{j-1} \rangle &= \langle \varphi'_{j-1}, \varphi'_{j} \rangle = - \int_{x_{j-1}}^{x_{j}} \frac{1}{h_{j}^{2}} dx = -\frac{1}{h_{j}}. \end{split}$$

For uniform partition:  $h_j = h = \frac{1-0}{M+1}$ . Then  $A\boldsymbol{\xi} = \boldsymbol{b}$  becomes

$$\frac{1}{h} \begin{bmatrix}
2 & -1 & 0 & \cdots & 0 \\
-1 & 2 & -1 & \ddots & 0 \\
0 & \ddots & \ddots & \ddots & 0 \\
\vdots & \ddots & -1 & 2 & -1 \\
0 & \cdots & 0 & -1 & 2
\end{bmatrix}
\begin{bmatrix}
\xi_1 \\ \xi_2 \\ \vdots \\ \xi_M
\end{bmatrix} = \begin{bmatrix}
\langle f, \varphi_1 \rangle \\ \langle f, \varphi_2 \rangle \\ \vdots \\ \langle f, \varphi_M \rangle
\end{bmatrix}$$

For  $j = 1, 2, \dots, M$ , we have

$$\begin{split} \langle \varphi'_{j}, \varphi'_{j} \rangle &= \int_{x_{j-1}}^{x_{j}} \varphi'_{j}(x)^{2} dx + \int_{x_{j}}^{x_{j+1}} \varphi'_{j}(x)^{2} dx \\ &= \int_{x_{j-1}}^{x_{j}} \frac{1}{h_{j}^{2}} dx + \int_{x_{j}}^{x_{j+1}} \frac{1}{h_{j+1}^{2}} dx = \frac{1}{h_{j}} + \frac{1}{h_{j+1}}, \\ \langle \varphi'_{j}, \varphi'_{j-1} \rangle &= \langle \varphi'_{j-1}, \varphi'_{j} \rangle = - \int_{x_{j-1}}^{x_{j}} \frac{1}{h_{j}^{2}} dx = -\frac{1}{h_{j}}. \end{split}$$

For uniform partition:  $h_j=h=rac{1-0}{M+1}$ . Then  $Aoldsymbol{\xi}=oldsymbol{b}$  becomes

$$\frac{1}{h} \begin{bmatrix} 2 & -1 & 0 & \cdots & 0 \\ -1 & 2 & -1 & \ddots & 0 \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & -1 & 2 & -1 \\ 0 & \cdots & 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} \xi_1 \\ \xi_2 \\ \vdots \\ \xi_M \end{bmatrix} = \begin{bmatrix} \langle f, \varphi_1 \rangle \\ \langle f, \varphi_2 \rangle \\ \vdots \\ \xi_M \end{bmatrix}.$$

For  $j = 1, 2, \dots, M$ , we have

$$\begin{split} \langle \varphi'_{j}, \varphi'_{j} \rangle &= \int_{x_{j-1}}^{x_{j}} \varphi'_{j}(x)^{2} dx + \int_{x_{j}}^{x_{j+1}} \varphi'_{j}(x)^{2} dx \\ &= \int_{x_{j-1}}^{x_{j}} \frac{1}{h_{j}^{2}} dx + \int_{x_{j}}^{x_{j+1}} \frac{1}{h_{j+1}^{2}} dx = \frac{1}{h_{j}} + \frac{1}{h_{j+1}}, \\ \langle \varphi'_{j}, \varphi'_{j-1} \rangle &= \langle \varphi'_{j-1}, \varphi'_{j} \rangle = - \int_{x_{j-1}}^{x_{j}} \frac{1}{h_{j}^{2}} dx = -\frac{1}{h_{j}}. \end{split}$$

For uniform partition:  $h_j = h = \frac{1-0}{M+1}$ . Then  $A\boldsymbol{\xi} = \boldsymbol{b}$  becomes

$$\frac{1}{h} \begin{bmatrix} 2 & -1 & 0 & \cdots & 0 \\ -1 & 2 & -1 & \ddots & 0 \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & -1 & 2 & -1 \\ 0 & \cdots & 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} u_h(\mathbf{x}_1) \\ u_h(\mathbf{x}_2) \\ \vdots \\ \vdots \\ u_h(\mathbf{x}_M) \end{bmatrix} = \begin{bmatrix} \langle f, \varphi_1 \rangle \\ \langle f, \varphi_2 \rangle \\ \vdots \\ \vdots \\ \langle f, \varphi_M \rangle \end{bmatrix}.$$

$$\frac{1}{h} \begin{bmatrix}
2 & -1 & 0 & \cdots & 0 \\
-1 & 2 & -1 & \ddots & 0 \\
0 & \ddots & \ddots & \ddots & 0 \\
\vdots & \ddots & -1 & 2 & -1 \\
0 & \cdots & 0 & -1 & 2
\end{bmatrix}
\begin{bmatrix}
u_h(x_1) \\
u_h(x_2) \\
\vdots \\
\vdots \\
u_h(x_M)
\end{bmatrix} = \begin{bmatrix}
\langle f, \varphi_1 \rangle \\
\langle f, \varphi_2 \rangle \\
\vdots \\
\vdots \\
\langle f, \varphi_M \rangle
\end{bmatrix}.$$

Taking into account that  $u_h(x_0) = u_h(x_{M+1}) = 0$ , the system above shows that

$$-\frac{1}{h^2} \big[ u_h(x_{j-1}) - 2u_h(x_j) + u_h(x_{j+1}) \big] = \frac{1}{h} \langle f, \varphi_j \rangle \qquad \forall \, 1 \leqslant j \leqslant M.$$

Since we expect that  $u_h$  approximates the solution u, the equality above implies that

$$-\frac{1}{h^2}\left[u(x_j-h)-2u(x_j)+u(x_j+h)\right]\approx \frac{1}{h}\int_{x_j-h}^{x_j+h}f(x)\varphi_j(x)\ dx\ \forall\ 1\leqslant j\leqslant M.$$

$$\frac{1}{h} \begin{bmatrix}
2 & -1 & 0 & \cdots & 0 \\
-1 & 2 & -1 & \ddots & 0 \\
0 & \ddots & \ddots & \ddots & 0 \\
\vdots & \ddots & -1 & 2 & -1 \\
0 & \cdots & 0 & -1 & 2
\end{bmatrix}
\begin{bmatrix}
u_h(x_1) \\
u_h(x_2) \\
\vdots \\
\vdots \\
u_h(x_M)
\end{bmatrix} = \begin{bmatrix}
\langle f, \varphi_1 \rangle \\
\langle f, \varphi_2 \rangle \\
\vdots \\
\vdots \\
\langle f, \varphi_M \rangle
\end{bmatrix}.$$

Taking into account that  $u_h(x_0) = u_h(x_{M+1}) = 0$ , the system above shows that

$$-\frac{1}{h^2} \big[ u_h(x_{j-1}) - 2u_h(x_j) + u_h(x_{j+1}) \big] = \frac{1}{h} \langle f, \varphi_j \rangle \qquad \forall \, 1 \leqslant j \leqslant M.$$

Since we expect that  $u_h$  approximates the solution u, the equality above implies that

$$-\frac{1}{h^2}\big[u(x_j-h)-2u(x_j)+u(x_j+h)\big]\approx \frac{1}{h}\int_{x_j-h}^{x_j+h}f(x)\varphi_j(x)\ dx\ \ \forall\ 1\leqslant j\leqslant M\ .$$

Recall that if  $u \in C^4([x_j - h, x_j + h])$ ,

$$u''(x_j) = \frac{1}{h^2} \left[ u(x_j + h) - 2u(x_j) + u(x_j - h) \right] - \frac{1}{12} h^2 u^{(4)}(\xi_j)$$

for some  $\xi_j \in (x_j - h, x_j + h)$ . Moreover, if  $f \in C^2([x_j - h, x_j + h])$ ,

by Taylor's Theorem there exists  $\eta_j \in (x_j - h, x_j + h)$  such that

$$f(x) = f(x_j) + f'(x_j)(x - x_j) + \frac{f''(\eta_j)}{2}(x - x_j)^2;$$

$$\frac{1}{h} \int_{x_{j-h}}^{x_{j+h}} f(x) \varphi_{j}(x) dx 
= \frac{1}{h} \int_{x_{j-h}}^{x_{j+h}} \left[ f(x_{j}) + f'(x_{j})(x - x_{j}) + \frac{f''(\eta_{j})}{2} (x - x_{j})^{2} \right] \varphi(x) dx 
= f(x_{j}) + \frac{1}{2h} \int_{x_{j-h}}^{x_{j+h}} f''(\eta_{j})(x - x_{j})^{2} \varphi(x) dx$$

Recall that if  $u \in C^4([x_j - h, x_j + h])$ ,

$$u''(x_j) = \frac{1}{h^2} \left[ u(x_j + h) - 2u(x_j) + u(x_j - h) \right] - \frac{1}{12} h^2 u^{(4)}(\xi_j)$$

for some  $\xi_j \in (x_j - h, x_j + h)$ . Moreover, if  $f \in C^2([x_j - h, x_j + h])$ ,

by Taylor's Theorem there exists  $\eta_j \in (x_j - h, x_j + h)$  such that

$$f(x) = f(x_j) + f'(x_j)(x - x_j) + \frac{f''(\eta_j)}{2}(x - x_j)^2;$$

$$\frac{1}{h} \int_{x_{j-h}}^{x_{j+h}} f(x) \varphi_{j}(x) dx 
= \frac{1}{h} \int_{x_{j-h}}^{x_{j+h}} \left[ f(x_{j}) + f'(x_{j})(x - x_{j}) + \frac{f''(\eta_{j})}{2} (x - x_{j})^{2} \right] \varphi(x) dx 
= f(x_{j}) + \frac{1}{2h} f''(\zeta_{j}) \int_{x_{j-h}}^{x_{j+h}} (x - x_{j})^{2} \varphi(x) dx$$

Recall that if  $u \in C^4([x_j - h, x_j + h])$ ,

$$u''(x_j) = \frac{1}{h^2} \left[ u(x_j + h) - 2u(x_j) + u(x_j - h) \right] - \frac{1}{12} h^2 u^{(4)}(\xi_j)$$

for some  $\xi_j \in (x_j - h, x_j + h)$ . Moreover, if  $f \in C^2([x_j - h, x_j + h])$ ,

by Taylor's Theorem there exists  $\eta_j \in (x_j - h, x_j + h)$  such that

$$f(x) = f(x_j) + f'(x_j)(x - x_j) + \frac{f''(\eta_j)}{2}(x - x_j)^2;$$

$$\frac{1}{h} \int_{x_{j-h}}^{x_{j+h}} f(x) \varphi_{j}(x) dx 
= \frac{1}{h} \int_{x_{j-h}}^{x_{j+h}} \left[ f(x_{j}) + f'(x_{j})(x - x_{j}) + \frac{f''(\eta_{j})}{2} (x - x_{j})^{2} \right] \varphi(x) dx 
= f(x_{j}) + \frac{1}{2h} f''(\zeta_{j}) \cdot \frac{h^{3}}{2}$$

Recall that if  $u \in C^4([x_j - h, x_j + h])$ ,

$$u''(x_j) = \frac{1}{h^2} \left[ u(x_j + h) - 2u(x_j) + u(x_j - h) \right] - \frac{1}{12} h^2 u^{(4)}(\xi_j)$$

for some  $\xi_j \in (x_j - h, x_j + h)$ . Moreover, if  $f \in C^2([x_j - h, x_j + h])$ ,

by Taylor's Theorem there exists  $\eta_j \in (x_j - h, x_j + h)$  such that

$$f(x) = f(x_j) + f'(x_j)(x - x_j) + \frac{f''(\eta_j)}{2}(x - x_j)^2;$$

$$\frac{1}{h} \int_{x_{j-h}}^{x_{j+h}} f(x) \varphi_{j}(x) dx 
= \frac{1}{h} \int_{x_{j-h}}^{x_{j+h}} \left[ f(x_{j}) + f'(x_{j})(x - x_{j}) + \frac{f''(\eta_{j})}{2} (x - x_{j})^{2} \right] \varphi(x) dx 
= f(x_{j}) + \frac{h^{2}}{4} f''(\zeta_{j}) \text{ for some } \zeta_{j} \in (x_{j} - h, x_{j} + h).$$

Therefore, we conclude that if  $f \in C^2([0,1])$  (so that  $u \in C^4([0,1])$ ), the finite element method introduced above indeed is a second order approximation of the equation

$$-u''(x) = f(x)$$

at each node  $x_j$ .

#### FDM different from the FEM methods in few aspects:

- In the FDM methods, the discretization of the domain is done as a set of nodes at which the results are determined, while in the FEM method the results are known in every point of the domain as the approximation is done with functions defined on small triangular (or quadrilateral) areas in 2D.
- ② Because of that, the algorithms used in FDM require generally less computational power to solve the equations, but it results also in less "refines" results (only at nodes).

#### The main difference between FEM and FDM (in simple terms):

- FDM is an older method than FEM that requires less computational power but is also less accurate in some cases where higher-order accuracy is required.
- PEM permits to get a higher order of accuracy, but requires more computational power and is also more exigent on the quality of the mesh.

#### Computer project

Consider the following one-dimensional convection-diffusion problem:

$$\begin{cases} -\varepsilon u''(x) + u'(x) = 0 & \text{for } x \in (0, 1), \\ u(0) = 1, \ u(1) = 0. \end{cases}$$
 (\*)

Write the computer codes for numerical solution of problem (\*) by using the finite difference methods on the uniform mesh of [0,1] with mesh size h:

- Replace  $u''(x_i) \approx \frac{U_{i+1} 2U_i + U_{i-1}}{h^2}$  and  $u'(x_i) \approx \frac{U_{i+1} U_{i-1}}{2h}$  and consider  $(\varepsilon, h) = (0.01, 0.1)$ ,  $(\varepsilon, h) = (0.01, 0.01)$ . Plot  $u_h$ .
- Replace  $u''(x_i) \approx \frac{U_{i+1} 2U_i + U_{i-1}}{h^2}$  and  $u'(x_i) \approx \frac{U_i U_{i-1}}{h}$  (upwinding) and consider  $(\varepsilon, h) = (0.01, 0.1)$ ,  $(\varepsilon, h) = (0.01, 0.01)$ . Plot  $u_h$ .