

To think and write portable programs we need an abstract model of their computational environment. Faithful models do exist, but they reveal that environment to be too diverse, forcing portable programmers to bloat even the simplest concrete tasks into abstract monsters. We need something simple or, if not so simple, not so capriciously complex.

----- William M. Kahan

(born June 5, 1933 in Toronto, Canada)

William Kahan is a mathematician and computer scientist who received the Turing Award in 1989 for “his fundamental contributions to numerical analysis”. Kahan was the primary architect behind the IEEE 754-1985 standard for floating-point computation. He has been called “The Father of Floating Point”.

Possible project topic: Go to Prof. Kahan’s personal webpage

<http://www.cs.berkeley.edu/~wkahan/>

and read his papers and give a lecture on his results. The interesting ones are “MATLAB’s loss is nobody’s gain”, “Beastly numbers”, etc.

Question:

How to calculate $f'(x)$ when the analytic formula is not available?

Numerical differentiation (Finite difference)

Goal: to calculate $f'(x)$.

A first order method for $f'(x)$

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = f'(x)$$

$$\implies \frac{f(x+h) - f(x)}{h} \approx f'(x) \quad \text{for small } h$$

That is, for small h , we can use $\frac{f(x+h) - f(x)}{h}$ to approximate $f'(x)$.

More precisely, we can write

$$\underbrace{\frac{f(x+h) - f(x)}{h}}_{\text{Numerical approximation}} = \underbrace{f'(x)}_{\text{Exact value}} + \underbrace{e(h)}_{\text{Discretization error}}$$

$$e(h) \rightarrow 0 \quad \text{as } h \rightarrow 0$$

Let us find out the order of $e(h)$. Taylor expansion of $f(x+h)$ around x :

$$f(x+h) = f(x) + f'(x)h + \frac{f''(x)}{2}h^2 + O(h^3)$$

$$\implies \frac{f(x+h) - f(x)}{h} = f'(x) + \frac{f''(x)}{2}h + O(h^2)$$

$$\implies e(h) = \frac{f''(x)}{2}h + O(h^2) = O(h)$$

$$\implies \frac{f(x+h) - f(x)}{h} \text{ is a first order method (forward difference) for calculating } f'(x).$$

“Big Oh” and “small oh” notation: Write $f = O(g)$ if $g(x) > 0$ for $x \in \mathbb{R}$ sufficiently large and $\frac{f(x)}{g(x)}$ is bounded for x sufficiently large. Write $f = o(g)$ if $\frac{f}{g}$ goes to zero as x goes to

$+\infty$. Also write $f \sim g$ (read f is asymptotic to g) if $\frac{f}{g} \rightarrow 1$ as $x \rightarrow \infty$. For example,

$$x^2 + x = O(x^2),$$

$$x^2 + x \sim x^2$$

$$e^{\sqrt{\ln x}} = o(x)$$

The proof goes like this. Since $\frac{x^2 + x}{x^2} = 1 + \frac{1}{x} \rightarrow 1$, as $x \rightarrow \infty$, the first two follow.

The third one is a little harder to prove. We note that

$$e^{\ln x} = x \implies \frac{e^{\sqrt{\ln x}}}{x} = \frac{e^{\sqrt{\ln x}}}{e^{\ln x}} = e^{\sqrt{\ln x} - \ln x}$$

Since $\ln x \rightarrow \infty$ as $x \rightarrow \infty$, $\frac{\sqrt{\ln x}}{\frac{1}{2} \ln x} \rightarrow 0$ as $x \rightarrow \infty$ (Apply L'Hopital's Rule to see this).

For x large enough, $\sqrt{\ln x} \leq \frac{1}{2} \ln x$.

Hence, for x large, $0 \leq \frac{e^{\sqrt{\ln x}}}{x} \leq \frac{e^{\frac{1}{2} \ln x}}{x} = \frac{e^{\ln x^{1/2}}}{x} = \frac{x^{1/2}}{x} = \frac{1}{\sqrt{x}} \rightarrow 0$, as $x \rightarrow \infty$.

By the squeeze theorem, $\frac{e^{\sqrt{\ln x}}}{x} \rightarrow 0$ as $x \rightarrow \infty$.

Note: (1) For a first order method, if we reduce h by a factor of 2, then $e(h)$ is approximately reduced by a factor of 2.

(2) The method is accurate (exact) when $f(x) = 1$ or $f(x) = x$.

A second order method for $f'(x)$

$$\frac{f(x+h) - f(x-h)}{2h} = f'(x) + e(h)$$

Taylor expansion around x

$$f(x+h) = f(x) + f'(x)h + \frac{f''(x)}{2}h^2 + \frac{f'''(x)}{6}h^3 + O(h^4)$$

$$f(x-h) = f(x) - f'(x)h + \frac{f''(x)}{2}h^2 - \frac{f'''(x)}{6}h^3 + O(h^4)$$

$$\implies \frac{f(x+h) - f(x-h)}{2h} = f'(x) + \frac{f'''(x)}{6}h^2 + O(h^3)$$

$$\implies e(h) = \frac{f'''(x)}{6}h^2 + O(h^3) = O(h^2)$$

$\implies \frac{f(x+h) - f(x-h)}{2h}$ is a second order method (central difference approximation) for calculating $f'(x)$.

Note: (1) For a second order method, if we reduce h by a factor of 2, then $e(h)$ is approximately reduced by a factor of $2^2 = 4$.

(2) The method is accurate (exact) for $f(x) = 1, x, x^2$. Let us check $f(x) = x^2$ case.

$$\frac{f(x+h) - f(x-h)}{2h} = \frac{(x+h)^2 - (x-h)^2}{2h} = \frac{4xh}{2h} = 2x$$

$$f'(x) = (x^2)' = 2x$$

A fourth order method for $f'(x)$

$$\frac{-f(x+2h) + 8f(x+h) - 8f(x-h) + f(x-2h)}{12h} = f'(x) + e(h)$$

$$e(h) = \frac{-f^{(5)}(x)}{30}h^4 + O(h^5) = O(h^4)$$

Note: (1) For a fourth order method, if we reduce h by a factor of 2, then $e(h)$ is approximately reduced by a factor of 16.

(2) The method is accurate for $f(x) = 1, x, x^2, x^3, x^4$.

Other finite difference methods:

$$f'(x) = \frac{-f(x+2h) + 4f(x+h) - 3f(x)}{2h} + \underbrace{\frac{h^2}{3} f'''(x) + \dots}_{O(h^2)}$$

$$f'(x) = \frac{-f(x+2h) + 4f(x+h) - 3f(x)}{2h} + O(h^2) \text{ (non-central finite difference approximation)}$$

For higher order derivatives, we have

$$f''(x) = \frac{f(x+h) - 2f(x) + f(x-h)}{h^2} + O(h^2)$$

$$f''(x) = \frac{f(x+2h) - 2f(x+h) + 2f(x-h) - f(x-2h)}{2h^3} + O(h^2)$$

$$f^{(4)}(x) = \frac{f(x+2h) - 4f(x+h) + 6f(x) - 4f(x-h) + f(x-2h)}{h^4} + O(h^2)$$

Total error in numerical differentiation

Consider the first order method

$$\underbrace{\frac{f(x+h) - f(x)}{h}}_{\substack{\text{Numerical value} \\ \text{obtained with} \\ \text{infinite precision}}} = \underbrace{f'(x)}_{\substack{\text{Exact} \\ \text{value}}} + \underbrace{e(h)}_{\substack{\text{Discretization} \\ \text{error}}}, \quad e(h) = O(h)$$

$\frac{f(x+h) - f(x)}{h}$ is NOT the numerical value we get from a computer!

The numerical value we get from a computer is

$$\frac{\text{fl}(f(x+h)) - \text{fl}(f(x))}{h}$$

where

$$\text{fl}(f(x+h)) = f(x+h)(1 + \varepsilon_1), \quad |\varepsilon_1| \sim 10^{-16} \text{ (IEEE double precision representation: machine precision)}$$

$$\text{fl}(f(x)) = f(x)(1 + \varepsilon_2), \quad |\varepsilon_2| \sim 10^{-16}$$

Thus, we have

$$\begin{aligned} & \frac{\text{fl}(f(x+h)) - \text{fl}(f(x))}{h} \\ &= \frac{f(x+h)(1 + \varepsilon_1) - f(x)(1 + \varepsilon_2)}{h} \end{aligned}$$

$$\begin{aligned}
 &= \frac{f(x+h) - f(x)}{h} + \frac{f(x+h)\varepsilon_1 - f(x)\varepsilon_2}{h} \\
 &= f'(x) + \underbrace{e(h)}_{\text{Discretization error}} + \underbrace{\frac{(f(x+h) + f(x))\varepsilon_3}{h}}_{\text{Effect of round-off error}}, \quad |\varepsilon_3| \sim 10^{-16}
 \end{aligned}$$

The discretization error:

$$e(h) \approx Ch \rightarrow 0 \quad \text{as } h \rightarrow 0$$

The effect of round-off error:

$$|f(x+h) + f(x)| \frac{|\varepsilon_3|}{h} \sim |f(x+h) + f(x)| \frac{10^{-16}}{h} \rightarrow \infty \quad \text{as } h \rightarrow 0$$

The total error is defined as

$$\begin{aligned}
 E_T(h) &= \left| \underbrace{\frac{\text{fl}(f(x+h)) - \text{fl}(f(x))}{h}}_{\substack{\text{Numerical value} \\ \text{obtained with} \\ \text{finite precision}}} - \underbrace{f'(x)}_{\substack{\text{Exact} \\ \text{value}}} \right| \\
 &= \left| e(h) + \frac{(f(x+h) + f(x))\varepsilon_3}{h} \right| \\
 &\leq Ch + |f(x+h) + f(x)| \cdot \frac{|\varepsilon_3|}{h}
 \end{aligned}$$

Consider the simple situation where the error is given by

$$E_T(h) = h + \frac{10^{-16}}{h} \Rightarrow E_T'(h) = 1 - \frac{10^{-16}}{h^2} = 0 \Rightarrow h^2 = 10^{-16} \Rightarrow h = 10^{-8}$$

At $h = 10^{-8}$, $E_T(h)$ attains the minimum.

$$\min_h E_T(h) = 10^{-8}$$

(Draw the graph of $E_T(h) = h + \frac{10^{-16}}{h}$)

Note: $\min_h E_T(h)$ is the minimum total error we can achieve using the numerical method.

For the second order method, we have

$$E_T(h) \leq Ch^2 + \frac{|f(x+h) + f(x-h)|}{2} \cdot \frac{|\varepsilon_3|}{h}$$

Consider the simple situation that all coefficients are one and $|\varepsilon_3| = 10^{-16}$.

$$E_T(h) = h^2 + \frac{10^{-16}}{h} \Rightarrow E_T'(h) = 2h - \frac{10^{-16}}{h^2} = 0 \Rightarrow 2h^3 = 10^{-16} \Rightarrow h \approx 10^{-16/3} = 10^{-5.3}$$

At $h = 10^{-5.3}$, $E_T(h)$ attains the minimum

$$\min_h E_T(h) = 10^{-10.7}$$

For the fourth order method, ...

$$E_T(h) = h^4 + \frac{10^{-16}}{h} \Rightarrow E_T'(h) = 4h^3 - \frac{10^{-16}}{h^2} = 0 \Rightarrow h = 10^{-3.2}$$

At $h = 10^{-3.2}$, $E_T(h)$ attains the minimum

$$\min_h E_T(h) = 10^{-12.8}$$

Note: first_ord.m: $h_{\min} = 10^{-8}$, $\min_h E_T(h) = 10^{-8}$

second_ord.m: $h_{\min} = 10^{-5.3}$, $\min_h E_T(h) = 10^{-10.7}$

fourth_ord.m: $h_{\min} = 10^{-3.2}$, $\min_h E_T(h) = 10^{-12.8}$

For a higher order method, the minimum total error is smaller and the minimum total error is achieved at a larger value of h (practically cheaper). This is why higher order methods are more desirable in practice.

See Codes/Total_error/first_ord.m, second_ord.m and fourth_ord.m.

Numerical integration

The definite integral $I = \int_a^b f(x) dx$ is defined in Calculus as a limit of what are called Riemann

sums. It is then proved that $I = \int_a^b f(x) dx = F(b) - F(a)$ where $F(x)$ is any antiderivative of

$f(x)$; this is the Fundamental Theorem of Calculus. Many integrals can be evaluated by using this formula, and a significant portion of Calculus books is devoted to this approach.

Nonetheless, most integrals cannot be evaluated by this formula because most integrands $f(x)$ do not have antiderivatives expressible in terms of elementary functions. Examples of such integrals are

$$\int_0^1 e^{-x^2} dx, \quad \int_0^\pi x^\pi \sin(\sqrt{x}) dx$$

Other methods are needed for evaluating such integrals.

Goal: To calculate (approximate) $\int_a^b f(x)dx$

We will introduce two of the oldest and most popular numerical methods: the trapezoidal rule and Simpson's rule.

Strategy: Divide $[a, b]$ into N subintervals of the size $h = \frac{b-a}{N}$.

Let $x_i = a + ih$, $i = 0, 1, 2, \dots, N$ (the points x_0, x_1, \dots, x_N are called the numerical integration node points)

(Draw the real axis to show $[a, b]$ and $x_i = a + ih$, $i = 0, 1, 2, \dots, N$)

$$\int_a^b f(x)dx = \sum_{i=1}^N \int_{x_{i-1}}^{x_i} f(x)dx$$

To approximate $\int_a^b f(x)dx$, we only need to approximate $\int_{x_{i-1}}^{x_i} f(x)dx$.

Trapezoidal rule:

$$\underbrace{\frac{h}{2}(f_{i-1} + f_i)}_{\text{Numerical approximation}} = \underbrace{\int_{x_{i-1}}^{x_i} f(x)dx}_{\text{Exact value}} + \underbrace{e_i(h)}_{\text{Error}},$$

$$f_i = f(x_i)$$

$$e_i(h) = O(h^3)$$

We will derive the trapezoidal rule later.

Composite trapezoidal rule:

Sum from $i = 1$ to $i = N$, we obtain

$$\begin{aligned} \frac{h}{2} \sum_{i=1}^N (f_{i-1} + f_i) &= \sum_{i=1}^N \int_{x_{i-1}}^{x_i} f(x)dx + \sum_{i=1}^N e_i(h) \\ \implies \underbrace{\frac{h}{2} \left(f_0 + f_N + 2 \sum_{i=1}^{N-1} f_i \right)}_{\substack{\text{Numerical} \\ \text{approximation} \\ T_N(f)}} &= \underbrace{\int_a^b f(x)dx}_{\text{Exact value}} + \underbrace{E(h)}_{\text{Error}}, \end{aligned}$$

$$E(h) = \sum_{i=1}^N e_i(h) = \sum_{i=1}^N O(h^3) = N O(h^3) = O(h^2) \text{ where we have used the fact that } Nh = (b-a).$$

When N is doubled, h is halved, and the term h^2 decreases by a factor of 4.

To illustrate how functions values are reused when N is doubled, consider $T_2(f), T_4(f)$:

$$T_2(f) = h \left[\frac{f(x_0)}{2} + f(x_1) + \frac{f(x_2)}{2} \right] = \frac{h}{2} [f_0 + f_2 + 2f_1],$$

$$h = \frac{b-a}{2}, x_0 = a, x_1 = \frac{a+b}{2}, x_2 = b$$

Also,

$$T_4(f) = h \left[\frac{f(x_0)}{2} + f(x_1) + f(x_2) + f(x_3) + \frac{f(x_4)}{2} \right] = \frac{h}{2} [f_0 + f_4 + 2(f_1 + f_2 + f_3)],$$

$$h = \frac{b-a}{4}, x_0 = a, x_1 = \frac{3a+b}{4}, x_2 = \frac{a+b}{2}, x_3 = \frac{a+3b}{4}, x_4 = b$$

Comparing T_2, T_4 , we can see that only $f(x_1), f(x_3)$ need to be evaluated for T_4 , as the other function values are known from T_2 .

The central idea behind most formulas for approximating $\int_{x_{i-1}}^{x_i} f(x) dx$ is to replace $f(x)$ by an approximating function whose integral can be evaluated.

If we approximate $f(x)$ by the linear polynomial

$$P_1(x) = \frac{(x_i - x)f(x_{i-1}) + (x - x_{i-1})f(x_i)}{x_i - x_{i-1}}$$

which interpolates $f(x)$ at x_{i-1} and x_i (i.e. $P_1(x_{i-1}) = f(x_{i-1}), P_1(x_i) = f(x_i)$). The integral of $P_1(x)$ over $[x_{i-1}, x_i]$ is the area of the trapezoid (draw a figure) and it is given by

$$\int_{x_{i-1}}^{x_i} P_1(x) dx = \int_{x_{i-1}}^{x_i} \frac{(x_i - x)f(x_{i-1}) + (x - x_{i-1})f(x_i)}{x_i - x_{i-1}} dx = (x_i - x_{i-1}) \left[\frac{f(x_{i-1}) + f(x_i)}{2} \right] \text{ (trapezoidal}$$

rule). This is a good approximation if $f(x)$ is almost linear on $[x_{i-1}, x_i]$.

To improve the above approach, use quadratic interpolation to approximate $f(x)$ on $[x_{i-1}, x_i]$. Let $P_2(x)$ be the quadratic polynomial that interpolates $f(x)$ at $x_{i-1}, x_{i-1/2}, x_i$. Using this to

approximate $\int_{x_{i-1}}^{x_i} f(x) dx$, we get

$$\int_{x_{i-1}}^{x_i} f(x) dx \approx \int_{x_{i-1}}^{x_i} P_2(x) dx$$

$$= \int_{x_{i-1}}^{x_i} \left[\frac{(x-x_{i-1/2})(x-x_i)}{(x_{i-1}-x_{i-1/2})(x_{i-1}-x_i)} f(x_{i-1}) + \frac{(x-x_{i-1})(x-x_i)}{(x_{i-1/2}-x_{i-1})(x_{i-1/2}-x_i)} f(x_{i-1/2}) + \frac{(x-x_{i-1/2})(x-x_{i-1})}{(x_i-x_{i-1/2})(x_i-x_{i-1})} f(x_i) \right] dx$$

The complete evaluation of this integral yields

$$\int_{x_{i-1}}^{x_i} f(x) dx \approx \frac{x_i - x_{i-1}}{6} [f(x_{i-1}) + 4f(x_{i-1/2}) + f(x_i)] \quad (\text{Simpson's rule})$$

Simpson's rule:

$$\underbrace{\frac{h}{6}(f_{i-1} + 4f_{i-1/2} + f_i)}_{\text{Numerical approximation}} = \underbrace{\int_{x_{i-1}}^{x_i} f(x) dx}_{\text{Exact value}} + \underbrace{e_i(h)}_{\text{Error}},$$

$$f_i = f(x_i), \quad f_{i-1/2} = f(x_{i-1/2}), \quad x_{i-1/2} = a + \left(i - \frac{1}{2}\right)h$$

$$e_i(h) = O(h^5)$$

We will derive the Simpson's rule later.

Composite Simpson's rule:

Sum from $i = 1$ to $i = N$, we obtain

$$\frac{h}{6} \sum_{i=1}^N (f_{i-1} + 4f_{i-1/2} + f_i) = \sum_{i=1}^N \int_{x_{i-1}}^{x_i} f(x) dx + \sum_{i=1}^N e_i(h)$$

$$\implies \underbrace{\frac{h}{6} \left(f_0 + f_N + 2 \sum_{i=1}^{N-1} f_i + 4 \sum_{i=1}^N f_{i-1/2} \right)}_{\text{Numerical approximation}} = \underbrace{\int_a^b f(x) dx}_{\text{Exact value}} + \underbrace{E(h)}_{\text{Error}},$$

$$E(h) = \sum_{i=1}^N e_i(h) = \sum_{i=1}^N O(h^5) = N O(h^5) = O(h^4)$$

Simpson's rule has been among the most popular numerical integration methods for more than two centuries.

Go through Codes\Num_integration\trapezoid.m which uses the composite trapezoidal rule to compute $\int_a^b f(x) dx$ with $a = 0.5, b = 2, N = 2^5 = 32, f(x) = \sin x$. The exact solution is known from calculus:

$$\int_a^b f(x) dx = \int_{0.5}^2 \sin x dx = -\cos x \Big|_{x=0.5}^{x=2} = \cos 0.5 - \cos 2 = 1.293729398437515$$

The absolute error from this computation is about 2.37×10^{-4} . You can try the codes with different values of N and see how the error is reduced.

Go through Code\Num_integration\simpson.m which uses the composite Simpson rule to solve the above problem. The absolute error now is reduced to 2.17×10^{-9} .

Talk about homework assignment 3.

Problem 2: use the exact value of $\int_a^b \sin x dx$ to compute the error. You can plot the error vs h or N.

Problem 3: compute $g(s) = \int_0^2 \left[\sqrt{1 + e^{-3\cos(sx)}} - 1.5 \right] dx$ for $s \in [0, 4]$ with the Composite Simpson's rule.

Note: MATLAB has built-in functions for numerical integration.

TRAPZ Trapezoidal numerical integration.

QUAD Numerically evaluate integral, adaptive Simpson quadrature.

QUADL Numerically evaluate integral, adaptive Lobatto quadrature.

QUAD2D Numerically evaluate double integral over a planar region.

DBLQUAD Numerically evaluate double integral over a rectangle.

TRIPLEQUAD Numerically evaluate triple integral.

For the model problem, we can try the following:

```
>>x=[0.5:0.01:2];
```

```
>>y=sin(x);
```

>>trapez(x,y)

Or

>>quad(@(x)sin(x),0.5,2)

>>quad(@(x)sin(x),0.5,2,1.e-10)

Derivation of trapezoidal rule:

For general interval $[x_{i-1}, x_i]$, we can make a change of variables to map it to $\left[-\frac{h}{2}, \frac{h}{2}\right]$. The procedure is as follows.

$$\tilde{x} = x - \underbrace{\left(a + ih - \frac{h}{2}\right)}_{\text{midpoint of the interval } [x_{i-1}, x_i] = [a+(i-1)h, a+ih]} = x - a - ih + \frac{h}{2}$$

$$\Rightarrow d\tilde{x} = dx, \quad x_{i-1} \rightarrow -\frac{h}{2}, \quad x_i \rightarrow \frac{h}{2}$$

$$\int_{x_{i-1}}^{x_i} f(x) dx = \int_{-h/2}^{h/2} f\left(\tilde{x} + a + ih - \frac{h}{2}\right) d\tilde{x}$$

So we only need to consider a special interval $\left[-\frac{h}{2}, \frac{h}{2}\right]$.

We use $I[f] \stackrel{\text{def}}{=} h\left(a_{-1}f\left(-\frac{h}{2}\right) + a_1f\left(\frac{h}{2}\right)\right)$ to approximate $\int_{-h/2}^{h/2} f(x)dx$. We have

$$\underbrace{I[f]}_{\text{Numerical approximation}} = \underbrace{\int_{-h/2}^{h/2} f(x) dx}_{\text{Exact value}} + \underbrace{e(h)}_{\text{Error}}$$

Question: How to determine a_{-1} and a_1 ?

Answer: We require that $I[f] = \int_{-h/2}^{h/2} f(x)dx$ for $f(x)=1$ and $f(x)=x$

$$I[1] = \int_{-h/2}^{h/2} 1 dx \implies h(a_{-1} + a_1) = h \implies a_{-1} + a_1 = 1$$

$$I[x] = \int_{-h/2}^{h/2} x \, dx \implies h \left(a_{-1} \left(-\frac{h}{2} \right) + a_1 \frac{h}{2} \right) = 0 \implies a_{-1} - a_1 = 0$$

$$\implies \begin{cases} a_{-1} + a_1 = 1 \\ a_{-1} - a_1 = 0 \end{cases} \implies \begin{cases} a_{-1} = \frac{1}{2} \\ a_1 = \frac{1}{2} \end{cases}$$

$$\implies I[f] = \frac{h}{2} \left(f \left(-\frac{h}{2} \right) + f \left(\frac{h}{2} \right) \right)$$

Question: What is the order of $e(h)$?

Answer: Let us do error analysis. We use the Taylor expansion of $f(x)$,

$$f(x) = f(0) + f'(0)x + \frac{f''(0)}{2}x^2 + \dots$$

Notice that both $I[f]$ and $\int_{-h/2}^{h/2} f(x)dx$ are linear in terms of $f(x)$.

$$\begin{aligned} e(h) &= I[f] - \int_{-h/2}^{h/2} f(x)dx \\ &= I \left[f(0) + f'(0)x + \frac{f''(0)}{2}x^2 + \dots \right] - \int_{-h/2}^{h/2} \left(f(0) + f'(0)x + \frac{f''(0)}{2}x^2 + \dots \right) dx \\ &= f(0) \left\{ I[1] - \int_{-h/2}^{h/2} 1 dx \right\} + f'(0) \left\{ I[x] - \int_{-h/2}^{h/2} x dx \right\} + \frac{f''(0)}{2} \left\{ I[x^2] - \int_{-h/2}^{h/2} x^2 dx \right\} + \dots \end{aligned}$$

The first two terms on the right hand side are zero. The third term is

$$I[x^2] - \int_{-h/2}^{h/2} x^2 dx = h \left(\frac{1}{2} \left(-\frac{h}{2} \right)^2 + \frac{1}{2} \left(\frac{h}{2} \right)^2 \right) - \frac{x^3}{3} \Big|_{-h/2}^{h/2} = \frac{h^3}{6}$$

$$\implies e(h) = \frac{f''(0)}{12} h^3 + \dots = O(h^3)$$

Thus, we derived

$$\underbrace{\frac{h}{2} \left(f \left(-\frac{h}{2} \right) + f \left(\frac{h}{2} \right) \right)}_{\text{Numerical approximation}} = \underbrace{\int_{-h/2}^{h/2} f(x)dx}_{\text{Exact value}} + \underbrace{e(h)}_{\text{Error}}$$

$$e(h) = O(h^3)$$

Apply this to $[x_{i-1}, x_i]$, we have

$$\underbrace{\frac{h}{2}(f_{i-1} + f_i)}_{\text{Numerical approximation}} = \underbrace{\int_{x_{i-1}}^{x_i} f(x) dx}_{\text{Exact value}} + \underbrace{e_i(h)}_{\text{Error}},$$

$$f_i = f(x_i)$$

$$e_i(h) = O(h^3)$$

This is the trapezoidal rule. A possible weakness in the trapezoidal rule can be inferred from the above analysis. If $f(x)$ does not have two continuous derivatives on $[a, b]$, then does $T_n(f)$ converge more slowly? The answer is yes for some functions, especially if the first derivative is not continuous. Same statement is true for Simpsons' rule, as we will see in the sample code (Codes\Estimate_order\example_4f.m)

Derivation of Simpson's rule:

We use $I[f] \stackrel{\text{def}}{=} h \left(a_{-1} f\left(-\frac{h}{2}\right) + a_0 f(0) + a_1 f\left(\frac{h}{2}\right) \right)$ to approximate $\int_{-h/2}^{h/2} f(x) dx$. We have

$$\underbrace{I[f]}_{\text{Numerical approximation}} = \underbrace{\int_{-h/2}^{h/2} f(x) dx}_{\text{Exact value}} + \underbrace{e(h)}_{\text{Error}}$$

Question: How to determine a_{-1}, a_0 and a_1 ?

Answer: we require that $I[f] = \int_{-h/2}^{h/2} f(x) dx$ for $f(x) = 1$, $f(x) = x$ and $f(x) = x^2$

$$I[1] = \int_{-h/2}^{h/2} 1 dx \implies h(a_{-1} + a_0 + a_1) = h \implies a_{-1} + a_0 + a_1 = 1$$

$$I[x] = \int_{-h/2}^{h/2} x dx \implies h \left(a_{-1} \left(-\frac{h}{2}\right) + a_1 \frac{h}{2} \right) = 0 \implies a_{-1} - a_1 = 0$$

$$I[x^2] = \int_{-h/2}^{h/2} x^2 dx \implies h \left(a_{-1} \left(-\frac{h}{2}\right)^2 + a_1 \left(\frac{h}{2}\right)^2 \right) = \frac{h^3}{12} \implies a_{-1} + a_1 = \frac{1}{3}$$

$$\implies \begin{cases} a_{-1} + a_0 + a_1 = 1 \\ a_{-1} - a_1 = 0 \\ a_{-1} + a_1 = \frac{1}{3} \end{cases} \implies \begin{cases} a_{-1} = \frac{1}{6} \\ a_0 = \frac{4}{6} \\ a_1 = \frac{1}{6} \end{cases}$$

$$\implies I[f] = \frac{h}{6} \left(f\left(-\frac{h}{2}\right) + 4f(0) + f\left(\frac{h}{2}\right) \right)$$

Let us do error analysis.

$$\begin{aligned} e(h) &= I[f] - \int_{-h/2}^{h/2} f(x) dx \\ &= I \left[f(0) + f'(0)x + \frac{f''(0)}{2}x^2 + \dots \right] - \int_{-h/2}^{h/2} \left(f(0) + f'(0)x + \frac{f''(0)}{2}x^2 + \dots \right) dx \\ &= f(0) \left\{ I[1] - \int_{-h/2}^{h/2} 1 dx \right\} + f'(0) \left\{ I[x] - \int_{-h/2}^{h/2} x dx \right\} + \frac{f''(0)}{2} \left\{ I[x^2] - \int_{-h/2}^{h/2} x^2 dx \right\} \\ &\quad + \frac{f^{(3)}(0)}{6} \left\{ I[x^3] - \int_{-h/2}^{h/2} x^3 dx \right\} + \frac{f^{(4)}(0)}{24} \left\{ I[x^4] - \int_{-h/2}^{h/2} x^4 dx \right\} + \dots \end{aligned}$$

The first three terms on the right hand side are zero. The fourth term is also zero because x^3 is an odd function. The fifth term is

$$I[x^4] - \int_{-h/2}^{h/2} x^4 dx = h \left(\frac{1}{6} \left(-\frac{h}{2}\right)^4 + \frac{1}{6} \left(\frac{h}{2}\right)^4 \right) - \frac{x^5}{5} \Big|_{-h/2}^{h/2} = \frac{h^5}{120}$$

$$\implies e(h) = \frac{f^{(4)}(0)}{2880} h^5 + \dots = O(h^5)$$

Thus, we derived

$$\implies \underbrace{\frac{h}{6} \left(f\left(-\frac{h}{2}\right) + 4f(0) + f\left(\frac{h}{2}\right) \right)}_{\text{Numerical approximation}} = \underbrace{\int_{-h/2}^{h/2} f(x) dx}_{\text{Exact value}} + \underbrace{e(h)}_{\text{Error}}$$

$$e(h) = O(h^5)$$

Applying this to $[x_{i-1}, x_i]$, we have

$$\underbrace{\frac{h}{6} (f_{i-1} + 4f_{i-1/2} + f_i)}_{\text{Numerical approximation}} = \underbrace{\int_{x_{i-1}}^{x_i} f(x) dx}_{\text{Exact value}} + \underbrace{e_i(h)}_{\text{Error}},$$

$$f_i = f(x_i), \quad f_{i-1/2} = f(x_{i-1/2}), \quad x_{i-1/2} = a + \left(i - \frac{1}{2}\right)h$$

$$e_i(h) = O(h^5)$$

This is the Simpson's rule.

The trapezoidal rule is the preferred integration rule when we are dealing with smooth periodic integrands.

Gaussian Numerical Integration

Gaussian numerical integration method is a numerical method that is based on the exact integration of polynomials of increasing degree; no subdivision of the integration interval is used. To illustrate the derivative of such integration formulas, we restrict our attend to the

integral $I(f) = \int_{-1}^1 f(x) dx$. Its relation to integrals over $[a, b]$ is as follows.

For an integral over $[a, b]$ $I(f) = \int_a^b f(x) dx$, introduce the linear change of

variable $x = \frac{b+a+t(b-a)}{2}$, $-1 \leq t \leq 1$, transforming the integral to

$$I(f) = \int_{-1}^1 \underbrace{f\left(\frac{b+a+t(b-a)}{2}\right)}_{\text{new integrand function}} \frac{b-a}{2} dt$$

So we consider $I(f) = \int_{-1}^1 f(x) dx$. The integration formula is to have the general form

$$I_n(f) = \sum_{j=1}^n w_j f(x_j) \text{ and we require that the nodes } [x_1, x_2, \dots, x_n] \text{ and weights } [w_1, w_2, \dots, w_n] \text{ be}$$

so chosen that $I_n(f) = I(f)$ for all polynomials $f(x)$ of as large a degree as possible.

Case 1 $n = 1$: The integration formula has the form

$$\int_{-1}^1 f(x) dx \approx w_1 f(x_1) \quad [\text{case1}]$$

We need to find w_1, x_1 .

It is to be exact for polynomials of as large a degree as possible.

Using $f(x) \equiv 1$ and forcing equality in [case1] give us

$$2 = w_1$$

Now use $f(x) = x$ and again force equality in (case1). Then

$$0 = w_1 x_1 \Rightarrow x_1 = 0 \text{ since } w_1 = 2 \neq 0. \text{ Thus, [case1] becomes}$$

$$\int_{-1}^1 f(x) dx \approx 2f(0) = I_1(f) \quad [\text{case1_2}]$$

This is the midpoint formula which is exact for all linear polynomials. To see that [case1] is not exact for quadratics, let $f(x) = x^2$. Then the error in [case1_2] is given by

$$\int_{-1}^1 x^2 dx - 2(0)^2 = \frac{2}{3} \neq 0.$$

Case 2: $n = 2$ The integration formula is

$$\int_{-1}^1 f(x) dx \approx w_1 f(x_1) + w_2 f(x_2)$$

Need to determine four unspecified quantities: w_1, w_2, x_1, x_2 . To determine these, we require it to be exact for the four monomials:

$$f(x) = 1, x, x^2, x^3$$

This leads to the four equations

$$2 = w_1 + w_2$$

$$0 = w_1 x_1 + w_2 x_2$$

$$\frac{2}{3} = w_1 x_1^2 + w_2 x_2^2$$

$$0 = w_1 x_1^3 + w_2 x_2^3$$

This is a nonlinear system in four unknowns; its solution can be shown to be

$w_1 = w_2 = 1, x_1 = -\frac{\sqrt{3}}{3}, x_2 = \frac{\sqrt{3}}{3}$ along with one based on reversing the signs of x_1 and x_2 . This yields the integration formula:

$$\int_{-1}^1 f(x) dx \approx f\left(-\frac{\sqrt{3}}{3}\right) + f\left(\frac{\sqrt{3}}{3}\right) = I_2(f)$$

From being exact for the monomials $1, x, x^2, x^3$, one can show this formula will be exact for all polynomials of degree ≤ 3 . It can also be shown by direct calculation to not be exact for the degree 4 polynomial $f(x) = x^4$.

Example: Approximate $I = \int_{-1}^1 e^x dx = e - e^{-1} \approx 2.3504024$

Using case 2 formula, we get $I_2 = e^{-\sqrt{3}/3} + e^{\sqrt{3}/3} \approx 2.3426961$

$I - I_2 \approx 0.00771$ The error is quite small for using such a small number of node points.

Case $n > 2$: We seek the formula $I_n(f) = \sum_{j=1}^n w_j f(x_j)$ (Gaussian numerical integration method), which has $2n$ unspecified parameters $w_1, \dots, w_n, x_1, \dots, x_n$, by forcing the integration formula to be exact for the $2n$ polynomials

$$f(x) = 1, x, x^2, \dots, x^{2n-1}$$

In turn, this forces $I_n(f) = I(f)$ for all polynomials f of degree $\leq 2n-1$. This leads to the following system of $2n$ nonlinear equations in $2n$ unknowns:

$$2 = w_1 + w_2 + \dots + w_n$$

$$0 = w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

$$\frac{2}{3} = w_1 x_1^2 + w_2 x_2^2 + \dots + w_n x_n^2$$

$$0 = w_1 x_1^3 + w_2 x_2^3 + \dots + w_n x_n^3$$

⋮

$$\frac{2}{2n-1} = w_1 x_1^{2n-2} + \dots + w_n x_n^{2n-2}$$

$$0 = w_1 x_1^{2n-1} + \dots + w_n x_n^{2n-1}$$

The resulting formula is of the order $(2n-1)$. Solving this system is a formidable problem. Thankfully, the nodes $\{x_i\}$ and weights $\{w_i\}$ have been calculated and collected in tables for the most commonly used values of n .

There is also another approach to the development of the numerical integration formula

$I_n(f) = \sum_{j=1}^n w_j f(x_j)$, using the theory of orthogonal polynomials. From that theory, it can be

shown that the nodes $[x_1, x_2, \dots, x_n]$ are the zeros of the Legendre polynomial $P_n(x)$ of degree n on the interval $[-1, 1]$ and its weights are given by

$$w_i = \frac{2}{(1-x_i^2) [P_n'(x_i)]^2}$$

Since the Legendre polynomials are well-known, the nodes $[x_1, x_2, \dots, x_n]$ can be found without any recourse to the nonlinear system. We list nodes and weights for several cases below.

Number of points, n	Points, x_i	Weights, w_i
1	0	2
2	$\pm \frac{\sqrt{3}}{3}$	1
3	0 $\pm \sqrt{\frac{3}{5}}$	$\frac{8}{9}$ $\frac{5}{9}$
4	$\pm \sqrt{(3-2\sqrt{6/5})/7}$ $\pm \sqrt{(3+2\sqrt{6/5})/7}$	$\frac{18+\sqrt{30}}{36}$ $\frac{18-\sqrt{30}}{36}$
5	0 $\pm \frac{1}{3}\sqrt{5-2\sqrt{10/7}}$ $\pm \frac{1}{3}\sqrt{5+2\sqrt{10/7}}$	128/225 $\frac{322+13\sqrt{70}}{900}$ $\frac{322-13\sqrt{70}}{900}$

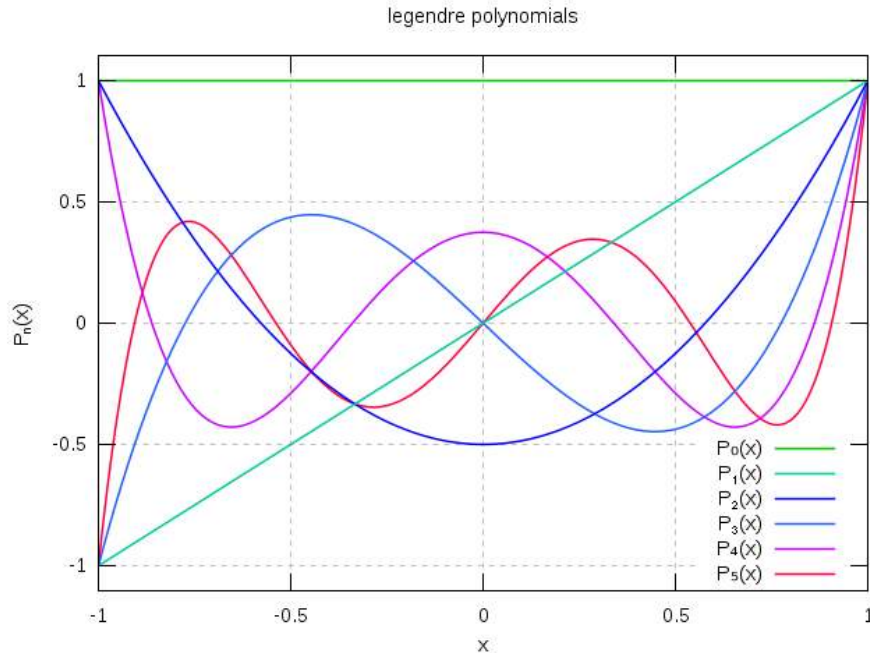
Legendre polynomials:

$P_0(x) = 1, P_1(x) = x$
 $(n+1)P_{n+1}(x) = (2n+1)xP_n(x) - nP_{n-1}(x)$ (Bonnet's recursion formula)

The first few Legendre polynomials are:

n	$P_n(x)$
0	1
1	x
2	$\frac{1}{2}(3x^2 - 1)$
3	$\frac{1}{2}(5x^3 - 3x)$
4	$\frac{1}{8}(35x^4 - 30x^2 + 3)$
5	$\frac{1}{8}(63x^5 - 70x^3 + 15x)$

The graphs of these polynomials (up to $n = 5$) are shown below:



Gauss–Lobatto rules

Also known as **Lobatto quadrature**, named after Dutch mathematician [Rehuel Lobatto](#).

It is similar to Gaussian quadrature with the following differences:

1. The integration points include the end points of the integration interval.
2. It is accurate for polynomials up to degree $2n-3$, where n is the number of integration points.

Lobatto quadrature of function $f(x)$ on interval $[-1, +1]$:

$$\int_{-1}^1 f(x) dx = \frac{2}{n(n-1)} [f(1) + f(-1)] + \sum_{i=2}^{n-1} w_i f(x_i) + R_n.$$

Abscissas: x_i is the $(i-1)$ st zero of $P'_{n-1}(x)$.

Weights:

$$w_i = \frac{2}{n(n-1)[P_{n-1}(x_i)]^2} \quad (x_i \neq \pm 1).$$

Remainder:

$$R_n = \frac{-n(n-1)^3 2^{2n-1} [(n-2)!]^4}{(2n-1)[(2n-2)!]^3} f^{(2n-2)}(\xi), \quad (-1 < \xi < 1)$$