

# CS6210 - Homework/Assignment-6

Arnab Das(u1014840)

December 8, 2016

---

**Question-1: Chapter-13, question-5**

---

(a) Given function,  $f(t) = e^{-3t}$  on the interval  $[0,3]$ . Let  $q_2(t)$  be the approximation of degree 2 and  $q_3(t)$  be the approximation of degree 3.

Let  $\phi_0(t), \phi_1(t), \phi_2(t), \phi_3(t)$  be the polynomial basis functions. We use the affine transformation to map them to the legendre polynomials as:

$$x = \frac{2t - a - b}{b - a} = \frac{2t - 3}{3}$$

Then , the basis functions in form of legendre polynomials become:

$$\phi_0(t) = 1$$

$$\phi_1(t) = \frac{2t - 3}{3}$$

$$\phi_2(t) = \frac{1}{2} \left[ 3 \left( \frac{2t - 3}{3} \right)^2 - 1 \right]$$

$$\phi_3(t) = \frac{1}{2} \left[ 5 \left( \frac{2t - 3}{3} \right)^3 - 3 \frac{2t - 3}{3} \right]$$

**Note:** The evaluations are done in matlab.

Next, we need to evaluate the matrix elements  $B(j, k) = \int_a^b \phi_j(x) \phi_k(x)$ . Since, the basis are orthogonal as we have chosen Legendre polynomials, hence we have only the diagonal elements: Thus:

$$B_{0,0} = \int_0^3 \phi_0(t) \phi_0(t) dt = 3$$

$$B_{1,1} = \int_0^3 \phi_1(t) \phi_1(t) dt = 1$$

$$B_{2,2} = \int_0^3 \phi_2(t) \phi_2(t) dt = 0.6$$

$$B_{3,3} = \int_0^3 \phi_3(t) \phi_3(t) dt = 0.4280$$

Next, evaluate the  $b$  vector of the rhs whose elements are defined as:  $b_j = \int_a^b \phi_j(t) f(t) dt$

$$b_0 = \int_0^3 \phi_0(t) f(t) dt = 0.3333$$

$$b_1 = \int_0^3 \phi_1(t) f(t) dt = -0.2593$$

$$b_2 = \int_0^3 \phi_2(t) f(t) dt = 0.1604$$

$$b_3 = \int_0^3 \phi_3(t) f(t) dt = -0.0811$$

Then, we can get the coefficients as:  $c_j = \frac{b_j}{B_{j,j}}$

$$c_0 = \frac{b_0}{B_{0,0}} = 0.111$$

$$c_1 = \frac{b_1}{B_{1,1}} = -0.2593$$

$$c_2 = \frac{b_2}{B_{2,2}} = 0.2674$$

$$c_3 = \frac{b_3}{B_{3,3}} = -0.1892$$

Hence, the quadratic will be:

$$q_2(t) = c_0\phi_0(t) + c_1\phi_1(t) + c_2\phi_2(t) = 0.111 - 0.2593\left(\frac{2t-3}{3}\right) + 0.2674\frac{1}{2}\left[3\left(\frac{2t-3}{3}\right)^2 - 1\right]$$

And the cubic will be:

$$q_3(t) = c_0\phi_0(t) + c_1\phi_1(t) + c_2\phi_2(t) + c_3\phi_3(t) \\ 0.111 - 0.2593\left(\frac{2t-3}{3}\right) + 0.2674\frac{1}{2}\left[3\left(\frac{2t-3}{3}\right)^2 - 1\right] - 0.1892\frac{1}{2}\left[5\left(\frac{2t-3}{3}\right)^3 - 3\frac{2t-3}{3}\right]$$

The code implementation is in Prob1.m

(b) The below figure shows the plots for the quadratic, cubic and original plot in the first subplot, and the error plots for quadratic and cubic with respect to the original, in the second subplot.

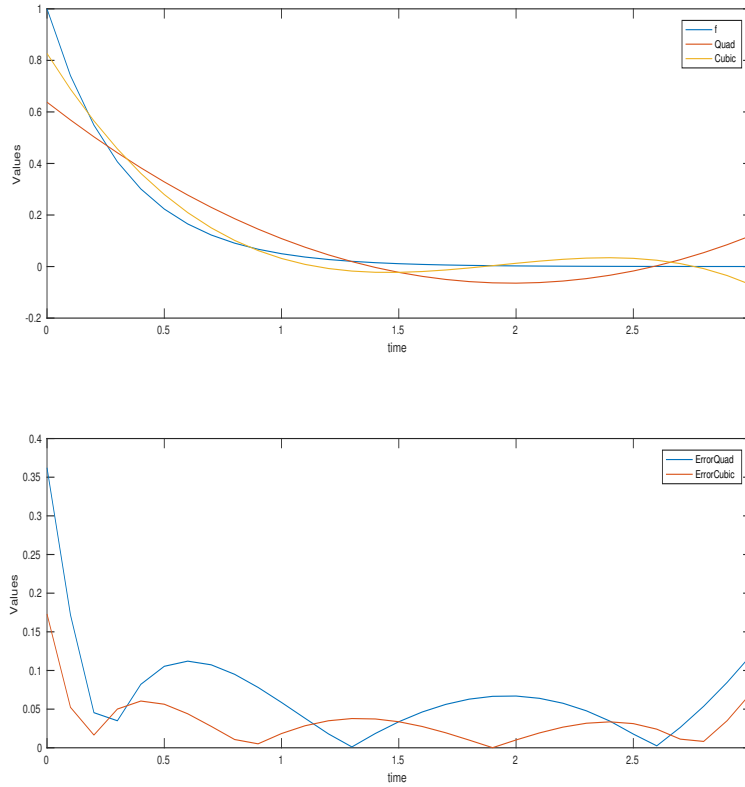


Figure 1: Quadratic and Cubic Polynomial Approximation

As the plot shows, the cubic polynomial (yellow line) very closely approximates the original than the quadratic one. The error plot also suggests, that the error magnitude mostly is much higher for the quadratic polynomial than the cubic polynomial. The evaluated l2norm of the error functions, gives **l2-norm of quad = 0.53115307** and **l2-norm of cubic = 0.25287738**, indicating that the cubic is a better approximating polynomial in this case.

(c) For orthogonal basis functions, we have  $\int_a^b \phi_j \phi_k = 0, \text{ for } j \neq k$ . Thus, each additional orthogonal basis added to the approximating polynomial aims to reduce the residual from the previous set of basis functions. Since, it always reduces the residue, it cannot make the approximating polynomial any worse than the lower degree polynomial.

---

**Question-2: Chapter-14, question-2**

---

Let  $f(x)$  be a given function that can be evaluated at points  $x_0 \pm jh, j = 0, 1, 2, \dots$  for any fixed value of  $h, 0 < h \ll 1$ .

(a) Using Taylor series expansion about  $x_0$ , we can write the following expansions about  $\pm h$ :

$$\begin{aligned} f(x_0 + h) &= f(x_0) + hf'(x_0) + \frac{h^2}{2}f''(x_0) + \frac{h^3}{6}f'''(x_0) + \frac{h^4}{24}f^{(4)}(x_0) + \frac{h^5}{120}f^{(5)}(\zeta_1) \\ f(x_0 - h) &= f(x_0) - hf'(x_0) + \frac{h^2}{2}f''(x_0) - \frac{h^3}{6}f'''(x_0) + \frac{h^4}{24}f^{(4)}(x_0) - \frac{h^5}{120}f^{(5)}(\zeta_2) \end{aligned}$$

where,  $x_0 \leq \zeta_1 \leq x_0 + h$ , and  $x_0 - h \leq \zeta_2 \leq x_0$  The difference of the above two equations gives the following:

$$f(x_0 + h) - f(x_0 - h) = 2\{hf'(x_0) + \frac{h^3}{6}f'''(x_0) + \frac{h^5}{120}f^{(5)}(\zeta_3)\} \quad (1)$$

where,  $x_0 - h \leq \zeta_3 \leq x_0 + h$

Similarly, we can get the following expansions using nearby points at  $x_0 \pm 2h$ :

$$\begin{aligned} f(x_0 + 2h) &= f(x_0) + 2hf'(x_0) + \frac{4h^2}{2}f''(x_0) + \frac{8h^3}{6}f'''(x_0) + \frac{16h^4}{24}f^{(4)}(x_0) + \frac{32h^5}{120}f^{(5)}(\zeta_4) \\ f(x_0 - 2h) &= f(x_0) - 2hf'(x_0) + \frac{4h^2}{2}f''(x_0) - \frac{8h^3}{6}f'''(x_0) + \frac{16h^4}{24}f^{(4)}(x_0) - \frac{32h^5}{120}f^{(5)}(\zeta_5) \end{aligned}$$

where,  $x_0 \leq \zeta_4 \leq x_0 + 2h$ , and  $x_0 - 2h \leq \zeta_5 \leq x_0$  The difference of the above two equations gives the following:

$$f(x_0 + 2h) - f(x_0 - 2h) = 2\{2hf'(x_0) + \frac{8h^3}{6}f'''(x_0) + \frac{32h^5}{120}f^{(5)}(\zeta_6)\} \quad (2)$$

where,  $x_0 - 2h \leq \zeta_6 \leq x_0 + 2h$

Then multiplying (1) by 2 and subtracting it from (2), we get:

$$f(x_0 + 2h) - f(x_0 - 2h) - 2f(x_0 + h) + 2f(x_0 - h) = 2h^3f'''(x_0) + \frac{h^5}{2}f^{(5)}(\zeta)$$

where,  $x_0 - 2h \leq \zeta \leq x_0 + 2h$  Rearranging the terms, we get:

$$f'''(x_0) = \frac{f(x_0 + 2h) - f(x_0 - 2h) - 2f(x_0 + h) + 2f(x_0 - h)}{2h^3} + \frac{-h^2}{4}f^{(5)}(\zeta) \quad (3)$$

(3) provides the formula for approximating the third derivative  $f'''(x_0)$ , with the bold test showing the component of the truncation error: Thus, truncation error =  $\frac{-h^2}{4}f^{(5)}(\zeta)$

(d) The below figure shows the plot of the approximated 3rd derivative derived as above.

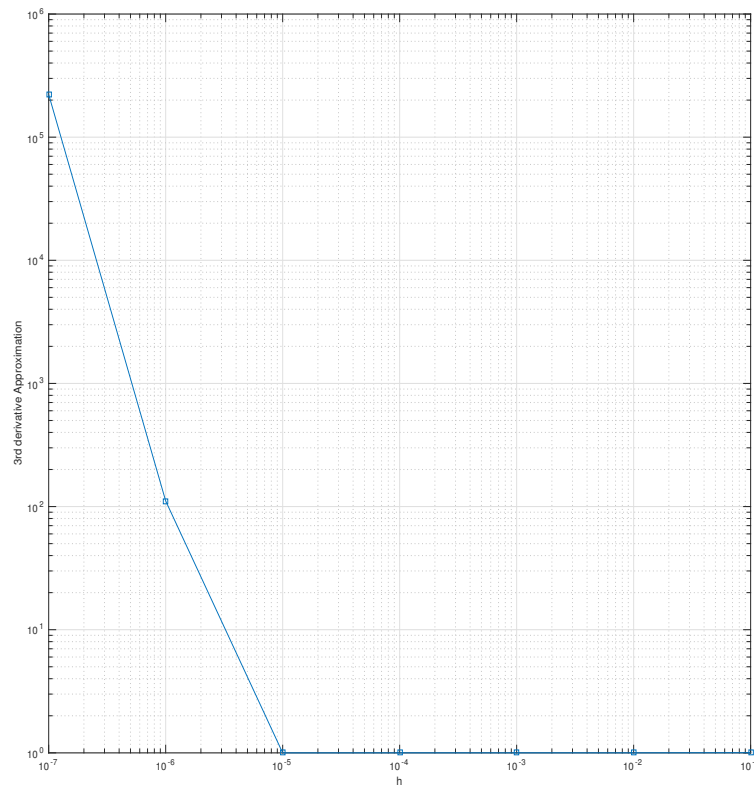


Figure 2: Approximated 3rd derivative

As the figure suggests, beyond  $h = 10^{-5}$ , the round-off error dominates. Since the true value of the function is 1 at  $x = 0$ , hence the deviation of the curve from the x-axis depicts the magnitude of the error due to increasing round-off with reducing values of  $h$ . However, for larger values of  $h$ , it is quite accurate. We plot an expanded plot in the range where it shows high accuracy as below:

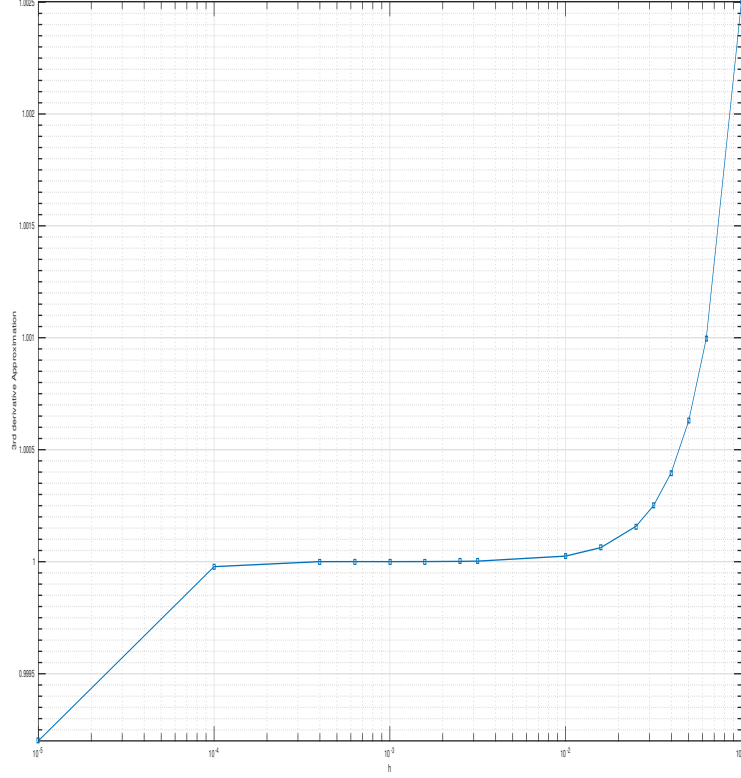


Figure 3: Improved accuracy at higher h

As the parabolic shape of the curve suggests that the approximation is indeed second order accurate for large values of h. Best value of approximation we get for  $h \approx 10^{-3}$ .

(c) To derive the expression for the round-off error, we denote the floating point evaluations of  $f$  by  $F$ . Thus the expression for round off error will be  $|F'''(x_0) - f'''(x_0)|$

$$= \left| \frac{F(x_0 + 2h) - F(x_0 - 2h) - 2F(x_0 + h) + 2F(x_0 - h)}{2h^3} - \frac{f(x_0 + 2h) - f(x_0 - 2h) - 2f(x_0 + h) + 2f(x_0 - h)}{2h^3} \right|$$

$$= \left| \frac{e_r(x_0 + 2h) - e_r(x_0 - 2h) - 2e_r(x_0 + h) + 2e_r(x_0 - h)}{2h^3} \right|$$

where, the round-off error term,  $e_r(x) = F(x) - f(x)$

$$\leq \left| \frac{e_r(x_0 + 2h)}{2h^3} \right| + \left| \frac{e_r(x_0 - 2h)}{2h^3} \right| + \left| \frac{2e_r(x_0 + h)}{2h^3} \right| + \left| \frac{2e_r(x_0 - h)}{2h^3} \right|$$

Assuming the each round-off error term is bounded by the machine precision,  $\epsilon$ , then we get:

$$|F'''(x_0) - f'''(x_0)| \leq \frac{6\epsilon}{2h^3} = \frac{3\epsilon}{h^3}$$

Thus, the round-off error grows inversely proportional to  $h$ . Hence, as h decreases, the round-off error increases, thus with  $h \rightarrow 0$ , the round-off error will tend to  $\infty$ .

(d) We can obtain a fourth order formula in a similar manner as we obtained the third order formula with some minor changes. First evaluate  $S1 : f(x_0 + h) + f(x_0 - h)$ , then evaluate  $S2 : f(x_0 + 2h) + f(x_0 - 2h)$ . Then subtract  $S2 - 4S1$ , to extract the fourth order formula.

---

**Question-3: Chapter-14: question-8**

---

Derivation of an approximate formulae for the second derivatives  $f''(x_0)$  of a smooth function  $f(x)$  using the three points  $x_{-1}, x_0 = x_{-1} + h_0, x_1 = x_0 + h_1$ , where  $h_0 \neq h_1$ .

(a) To prove both the given approximations are the same, we need to show:  $2[f_{x-1, x_0, x_1}] = \frac{g_{1/2} - g_{-1/2}}{(h_0 + h_1)/2}$ .

**Proof:**

Using divided difference, we can write:

$$2[f_{x-1, x_0, x_1}] = 2 \times \frac{\frac{f(x_1) - f(x_0)}{x_1 - x_0} - \frac{f(x_0) - f(x_{-1})}{x_0 - x_{-1}}}{x_1 - x_{-1}}$$

$$2[f_{x-1, x_0, x_1}] = 2 \times \frac{\frac{f(x_1) - f(x_0)}{h_1} - \frac{f(x_0) - f(x_{-1})}{h_0}}{h_1 + h_0}$$

Given, that:  $g_{1/2} = \frac{f(x_1) - f(x_0)}{h_1}$  and  $g_{-1/2} = \frac{f(x_0) - f(x_{-1})}{h_0}$ , we can replace them to get:

$$2[f_{x-1, x_0, x_1}] = 2 \times \frac{g_{1/2} - g_{-1/2}}{h_1 + h_0}$$

$$2[f_{x-1, x_0, x_1}] = \frac{g_{1/2} - g_{-1/2}}{\frac{h_1 + h_0}{2}}$$

(b) To show that the method is only first order accurate, we use the Taylor series expansion upto second order with third order error term to derive the formulation:

$$f(x_0 + h_1) = f(x_0) + h_1 f'(x_0) + \frac{h_1^2}{2} f''(x_0) + \frac{h_1^3}{3!} f'''(\zeta_1)$$

$$f(x_0 - h_0) = f(x_0) - h_0 f'(x_0) + \frac{h_0^2}{2} f''(x_0) - \frac{h_0^3}{3!} f'''(\zeta_2)$$

Which can be further simplified to:

$$\frac{f(x_1) - f(x_0)}{h_1} = f'(x_0) + \frac{h_1}{2} f''(x_0) + \frac{h_1^2}{3!} f'''(\zeta_1)$$

$$\frac{f(x_{-1}) - f(x_0)}{h_0} = -f'(x_0) + \frac{h_0}{2} f''(x_0) - \frac{h_0^2}{3!} f'''(\zeta_2)$$

Adding these two, we get:

$$\frac{f(x_1) - f(x_0)}{h_1} - \frac{f(x_0) - f(x_{-1})}{h_0} = \frac{(h_0 + h_1)}{2} f''(x_0) + \frac{h_1^2}{3!} f'''(\zeta_1) - \frac{h_0^2}{3!} f'''(\zeta_2)$$

For the error term, replacing with a  $\zeta$ , such that  $x_{-1} \leq \zeta \leq x_1$ :

$$\frac{f(x_1) - f(x_0)}{h_1} - \frac{f(x_0) - f(x_{-1})}{h_0} = \frac{(h_0 + h_1)}{2} f''(x_0) + f'''(\zeta) \left( \frac{h_1^2}{3!} - \frac{h_0^2}{3!} \right)$$

Rearranging to get the second order term:

$$f''(x_0) = \frac{\frac{f(x_1) - f(x_0)}{h_1} - \frac{f(x_0) - f(x_{-1})}{h_0}}{(h_0 + h_1)/2} - \frac{\mathbf{f'''(\zeta)(\mathbf{h_1} - \mathbf{h_0})}}{\mathbf{3}}$$

The error term shown in bold is thus only first order,  $O(h_0 + h_1)$ , accurate.

(c) The matlab code file is *Prob3.m*, where we run the two methods above for second derivative approximation of first order accuracy for the given values of h, around the point  $x_0 = 0$ . Both methods do not show any deviation from each other, that is, their mutual difference is zero. Thus, the table will be :

Table 1: Second derivative approximation result for  $e^x$  at  $x = 0$

h	firstMethod	secondMethod	difference
1.0e-1	14.0227	14.0227	0
1.0e-2	134.00	134.00	0
1.0e-3	-1334.00	-1334.00	0
1.0e-4	13334.00	13334.00	0
1.0e-5	133334.00	133334.00	0

The results are extremely inaccurate for an approximation of second order derivative of  $e^x$  at  $x = 0$ . This happens mainly because of round-off as we keep reducing h. Since, there are subtracting terms of very close values, round-off escalates the error to large values.

---

#### Question-4: Chapter-14, question-13

---

(a)

$$f_{pp0} = \frac{(g_{1/2} - g_{-1/2})}{h} = \frac{\frac{f_1 - f_0}{h} - \frac{f_0 - f_{-1}}{h}}{h} = \frac{f_1 - 2f_0 + f_{-1}}{h^2}$$

Hence, it is same as the  $f_{pp0}$  defined in Example 12.

(b) Here we denote the two methods as:

**First Method:**  $(f_1 - 2f_0 + f_{-1})/h^2$

**Second Method:**  $\frac{g_{1/2} - g_{-1/2}}{h}$ , where  $g_{1/2} = \frac{f_1 - f_0}{h}$  and  $g_{-1/2} = \frac{f_0 - f_{-1}}{h}$ .

The below table reports the error/deviations of these methods from the true second order derivative for  $f(x) = \sin(x)$  at  $x_0 = 1.2$

Table 2: Second derivative approximation result for  $\sin x$  at  $x = 1.2$

h	firstMethod	secondMethod	difference
$10^0$	0.075126648647425	0.075126648647425	0
$10^{-0.5}$	0.007741148589218	0.007741148589218	0
$10^{-1.0}$	0.000776440384795	0.000776440384795	0.000000222044605
$10^{-1.5}$	0.000077667334949	0.000077667334949	0.000000777156117
$10^{-2.0}$	0.000007766966345	0.000007766966345	0
$10^{-2.5}$	0.000000776702422	0.000000776702429	0.000006661338148
$10^{-3.0}$	0.000000077772619	0.000000077772619	0
$10^{-3.5}$	0.000000008938792	0.000000008938954	0.000162203583898
$10^{-4.0}$	0.000000006718345	0.000000006718345	0



$10^{-4.5}$	0.000000084433957	0.000000084433986	0.000028532731733
$10^{-5.0}$	0.000001305679284	0.000001305679284	0.000000111022302
$10^{-5.5}$	0.000006856794407	0.000006856780840	0.013567702517037
$10^{-6.0}$	0.000117879096870	0.000117879096870	0
$10^{-6.5}$	0.001672191331345	0.001672191233429	0.097916452723723
$10^{-7.0}$	0.010553975528346	0.010553975528346	0.000000111022302
$10^{-7.5}$	0.154882968729617	0.154882969473164	0.743547445836157
$10^{-8.0}$	0.932039085967226	0.932039085967226	0

The below plots figure shows the error plots using the first and second method and the difference between the two methods(mainly due to round-off issues).

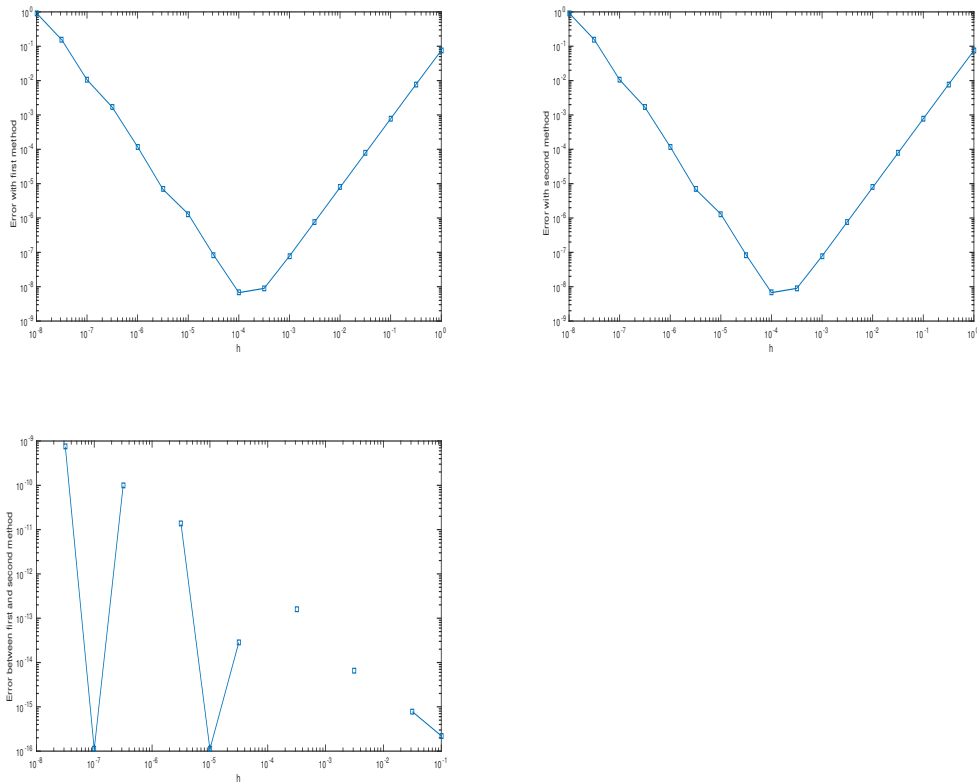


Figure 4: Error Plot for first and second method

Both the plots have the similar feature of reducing truncation error with reducing  $h$  until it reaches a critical point where round off error starts to dominate and the error grow larger with reducing  $h$ . The difference plot shows the variation in the two methods. The second method is more prone to round-off errors due to more subtractions between very close numbers.

---

**Question-5: Chapter-14, question-15**

---

Consider the numerical differentiation of the function,  $f(x) = c(x)e^{x/\pi}$ , defined on  $[0, \pi]$ , where

$$c(x) = j; \frac{\pi}{4}(j-1) \leq x < \frac{\pi}{4}j$$

for  $j = 1, 2, 3, 4, \dots$  (a) We want a difference approximation with step size,  $h = \frac{n}{\pi}$ . Suppose, we have  $n$  as

multiple of 4, that is,  $n = 4l$ , where  $l$  is an integer. Since, the points are equispaced, thus the entire  $[0, \pi]$  region can be seen as partitioned into 4 regions each having  $l$  points as described below:

**L1** :  $0 \leq x < \pi/4$  :  $l$  points at  $0, \frac{\pi}{4l}, \frac{2\pi}{4l}, \dots, \frac{(l-1)\pi}{4l}$

**L2** :  $\pi/4 \leq x < \pi/2$  :  $l$  points at  $\frac{\pi}{4}, \frac{\pi}{4} + \frac{\pi}{4l}, \frac{\pi}{4} + \frac{2\pi}{4l}, \dots, \frac{\pi}{4} + \frac{(l-1)\pi}{4l}$

**L3** :  $\pi/2 \leq x < 3\pi/4$  :  $l$  points at  $\frac{\pi}{2}, \frac{\pi}{2} + \frac{\pi}{4l}, \frac{\pi}{2} + \frac{2\pi}{4l}, \dots, \frac{\pi}{2} + \frac{(l-1)\pi}{4l}$

**L4** :  $3\pi/4 \leq x < \pi$  :  $l$  points at  $\frac{3\pi}{4}, \frac{3\pi}{4} + \frac{\pi}{4l}, \frac{3\pi}{4} + \frac{2\pi}{4l}, \dots, \frac{3\pi}{4} + \frac{(l-1)\pi}{4l}$

The last point is at  $\pi$

Now,  $c(x)$  is discontinuous at the points of integer multiples of  $\pi/4$ . To make sure that the discontinuities do not affect our difference formulation it is beneficial to have some of the interval end points at these discontinuities. For example, if we want to approximate the difference value at a point  $t$  in the interval  $[x_{i-1}, x_i]$ , there will be two possibilities while taking the difference. If  $x_{i-1}$  and  $x_i$  lie in the same region described above, then  $c(x)$  has the same values at both the points and can be approximated with the value of  $c(x_{i-1})$ . If on the other hand,  $x_{i-1}$  and  $x_i$  lie in two different regions, the only possibility is the  $x_{i-1}$  is the last point in one region while  $x_i$  is the first point of the other region and hence we have a discontinuity at  $x_i$ , however,  $t$  will be in the region of  $x_{i-1}$  and hence to avoid the discontinuity in the calculation, we can take the left continuous value at  $x_i$ , which will be equal to the value at  $t_i$ . Hence, with  $n$  being multiple of 4, we have the possibility of avoiding the discontinuities resulting in better values. Consider on the other hand, if  $n$  was not a multiple of 4, that is, we did not have this nice region partition, then the point  $t$  could itself have a discontinuity which we could not be able to derive from the sampling points  $[x_{i-1}, x_i]$ . Moreover, suppose we have  $t$  in  $[x_{i-1}, x_i]$ , where  $c(x_{i-1}) \neq c(x)$ , and  $t$  lies not at the point of discontinuity, but either in the step with  $x_{i-1}$  or with  $x_i$ . All these points lying on either side of the step, will get the same difference value which will be highly inaccurate. These issues we can solve by having  $n$  has multiple of 4, and having the nice partitioning described above.

**(b) Show that:**  $h^{-1}c(t_i)(e^{x_{i+1}/\pi} - e^{x_i/\pi})$  provides a second order approximation of  $f'(t_i)$

**Proof:**  $t_i = x_i + \frac{h}{2} = ih + \frac{h}{2} = (i + \frac{1}{2})h$ ;  $i = 0, 1, \dots, (n-1)$

Using Taylor series expansion around the point  $t_i = x_i + \frac{h}{2}$ , we can write:

$$f(x_i + \frac{h}{2} - \frac{h}{2}) = f(x_i + \frac{h}{2}) - \frac{h}{2}f'(x_i + \frac{h}{2}) + \frac{h^2}{4}f''(x_i + \frac{h}{2}) - \frac{h^3}{8}f'''(\zeta_1)$$

and,

$$f(x_i + \frac{h}{2} + \frac{h}{2}) = f(x_i + \frac{h}{2}) + \frac{h}{2}f'(x_i + \frac{h}{2}) + \frac{h^2}{4}f''(x_i + \frac{h}{2}) + \frac{h^3}{8}f'''(\zeta_2)$$

Subtracting  $f(x_{i+1} = x_i + h) - f(x_i)$  gives,

$$f(x_{i+1}) - f(x_i) = hf'(x_i + \frac{h}{2}) + \frac{h^3}{8}\{f'''(\zeta_1) + f'''(\zeta_2)\}$$

where  $t_i - h/2 \leq \zeta_1 \leq t_i$  and  $t_i \leq \zeta_2 \leq t_i + h/2$ . Choosing a  $\zeta$ , such that,  $t_i - h/2 \leq \zeta \leq t_i + h/2$ , we can then write:

$$f(x_{i+1}) - f(x_i) = hf'(x_i + \frac{h}{2}) + \frac{h^3}{4}\{f'''(\zeta)\}$$

Since,  $t_i = x_i + h/2$ ,

$$\begin{aligned} f(x_{i+1}) - f(x_i) &= hf'(t_i) + \frac{h^3}{4}\{f'''(\zeta)\} \\ f'(t_i) &= \frac{f(x_{i+1}) - f(x_i)}{h} - \frac{h^2}{4}f'''(\zeta) \\ f'(t_i) &= \frac{c(x_{i+1})e^{x_{i+1}/\pi} - c(x_i)e^{x_i/\pi}}{h} - \frac{h^2}{4}f'''(\zeta) \end{aligned}$$

As, described in first part of this question, if the end-points of the interval  $[x_i, x_{i+1}]$ ,  $t_i$  belongs to falls completely in one of the regions, then we have  $c(x_i) = c(t_i) = c(x_{i+1})$ . If  $x_i$  and  $x_{i+1}$  fall in different region, then  $x_{i+1}$  falls in the point of discontinuity, so we take the left continuous value at  $x_{i+1}$  which will be equal to  $c(t_i) = c(x_i)$ . Thus, in the above expression we can replace  $c(x_i), c(x_{i+1})$  with  $c(t_i)$ . Thus we get the following expression:

$$\begin{aligned} f'(t_i) &= \frac{c(t_i)e^{x_{i+1}/\pi} - c(t_i)e^{x_i/\pi}}{h} - \frac{h^2}{4}f'''(\zeta) \\ f'(t_i) &= \frac{c(t_i)(e^{x_{i+1}/\pi} - e^{x_i/\pi})}{h} - \frac{h^2}{4}f'''(\zeta) \end{aligned} \quad (4)$$

(4) shows the second order approximation of  $f'(t_i)$

#### Question-6: Chapter-15, question-4

(a) **Prove:** Error in basic corrected trapezoidal rule in the interval  $[a, b]$  can be estimated by:

$$E(f) = \frac{f'''' * (\eta)}{720} (b - a)^5$$

**Proof:** The osculating polynomial formula for the basic corrected trapezoidal rule is written as:

$$p_3(x) = f(a) + f'(a)(x - a) + f[a, a, b](x - a)^2 + f[a, a, b, b](x - a)^2(x - b)$$

The error in the polynomial interpolation in that case will be given by:

$$f[a, a, b, b, x](x - a)(x - a)(x - b)(x - b)$$

Then, to find the error in the integral of the polynomial, we can integrate the error of polynomial described above, in the interval  $[a, b]$

$$E(f) = \int_a^b f[a, a, b, b, x]\psi(x)$$

where  $\psi(x) = \prod_{i=0}^3 (x - x_i) = (x - a)(x - a)(x - b)(x - b)$

Notice that, since  $x$  lies in the interval  $[a, b]$ , hence  $(x - a \geq 0)$  and  $(x - b \leq 0)$ . In any case, the square of the terms will be greater than equal to 0. So, in the given interval  $\psi(x) \geq 0$  always. Because  $\psi(x)$  does not changes sign in the interval, then using the intermediate value theorem, there exists  $a \leq \eta \leq b$ , such that:

$$E(f) = \int_a^b f[a, a, b, b, x]\psi(x) = \int_a^b f[a, a, b, b, \eta]\psi(x)$$

where,

$$f[a, a, b, b, \eta] = \frac{f''''(\eta)}{4!}$$

which is a constant, say  $K$ .

Then we can write the error integral as:

$$E(f) = K \int_a^b (x - a)^2(x - b)^2$$

Doing integration by parts:

$$E(f) = K \left[ \frac{(x-a)^2(x-b)^3}{3} - \frac{(x-a)(x-b)^4}{6} + \frac{(x-b)^5}{30} \right]_a^b = \frac{-(a-b)^5}{30}$$

Replacing back K, we get:

$$E(f) = \frac{f''''(\eta)}{4!} \frac{-(a-b)^5}{30} = \frac{f''''(\eta)(b-a)^5}{720}$$

(b) The integral for the corrected trapezoidal is written is:

$$I_f \approx \int_a^b p_3(x)dx = \frac{(b-a)}{2} [f(a) + f(b)] + \frac{(b-a)^2}{12} [f'(a) - f'(b)]$$

**part-1:** For the integral  $\int_0^1 e^x dx$ , thus  $a = 0, b = 1$ , and  $f(x) = e^x, f'(x) = e^x$ . So,  $f(a) = 1, f(b) = e, f'(a) = 1, f'(b) = e$

Using the basic corrected trapezoidal, we get:

$$\int_0^1 e^x dx = 1.71595$$

the actual evaluation is 1.7183... while the basic trapezoidal evaluation from Example-15.2 is 1.7183.... Hence, the evaluation using the basic corrected trapezoidal is more accurate than the basic trapezoidal.

**part-2:** For the integral  $\int_{0.9}^1 e^x dx$ , thus  $a = 0.9, b = 1$ , and  $f(x) = e^x, f'(x) = e^x$ . So,  $f(a) = e^{0.9}, f(b) = e, f'(a) = e^{0.9}, f'(b) = e$

Using the basic corrected trapezoidal, we get:

$$\int_{0.9}^1 e^x dx = 0.258678$$

The actual evaluation is 0.2586787171..., while the basic trapezoidal evaluation from Example-15.2 is 0.258894.... hence, the evaluation using the basic corrected trapezoidal is more accurate than the basic trapezoidal.

### Question-7: Chapter-15, question-5

(a) In the interval  $[a, b]$ , the basic midpoint rule is given as :

$$I_f \approx (b-a)f\left(\frac{a+b}{2}\right) \quad (5)$$

For the composite midpoint rule, we consider  $r$  subintervals in the original interval  $[a, b]$  and apply the basic midpoint rule to each subinterval and then sum over all the subintervals to get the composite integral. The rule applied to an interval  $[t_{i-1}, t_i]$ , such that the interval widths are uniform and  $t_i - t_{i-1} = h = \frac{b-a}{r}$ , will be:

$$\int_{t_{i-1}}^{t_i} f(x)dx \approx hf\left(\frac{t_{i-1} + t_i}{2}\right)$$

Summing over all the subintervals to get the complete composite integral:

$$\int_a^b f(x)dx = h \sum_{i=1}^r f\left(\frac{t_{i-1} + t_i}{2}\right)$$

For,  $r$  equispaced intervals over  $[a, b]$ , we have the interval width as  $h = \frac{b-a}{r}$ . Then,  $t_0 = a, t_1 = a + h, t_2 = a + 2h, \dots, t_i = a + ih$ . So,

$$\frac{t_{i-1} + t_i}{2} = \frac{a + (i-1)h + a + ih}{2} = a + (i - \frac{1}{2})h$$

Replacing it in the original integral, we get the final form for the composite midpoint as:

$$\int_a^b f(x)dx \approx h \sum_{i=1}^r f(a + (i - \frac{1}{2})h)$$

From the above expression, we can see that there is one function evaluation per subinterval. Hence, the number of function evaluations is  $r = \frac{b-a}{h}$

### (b) Derive an expression for the error in composite midpoint rule

For the basic midpoint rule, the error expression in the interval  $[a, b]$  is given by:

$$E(f) = \frac{f''(\eta)}{24} (b-a)^3$$

This comes from doing the following integral:

$$\frac{f''}{2!}(\eta) \int_a^b (x - \frac{a+b}{2})(x - \frac{a+b}{2}) dx$$

The reason for adding the second term of  $(x - \frac{a+b}{2})$  even though there is only a single point, is that  $(x - \frac{a+b}{2})$  can change signs within the interval  $[a, b]$ , and hence we cannot apply the intermediate value theorem to take out  $f''(\eta)$  as a constant, for  $a \leq \eta \leq b$ . So, we duplicate the point  $\frac{a+b}{2}$  as a dummy interpolation point, since it does not change the area evaluated and hence the error term should remain the same.

Next, we come to the derivation of the error for the composite midpoint rule. In the composite cases, we have divided the original interval,  $[a, b]$ , into  $r$  equispaced sub-intervals of width  $h$ , such that  $r = \frac{b-a}{h}$ . For evaluating the composite integral using mid-point we applied the basic midpoint to each of these sub-intervals and summed them. Similarly, now each evaluation of the basic midpoint in these intervals will result in an error, which we can further sum up to get the expression for the error in composite mid-point. The error term in the interval  $[t_{i-1}, t_i]$  will be  $\frac{f''(\eta_i)}{24} h^3$ , where  $\frac{f''(\eta_i)}{24}$  is a constant in the interval  $[t_{i-1}, t_i]$

Hence, the expression of error for the composite midpoint will be:

$$E_{CM}(f) = \sum_{i=1}^r \frac{f''(\eta_i)}{24} h^3$$

Since,  $f''(\eta_i)$  is a constant in the interval  $[t_{i-1}, t_i]$ , we can generalize it with an appropriate constant  $f''(\eta)$  where  $a \leq \eta \leq b$ . Then the error expression comes to be:

$$E_{CM}(f) = \frac{f''(\eta)}{24} \sum_{i=1}^r h \cdot h^2$$

$$E_{CM}(f) = \frac{f''(\eta)}{24} h^2 \cdot rh$$

Since,  $h = \frac{b-a}{r}$ , so replacing  $(b-a) = rh$ , we get:

$$E_{CM}(f) = \frac{f''(\eta)}{24}(b-a)h^2 \quad (6)$$

This is the final expression for the error in the composite mid point rule.

(6) Suggests that the error varies proportional to  $h^2$ , hence it is second order accurate.

### Question-8: Chapter-15, question-13

Given that the interval of integration,  $[a,b]$ , is divided into equal sub-intervals of length  $h$ , such that  $r = \frac{b-a}{h}$

**Composite Simpson:**

$$\int_a^b f(x)dx \approx \frac{h}{3} \left[ f(a) + 2 \sum_{k=1}^{\frac{r}{2}-1} f(t_{2k}) + 4 \sum_{k=1}^{\frac{r}{2}} f(t_{2k-1}) + f(b) \right] \quad (7)$$

The expression for composite trapezoidal with step size  $h$  is given by:

**R2: Composite trapezoidal rule of step size  $h$**

$$\int_a^b f(x)dx \approx \frac{h}{2} \sum_{i=1}^r f(t_{i-1}) + f(t_i)$$

$$R_2 = \frac{h}{2} [f(a) + 2f(t_1) + 2f(t_2) + \cdots + 2f(t_{r-1}) + f(b)]$$

**R1: Composite trapezoidal rule of step size  $2h$**  For step-size of  $2h$ , we require even number of subintervals. In the above expression for summation, thus we change the summing variable  $i$  to  $2k$ , and the limit become  $\frac{r}{2}$ . Hence, we have:

$$R_1 = \frac{2h}{2} \sum_{k=1}^{\frac{r}{2}} f(t_{2k-2}) + f(t_{2k})$$

$$R_1 = h[\{f(t_0) + f(t_2) + \cdots + f(t_{r-2})\} + \{f(t_2) + f(t_4) + \cdots + f(t_r)\}]$$

Since,  $t_0$  and  $t_r$  are the two extreme end points of the interval, hence  $t_0 = a$  and  $t_r = b$  Thus, we get:

$$R_1 = h[f(a) + 2f(t_2) + 2f(t_4) + \cdots + 2f(t_{r-2}) + f(b)]$$

Hence, evaluating  $S = \frac{4R_2 - R_1}{3}$

$$4R_2 - R_1 = h[2f(a) + 4f(t_1) + 4f(t_2) + \cdots + 4f(t_{r-1}) + 2f(b)] - h[f(a) - 2f(t_2) - 2f(t_4) - \cdots - f(b)]$$

$$4R_2 - R_1 = h[f(a) + \{2f(t_2) + 2f(t_4) + \cdots + 2f(t_{r-2})\} + \{4f(t_1) + 4f(t_3) + \cdots + 4f(t_{r-1})\} + f(b)]$$

$$4R_2 - R_1 = h \left[ f(a) + 2 \sum_{k=1}^{\frac{r}{2}-1} f(t_{2k}) + 4 \sum_{k=1}^{\frac{r}{2}} f(t_{2k-1}) + f(b) \right]$$

$$\frac{4R_2 - R_1}{3} = \frac{h}{3} \left[ f(a) + 2 \sum_{k=1}^{\frac{r}{2}-1} f(t_{2k}) + 4 \sum_{k=1}^{\frac{r}{2}} f(t_{2k-1}) + f(b) \right]$$

The rhs of the above is exactly the expression for the composite Simpson's rule (7)

---

**Question-9: Chapter-15, question-14**


---

Code available at Prob9.m.

Here we implement the Romberg algorithm provided in the book. One heuristic choice made was to set the initial number of intervals to  $r = 1$ , such that it starts with the basic form of the trapezoid with  $h = (b - a)$ . Since, the algorithm is adaptive, refining the mesh into finer ones as at every stage  $h$  gets halved, so we started with a max  $h$ . Some experiments starting with a small  $h$  did not converge. We use the *roundn* function to get the desired 8digit limit for  $\pi$  and then keep doing the recursive step until that 8digit accuracy is met within the given tolerance. Below is the table given from the Romberg implementation:

3.0000000000000000	0	0	0	0
3.1000000000000000	3.13333333333333	0	0	0
3.131176470588235	3.141568627450980	3.142117647058823	0	0
3.138988494491089	3.141592502458707	3.141594094125888	3.141585783761874	0
3.140941612041389	3.141592651224822	3.141592661142563	3.141592638396796	3.141592665277717

Hence, the approximation upto 8 digits will be 3.1415926 .

The number of iterations required here is 4. Thus the number of function evaluations is roughly  $\sum_{j=1}^s 2^{j-1} \approx 2^s + 1 = 17$ .

The error for the composite trapezoid is second order accurate, that is, it goes down as  $h^2$ . Using , the extrapolation between different step sizes, we are able to cancel out leading error terms further, which results in faster convergence or less initial error. That is why, in some entries of the table, during the first fill-up of those entries, we reached the necessary convergence, thus no further improvements were seen in next iterations.