

## Q2

a)

We have that  $k(h) = \left| \frac{xg'(a, h)}{g(a, h)} \right|$  where the derivative of  $g(a, h)$  is  $g'(a, h) = \frac{1}{2\sqrt{a+h}} + \frac{1}{2\sqrt{a-h}} - \frac{\sqrt{a+h} - \sqrt{a-h}}{2h^2}$ . So  $k(h)$  is

$$\begin{aligned}
 k(h) &= \left| \frac{h \left( \frac{1}{2\sqrt{a+h}} + \frac{1}{2\sqrt{a-h}} - \frac{\sqrt{a+h} - \sqrt{a-h}}{2h^2} \right)}{\frac{\sqrt{a+h} - \sqrt{a-h}}{2h}} \right| \\
 &= \left| \frac{\frac{h}{2\sqrt{a+h}} + \frac{h}{2\sqrt{a-h}}}{\sqrt{a+h} - \sqrt{a-h}} - 1 \right| \\
 &= \left| \frac{\frac{h\sqrt{a-h} + h\sqrt{a+h}}{2\sqrt{(a-h)(a+h)}}}{\sqrt{a+h} - \sqrt{a-h}} - 1 \right| \\
 &= \left| \frac{h(\sqrt{a-h} - \sqrt{a+h})}{2(a+h)\sqrt{a-h} - 2(a-h)\sqrt{a+h}} - 1 \right| \\
 &= \left| \frac{h(\sqrt{a-h} - \sqrt{a+h})}{2((a+h)\sqrt{a-h} - (a-h)\sqrt{a+h})} \cdot \frac{(a+h)\sqrt{a-h} + (a-h)\sqrt{a+h}}{(a+h)\sqrt{a-h} + (a-h)\sqrt{a+h}} - 1 \right| \\
 &= \left| \frac{h(2(a^2 - h^2) + 2a\sqrt{a^2 - h^2})}{2(2a^2h - 2h^3)} - 1 \right| \\
 &= \left| \frac{a^2 - h^2 + a\sqrt{a^2 - h^2}}{2(a^2 - h^2)} - 1 \right| \\
 &= \left| \frac{1}{2} + \frac{a\sqrt{a^2 - h^2}}{2(a^2 - h^2)} - 1 \right| \\
 &= \left| \frac{a\sqrt{a^2 - h^2}}{2(a^2 - h^2)} - \frac{1}{2} \right| \\
 &= \left| \frac{a}{2\sqrt{a^2 - h^2}} - \frac{1}{2} \right|
 \end{aligned}$$

Here we can drop the absolute value sign as  $a > 0$  and  $\frac{a}{2\sqrt{a^2 - h^2}} - \frac{1}{2}$  doesn't go negative because the largest positive value of  $\sqrt{a^2 - h^2}$  is  $a$  itself (when  $h = 0$ ) so it gives that the lowest possible value for any given  $a$  is 0.

To study for what ranges of  $h$ ,  $k(h)$  becomes large. we'll consider a large denominator and then small denominator as the numerator does not change for a fixed  $a$ .

If the denominator is small (ie  $|h| \rightarrow a$ ), then  $\frac{a}{2\sqrt{a^2 - h^2}}$  is large (goes to infinity as  $h$  approaches  $|a|$ ), so  $\frac{a}{2\sqrt{a^2 - h^2}} - \frac{1}{2}$  is subsequently large (approaches to infinity as  $|h|$  approaches  $a$ ).

On the other hand, if the denominator is large (ie  $|h| \rightarrow 0$ ), then  $\frac{a}{2\sqrt{a^2 - h^2}}$  is small (goes to 0 as  $h$  approaches 0), so  $\frac{a}{2\sqrt{a^2 - h^2}} - \frac{1}{2}$  is subsequently small (approaches to 0 as  $|h|$  approaches 0).

Hence, the conditional number is large as  $|h| \rightarrow a$ , and  $k(h) \rightarrow 0$  for  $|h| \rightarrow 0$ , because for  $h = 0$ ,  $\frac{a}{2\sqrt{a^2 - h^2}} = \frac{a}{2\sqrt{a^2 - 0}} =$

$$\frac{a}{2\sqrt{a^2}} = \frac{a}{2a} = \frac{1}{2}.$$

And lastly, for a fixed  $h$  and for  $a \rightarrow \infty$ . As  $\lim_{a \rightarrow \infty} \sqrt{a^2 - h^2} = \lim_{a \rightarrow \infty} \sqrt{a^2} = \lim_{a \rightarrow \infty} a$ , hence  $\lim_{a \rightarrow \infty} \frac{a}{2\sqrt{a^2 - h^2}} - 1 = \lim_{a \rightarrow \infty} \frac{a}{2a} - 1 = \frac{1}{2} - \frac{1}{2} = 0$ . So  $k(h)$  approaches 0 for a fixed  $h$  and  $a$  that approaches infinity.

**b)**

The problematic ranges of  $h$  are:

1.  $|h| < \epsilon_{\text{mach}}$ . In this case, we have that  $h$  will be rounded down to 0, in which case, the function  $g(x)$  will be  $\frac{\sqrt{x+h} - \sqrt{x-h}}{2h} = \frac{\sqrt{a+0} - \sqrt{a-0}}{2 \cdot 0} = \frac{0}{0}$  which will result in an error and no computation of  $g(x)$  will be done, even though  $h$  itself is not 0.
2.  $x - h < \epsilon_{\text{mach}}$  or  $x + h < \epsilon_{\text{mach}}$ . Similar situation will happen here, when  $x - h$  or  $x + h$  will be rounded down to 0 if  $x$  is close enough to  $h$  as  $|x|$  will be rounded to  $|h|$ . So  $g(x)$  will be  $\frac{\sqrt{2h}}{2h}$  which mathematically is not the same as  $\frac{\sqrt{x+h} - \sqrt{x-h}}{2h}$  as  $|x| \neq |h|$  and would result in a wrong output.

**c)**

There is a function that approximates the wanted function  $g(x)$  better than just using the derivative while or giving rise to the issues mentioned in part (b). The function is

$$\begin{aligned} f(x) &= \frac{|h|}{x} \left( x - \frac{|h|}{x} \right) \left( x + \frac{|h|}{x} \right) \cdot \frac{\sqrt{x + \frac{|h|}{x}} - \sqrt{\max\left(0, x - \frac{|h|}{x}\right)}}{2 \frac{|h|}{x}} \\ &\quad + \left[ 1 - \frac{|h|}{x} \left( x - \frac{|h|}{x} \right) \left( x + \frac{|h|}{x} \right) \right] \cdot \frac{1}{2\sqrt{x}} \\ &= \left( x^2 - \frac{|h|^2}{x^2} \right) \cdot \frac{\sqrt{x + \frac{|h|}{x}} - \sqrt{\max\left(0, x - \frac{|h|}{x}\right)}}{2} + \left[ 1 - \frac{|h|}{x} \left( x^2 - \frac{|h|^2}{x^2} \right) \right] \cdot \frac{1}{2\sqrt{x}} \end{aligned}$$

Here, it's important to note that in the original function  $g$ , we have that if  $h = -h$ , then we have that  $\frac{\sqrt{x+(-h)} - \sqrt{x-(-h)}}{-2h} = \frac{\sqrt{x-h} - \sqrt{x+h}}{-2h} = \frac{\sqrt{x+h} - \sqrt{x-h}}{2h}$ .

For a fixed  $x$  value and when  $|h| \rightarrow 0$ , we have  $\lim_{|h| \rightarrow 0} f(x) = \frac{1}{2\sqrt{x}}$  because

$$\begin{aligned} &\frac{\sqrt{x + \frac{|h|}{x}} - \sqrt{\max\left(0, x - \frac{|h|}{x}\right)}}{2} \text{ will be 0 as } \sqrt{x + \frac{|h|}{x}} - \sqrt{\max\left(0, x - \frac{|h|}{x}\right)} \\ &= \frac{\sqrt{x} - \sqrt{x}}{2} = 0 \end{aligned}$$

For a fixed  $h$  and  $a \rightarrow \infty$ , then  $g(x) \rightarrow 0$ , and similarly,  $f(x) \rightarrow 0$ , because  $\lim_{x \rightarrow 0} \left( x^2 - \frac{|h|^2}{x^2} \right) = 0$  so the first term in  $f$  will be 0. On the other hand,  $\lim_{x \rightarrow 0} \left[ 1 - \frac{|h|}{x} \left( x^2 - \frac{|h|^2}{x^2} \right) \right] \cdot \frac{1}{2\sqrt{x}} = 1 \cdot \frac{1}{2\sqrt{x}} = 0$

Lastly, if  $x$  is fixed and  $|h| \rightarrow x$ , then we either have  $\lim_{|h| \rightarrow x} f(x) = (x^2 - 1) \frac{\sqrt{x+1} - \sqrt{x-1}}{2} + (2 - x^2) \cdot \frac{1}{2\sqrt{x}}$  if  $\max(0, x - h) = x - h$ , which happens when  $x$  is large, then  $(2 - x^2) \cdot \frac{1}{2\sqrt{x}}$  is negligible and approaches

0 if  $x$  approaches 0, so  $\lim_{|h| \rightarrow x} f(x) = (x^2 - 1) \frac{\sqrt{x+1} - \sqrt{x-1}}{2}$  and is an approximation for  $g(x)$ .

On the other hand, if  $x$  is small (less than 1), then  $\max(0, x - h) = 0$  and we have  $\lim_{|h| \rightarrow x} f(x) = (x^2 - 1) \frac{\sqrt{x+1}}{2} + (2 - x^2) \cdot \frac{1}{2\sqrt{x}}$  where,  $(x^2 - 1) \frac{\sqrt{x+1}}{2}$  is negligible as  $\left| (x^2 - 1) \frac{\sqrt{x+1}}{2} \right| < 1$ , while on the other hand,  $(2 - x^2) \cdot \frac{1}{2\sqrt{x}}$  will be large since  $\frac{1}{2\sqrt{x}} \rightarrow \infty$  as  $x \rightarrow 0$  and  $2 > (2 - x^2) > 1$ . Hence,  $\lim_{|h| \rightarrow x} f(x)$  for small  $x$  is an approximation for  $\frac{1}{2\sqrt{x}}$  which itself is an approximation for  $g(x)$ .

Lastly, we need to explain why  $h$  is in absolute value and why we have  $\max\left(0, x - \frac{|h|}{x}\right)$ .

1. We have  $|h|$  because as I had stated before, the value of  $g(x)$  doesn't change if we have  $h$  versus  $-h$ . And as I showed above, when  $h \geq 0$ ,  $f$  is a good approximation for  $g(x)$  with both positive and negative  $h$  values. So it suffices to take  $|h|$  to have  $f$  be a reasonable approximation of  $g(x)$ .
2. We have  $\max\left(0, x - \frac{|h|}{x}\right)$  because for certain reasonable values of  $x$  and  $h$ ,  $x - \frac{|h|}{x}$  is negative so the function is undefined (e.g.  $x = 0.7, h = 0.6$ ). This happens only for relatively small  $x$  values where  $|h| \rightarrow x$  but as discussed above, the first term of  $f$  becomes small and negligible, so the function value becomes close to  $\frac{1}{2\sqrt{x}}$  which is an approximation of  $g$ . So to combat the issue of  $f$  being undefined for those values of  $x$ , and  $h$ , we just set the value of the square root to 0.

Hence,  $f(x) = \left(x^2 - \frac{|h|^2}{x^2}\right) \cdot \frac{\sqrt{x + \frac{|h|}{x}} - \sqrt{\max\left(0, x - \frac{|h|}{x}\right)}}{2} + \left[1 - \frac{|h|}{x} \left(x^2 - \frac{|h|^2}{x^2}\right)\right] \cdot \frac{1}{2\sqrt{x}}$  is another way of reasonable approximation.

On another note,  $f$  is not subject to issues noted in part *b* with  $g$  as if  $|h| < \epsilon_{\text{mach}}$  then the value of  $f$  will just be  $\frac{1}{2\sqrt{x}}$  which we know is a reasonable approximation for  $g$ , and similarly, if  $x - h < \epsilon_{\text{mach}}$  or  $x + h < \epsilon_{\text{mach}}$ , then the value of  $f$  (and subsequently the approximate value of  $g$ ) would be calculated using the same logic as in the case of  $\lim_{|h| \rightarrow x} f(x)$ .

So neither one of the issues from (b) happen with  $f$ .

d)

All code in this assignment is written in Python.

```
import matplotlib.pyplot as plt
import numpy as np
from math import sqrt

h_list = []
for i in range(1, 18):
    h_list.append( 10**(-1*i) )

def der_sqrt(x):
    '''derivative of sqrt(x) '''
    return 1/(2 * sqrt(x))

def g(x,h):
    '''straight away calculation of g(x)'''
    g = (sqrt(x+h) - sqrt(x-h)) / (2*h)
    return g

def f(x, h, part_term_2 = der_sqrt):
    '''
    Function built according to Q2 part c
    0 <= |h| < x
    '''
    if 0 <= abs(h) < x:
        part_coeff = x**2 - (abs(h)/x)**2
        rad = sqrt(x + abs(h)/x)
        max_rad = sqrt(max(0, x - abs(h)/x ))
        term_1 = part_coeff * (rad - max_rad)/2
        term_2 = (1 - (abs(h)/x)*part_coeff)*part_term_2(x)
        tot = term_1 + term_2
        return term_1 + term_2
    else:
        raise ValueError('Incorrect input')

def output(x, h_list=h_list, der=der_sqrt, g=g, f=f):
    '''Returns the needed result for each
    i as a tuple '''
    for h in h_list:
        i = h_list.index(h) + 1
        d0 = der(x)
        d1 = g(x, h)
        d2 = f(x, h)
        delta1 = d0-d1
        delta2 = d0-d2

        yield i, d0, d1, d2, abs(delta1), abs(delta2)

def print_result(num, gen=output):
    '''Prints every tuple from output
    and plots the graphs'''
    delta1_list = []
    delta2_list = []

    for e in gen(num):
```

```

print(e)
delta1_list.append(abs(e[4]))
delta2_list.append(abs(e[5]))

```

```

delta1_list.reverse()
delta2_list.reverse()
h_list.reverse()

```

```

x = np.array(h_list)
y = np.array(delta1_list)
z = np.array(delta2_list)
default_x_ticks = range(len(x))
plt.plot(default_x_ticks, y, linestyle='solid')
plt.plot(default_x_ticks, z, linestyle='dashed')
plt.xticks(default_x_ticks, x)
plt.xlabel('h values')
plt.ylabel('Error magnitude')
plt.legend(['d0-d1', 'd0-d2'])
plt.yscale("log")
plt.show()

```

# To run code, call  
print\_result() # input some number e.g. 1 or 10

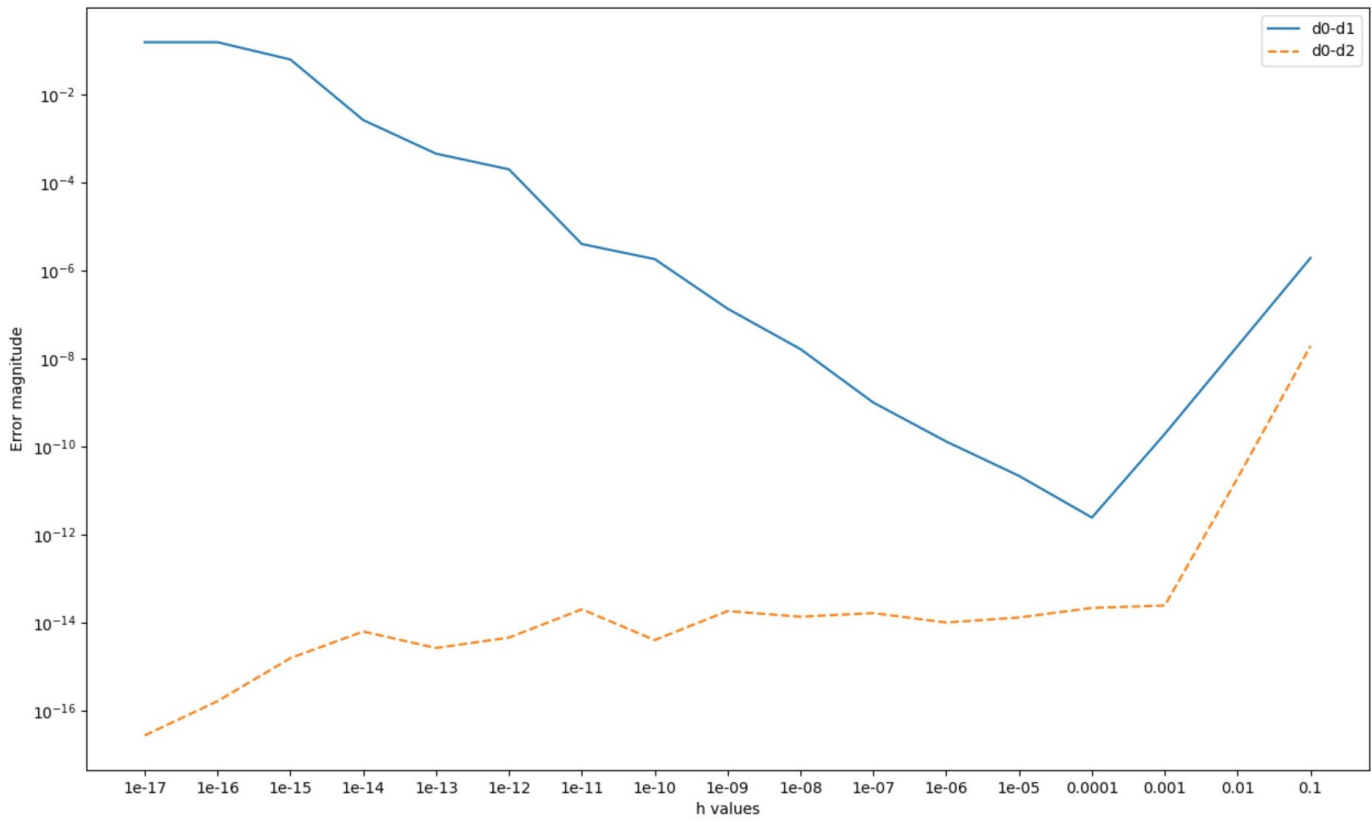
# Run each print\_result separately.  
# Do not write one after the other.

The results for  $x = 10$  are:

```

(1, 0.15811388300841897, 0.15811585951842844, 0.15811390277265214, 1.9765100094659704e-06, 1.9764233166741363e-08)
(2, 0.15811388300841897, 0.1581139027726719, 0.1581138830281628, 1.97642529287112e-08, 1.9743817691875165e-11)
(3, 0.15811388300841897, 0.15811388320585706, 0.15811388300844378, 1.974380936520248e-10, 2.481348460037225e-14)
(4, 0.15811388300841897, 0.1581138830109019, 0.15811388300839702, 2.4829305278473157e-12, 2.195466031196247e-14)
(5, 0.15811388300841897, 0.158113882986477, 0.15811388300840568, 2.1941976013906128e-11, 1.329492071988625e-14)
(6, 0.15811388300841897, 0.1581138828754547, 0.15811388300840876, 1.3296427847642178e-10, 1.021405182655144e-14)
(7, 0.15811388300841897, 0.15811388198727627, 0.1581138830084357, 1.021142698176547e-09, 1.6736612096224235e-14)
(8, 0.15811388300841897, 0.15811389975084467, 0.15811388300840512, 1.6742425695825958e-08, 1.3850032232198828e-14)
(9, 0.15811388300841897, 0.15811374431962122, 0.15811388300843757, 1.3868879775169596e-07, 1.8596235662471372e-14)
(10, 0.15811388300841897, 0.15811574272106554, 0.15811388300842305, 1.8597126465735858e-06, 4.08006961549745e-15)
(11, 0.15811388300841897, 0.1581179631671148, 0.15811388300843937, 4.080158695823899e-06, 2.040034807748725e-14)
(12, 0.15811388300841897, 0.15831780331154732, 0.15811388300841434, 0.00020392030312835208, 4.6351811278100286e-15)
(13, 0.15811388300841897, 0.15765166949677223, 0.15811388300841628, 0.00046221351164674185, 2.6922908347160046e-15)
(14, 0.15811388300841897, 0.15543122344752192, 0.15811388300842535, 0.002682659560897055, 6.38378239159465e-15)
(15, 0.15811388300841897, 0.22204460492503128, 0.1581138830084174, 0.06393072191661231, 1.582067810090848e-15)
(16, 0.15811388300841897, 0.0, 0.1581138830084188, 0.15811388300841897, 1.6653345369377348e-16)
(17, 0.15811388300841897, 0.0, 0.15811388300841894, 0.15811388300841897, 2.7755575615628914e-17)

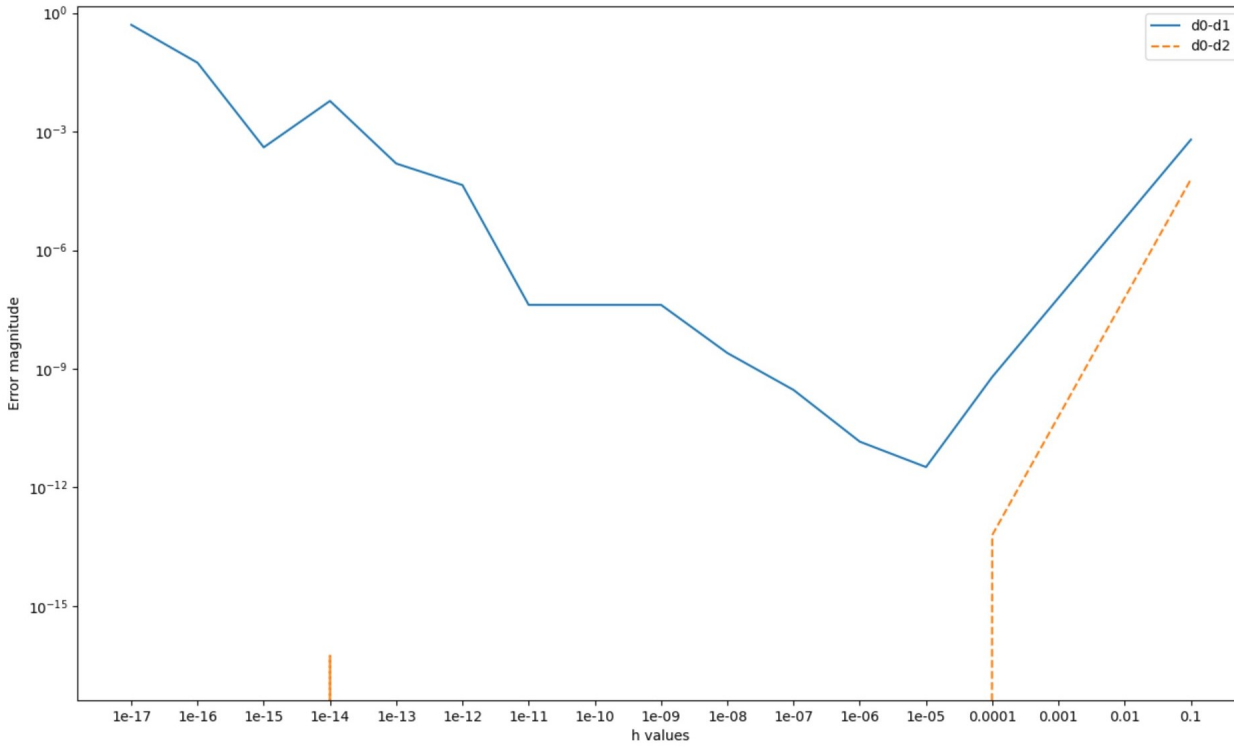
```



While the the results for  $x = 1$  are:

Note that the dotted graph extends down to 0 for  $h$  values of  $10^{-17}$ ,  $10^{-16}$ ,  $10^{-15}$  and from  $10^{-13}$   $10^{-5}$ . But due to the specifications of the graph plotting python library, zooming out is not an option.

```
(1, 0.5, 0.5006277505981893, 0.5000621473092207, 0.0006277505981893139, 6.214730922071698e-05)
(2, 0.5, 0.5000062502734492, 0.5000000624964842, 6.250273449248667e-06, 6.249648420997289e-08)
(3, 0.5, 0.5000000625000056, 0.5000000000624999, 6.25000056153624e-08, 6.249989414897072e-11)
(4, 0.5, 0.5000000006244454, 0.5000000000000624, 6.244453842896291e-10, 6.23945339839338e-14)
(5, 0.5, 0.5000000000032756, 0.5, 3.275602011854062e-12, 0.0)
(6, 0.5, 0.5000000000143778, 0.5, 1.4377832258105627e-11, 0.0)
(7, 0.5, 0.5000000002919336, 0.5, 2.9193358841439476e-10, 0.0)
(8, 0.5, 0.5000000025123796, 0.5, 2.512379637664708e-09, 0.0)
(9, 0.5, 0.5000000413701855, 0.5, 4.137018549954519e-08, 0.0)
(10, 0.5, 0.5000000413701855, 0.5, 4.137018549954519e-08, 0.0)
(11, 0.5, 0.5000000413701855, 0.5, 4.137018549954519e-08, 0.0)
(12, 0.5, 0.5000444502911705, 0.5, 4.445029117050581e-05, 0.0)
(13, 0.5, 0.500155472593633, 0.5, 0.00015547259363302146, 0.0)
(14, 0.5, 0.49404924595819466, 0.49999999999999994, 0.0059507540418053395, 5.551115123125783e-17)
(15, 0.5, 0.4996003610813204, 0.5, 0.0003996389186796123, 0.0)
(16, 0.5, 0.5551115123125783, 0.5, 0.05511151231257827, 0.0)
(17, 0.5, 0.0, 0.5, 0.5, 0.0)
```



As we can see from the graph of  $|d0 - d1|$  for  $x = 10$ , the magnitude of error decreases sharply from about 0.15811 at  $10^{-17}$  and  $10^{-16}$  down to about  $2.4829 \times 10^{-12}$  at  $10^{-4}$  and from there it later increases to  $1.9765 \times 10^{-6}$ . While  $|d0 - d2|$  stays under the graph of  $|d0 - d1|$  at all times. Its values increase from  $2.7755 \times 10^{-17}$  at  $10^{-17}$  to about  $6.3837 \times 10^{-15}$  at  $10^{-14}$ , stays stagnant until 0.001 and after it increases to  $1.9764 \times 10^{-8}$ .

On the other hand, when  $x = 1$ , exactly the same behaviour is noticed for  $|d0 - d1|$  with the only difference with  $|d0 - d2|$  being that the error magnitude decreases less smoothly with certain  $h$  values where it falls and increases sharply, notable at  $h = 10^{-16}, 10^{-15}, 10^{-12}$ .

For the graph of  $|d0 - d2|$ , the situation is significantly more interesting as most of the values of the graph are in fact 0, meaning that there is no error. This is because when the value of  $h$  is too small that it ends up being rounded down to 0, or the magnitude of the error is rounded to 0. When the value of  $h$  is being rounded to 0, we end up having that is just  $f(x) = \frac{1}{2\sqrt{x}}$  itself, so there no error in this case.