



UNIVERSITY OF WATERLOO
FACULTY OF ENGINEERING
Department of Electrical &
Computer Engineering



ECE 204 *Numerical methods*

Tools for numerical algorithms

Douglas Wilhelm Harder, LEL, M.Math.

dwharder@uwaterloo.ca

dwharder@gmail.com





Introduction

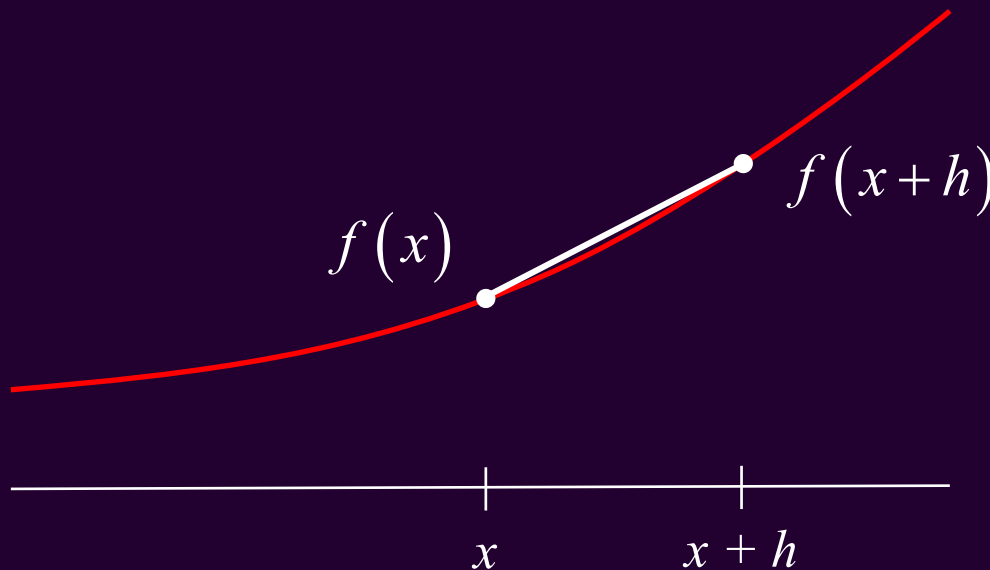
- In this topic, we will
 - Review the issues with using analytic formulas with numerical values
 - Look at two formulas that approximate the same value, but where one is significantly better
 - Describe the tools that we will use in this course to come up with numerical algorithms



Approximating the derivative

- In calculus, you saw the definition of the derivative is

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

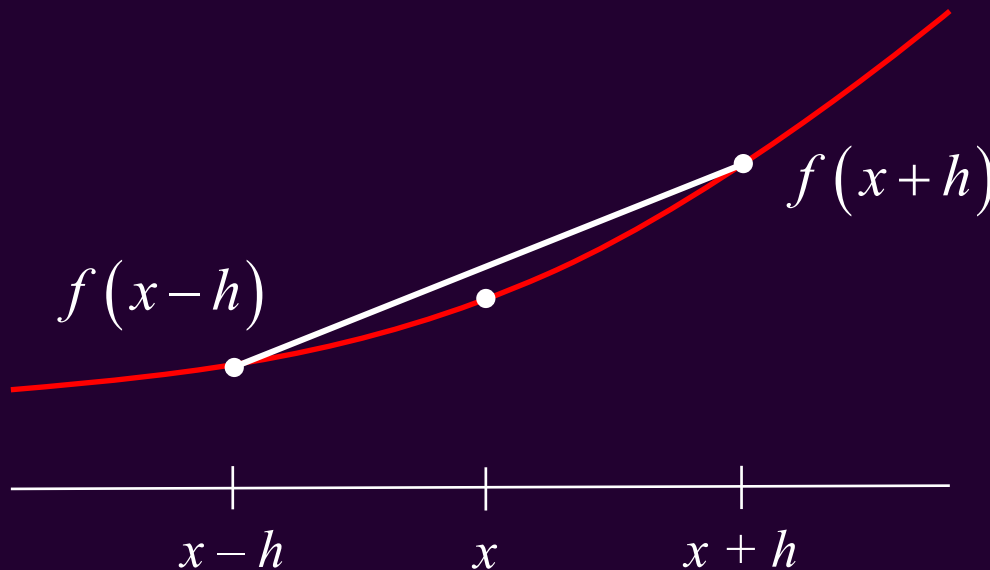




Approximating the derivative

- Claim, here is an alternate definition of the derivative:

$$\frac{d}{dx} f(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x-h)}{2h}$$





Approximating the derivative

- You could prove the definitions are equivalent mathematically:

$$\begin{aligned}\frac{d}{dx} f(x) &= \lim_{h \rightarrow 0} \frac{f(x+h) - f(x-h)}{2h} \\ &= \lim_{h \rightarrow 0} \frac{f(x+h) - f(x) + f(x) - f(x-h)}{2h} \\ &= \frac{1}{2} \lim_{h \rightarrow 0} \left(\frac{f(x+h) - f(x)}{h} + \frac{f(x) - f(x-h)}{h} \right) \\ &= \frac{1}{2} \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} + \frac{1}{2} \lim_{h \rightarrow 0} \frac{f(x) - f(x-h)}{h} \\ &= \frac{1}{2} \frac{d}{dx} f(x) + \frac{1}{2} \frac{d}{dx} f(x) = \frac{d}{dx} f(x)\end{aligned}$$



Approximating the derivative

- Let us approximate the derivative of $\sin(x)$ at $x = 1$ and see what happens in MATLAB:

```
>> cos(1) % The correct answer
```

```
ans =
```

```
0.540302305868140
```

```
>> % This is a for loop with 'n'
```

```
>> % going from 0 to 17
```

```
>> for n = 0:17
```

```
    h = 10.0^-n; % h = 1.0, 0.1, 0.01, ..
```

```
    (sin(1.0 + h) - sin(1.0 - h))/(2.0*h)
```

```
end
```

```
0.454648713412841  
0.539402252169760  
0.540293300874733  
0.540302215817690  
0.540302304967710  
0.540302305856999  
0.540302305895857  
0.540302305673812  
0.540302308449370  
0.540302302898255  
0.540302247387103  
0.540301137164079  
0.540290034933832  
0.540123501480139  
0.544009282066327  
0.555111512312578  
0.555111512312578  
0
```



Approximating the derivative

- Issue: this tells us something is happening, but not what is happening...
 - If we use `format hex`, we can look at the individual bits of the approximation
 - To most clearly see what is going on, however, it's better to use $h = 2^{-n}$

```
>> format bin
>> for n = 1:53
    h = 2.0^-n; % h = 0.5, 0.25, 0.125, ..
    (sin(1.0 + h) - sin(1.0))/h
end
```

```
0.454648713412841
0.539402252169760
0.540293300874733
0.540302215817690
0.540302304967710
0.540302305856999
0.540302305895857
0.540302305673812
0.540302308449370
0.540302302898255
0.540302247387103
0.540301137164079
0.540290034933832
0.540123501480139
0.544009282066327
0.555111512312578
0.555111512312578
0
```



2^{-1} 0.01001111111000100110000011000010001101100000100111010
 2^{-2} 0.01101110000110000000110111100001000001101001001111000
 2^{-3} 0.01111100100000101110111011000100111010000011100000000
 2^{-4} 0.100000011011111101101110100011011101011101001011110000
 2^{-5} 0.1000011011101101111001001000111011101111011010000000
 2^{-6} 0.10001000101000001111100100110111010001101000010000000
 2^{-7} 0.1000100101111001011110011110101011100101100011000000
 2^{-8} 0.1000100111100101011101000010011011011100010000000000
 2^{-9} 0.10001010000110110110000000010010010001101100000000000
 2^{-10} 0.1000101000110110010100011011100001100100011000000000
 2^{-11} 0.100010100100001111001001011101100101111000000000000
 2^{-12} 0.100010100100101010000101000100010111011100000000000
 2^{-13} 0.100010100100110111100010110011010101001100000000000
 2^{-14} 0.100010100100111110010001101001101110111000000000000
 2^{-15} 0.100010100101000001101001000100101010100000000000000
 2^{-16} 0.100010100101000011010101100100001000000000000000000
 2^{-17} 0.100010100101000100001010101000110000000000000000000
 2^{-18} 0.100010100101000100100101100100000110000000000000000
 2^{-19} 0.100010100101000100110011000001110000000000000000000
 2^{-20} 0.100010100101000100111001110000101000000000000000000
 2^{-21} 0.1000101001010001001110100100000000000000000000000000
 2^{-22} 0.1000101001010001001111101100111000000000000000000000
 2^{-23} 0.1000101001010001001111111010100000000000000000000000
 2^{-24} 0.1000101001010001010000000010000000000000000000000000
 2^{-25} 0.1000101001010001010000000101000000000000000000000000
 2^{-26} 0.10001010010100010100000001100000000000000000000000000
 2^{-27} 0.10001010010100010100000010000000000000000000000000000
 2^{-28} 0.10001010010100010100000010000000000000000000000000000
 2^{-29} 0.10001010010100010100000000000000000000000000000000000
 2^{-30} 0.10001010010100010100000000000000000000000000000000000
 2^{-31} 0.10001010010100010100000000000000000000000000000000000
 2^{-32} 0.10001010010100010100000000000000000000000000000000000
 2^{-33} 0.10001010010100010100000000000000000000000000000000000
 2^{-34} 0.10001010010100010100000000000000000000000000000000000
 2^{-35} 0.10001010010100010100000000000000000000000000000000000
 2^{-36} 0.10001010010100011000000000000000000000000000000000000
 2^{-37} 0.10001010010100010000000000000000000000000000000000000
 2^{-38} 0.10001010010100010000000000000000000000000000000000000
 2^{-39} 0.100010100101000
 2^{-40} 0.100010100101000
 2^{-41} 0.100010100101000
 2^{-42} 0.1000101001100
 2^{-43} 0.1000101001000
 2^{-44} 0.100010100
 2^{-45} 0.100010100
 2^{-46} 0.100010100
 2^{-47} 0.100011000
 2^{-48} 0.1000100
 2^{-49} 0.1001000
 2^{-50} 0.100
 2^{-51} 0.100
 2^{-52} 0.100
 2^{-53} 0.000

$$\frac{\sin(1+h) - \sin(1)}{h}$$



Tools for numerical algorithms



2⁻¹ 0.010011111100010011000011000010001101100000100111010
 2⁻² 0.011011100001100000001101111000010000011010010011111000
 2⁻³ 0.0111110010000010111011101100010011101000001110000000
 2⁻⁴ 0.1000001101111110110111010001011101011101001011110000
 2⁻⁵ 0.1000011011101101111001001000111011101111011010000000
 2⁻⁶ 0.1000100010100000111110010011011101000110100001000000
 2⁻⁷ 0.10001001011110010111100111101010111001011000110000000
 2⁻⁸ 0.1000100111100101011101000010011011011100010000000000
 2⁻⁹ 0.10001010000110110110000000010010010001101100000000000
 2⁻¹⁰ 0.10001010001101100101000110111000011001000110000000000
 2⁻¹¹ 0.10001010010000111100100111100100111100000000000000
 2⁻¹² 0.1000101001001010100001010001000101110111100000000000
 2⁻¹³ 0.1000101001001101111000101100110101010011000000000000
 2⁻¹⁴ 0.100010100100111110010001101001101110111000000000000
 2⁻¹⁵ 0.100010100101000001101001000100101010100000000000000
 2⁻¹⁶ 0.100010100101000011010101100100011000000000000000000
 2⁻¹⁷ 0.100010100101000010001010101000110000000000000000000
 2⁻¹⁸ 0.100010100101000010010110010000011000000000000000000
 2⁻¹⁹ 0.100010100101000010011001100000111000000000000000000
 2⁻²⁰ 0.10001010010100001001110011100001010000000000000000
 2⁻²¹ 0.10001010010100001001110100100000000000000000000000
 2⁻²² 0.10001010010100001001111011001110000000000000000000
 2⁻²³ 0.10001010010100001001111110101000000000000000000000
 2⁻²⁴ 0.10001010010100001000100000000000000000000000000000
 2⁻²⁵ 0.10001010010100001000000000000000000000000000000000
 2⁻²⁶ 0.10001010010100001000000000000000000000000000000000
 2⁻²⁷ 0.10001010010100001000000100000000000000000000000000
 2⁻²⁸ 0.100010100101000010000001000000000000000000000000000
 2⁻²⁹ 0.100010100101000010000000000000000000000000000000000
 2⁻³⁰ 0.100010100101000010000000000000000000000000000000000
 2⁻³¹ 0.100010100101000010000000000000000000000000000000000
 2⁻³² 0.100010100101000010000000000000000000000000000000000
 2⁻³³ 0.100010100101000010000000000000000000000000000000000
 2⁻³⁴ 0.100010100101000010000000000000000000000000000000000
 2⁻³⁵ 0.100010100101000010000000000000000000000000000000000
 2⁻³⁶ 0.1000101001010001100000000000000000000000000000000000
 2⁻³⁷ 0.1000101001010001000000000000000000000000000000000000
 2⁻³⁸ 0.1000101001010010000000000000000000000000000000000000
 2⁻³⁹ 0.10001010010100
 2⁻⁴⁰ 0.10001010010100
 2⁻⁴¹ 0.10001010010100
 2⁻⁴² 0.10001010011000
 2⁻⁴³ 0.100010100100
 2⁻⁴⁴ 0.1000101000
 2⁻⁴⁵ 0.1000101000
 2⁻⁴⁶ 0.1000101000
 2⁻⁴⁷ 0.10001100
 2⁻⁴⁸ 0.10001000
 2⁻⁴⁹ 0.100100
 2⁻⁵⁰ 0.1000
 2⁻⁵¹ 0.1000
 2⁻⁵² 0.1000
 2⁻⁵³ 0.00

2⁻¹ 0.100001001010000000110011000010000011011101010010011
 2⁻² 0.10001000111000011000111001110011010100011011001001111110
 2⁻³ 0.10001001111101010001110011000010001001000000010110000
 2⁻⁴ 0.10001010001110100011010000011001111101010100111010000
 2⁻⁵ 0.10001010010010101101111010010110101101010101101100000
 2⁻⁶ 0.1000101001001111110011111001100010001101001110000000
 2⁻⁷ 0.1000101001010000111001000100011110001111101010000000
 2⁻⁸ 0.1000101001010001001010010111000000111101001110000000
 2⁻⁹ 0.10001010010100010011101010111010010001011001100000000
 2⁻¹⁰ 0.1000101001010001001111110000110011001111100000000000
 2⁻¹¹ 0.1000101001010001000101000000001000001011100100000000000
 2⁻¹² 0.100010100101000101000000000000000011001101001101011000000000
 2⁻¹³ 0.1000101001010001010000000000000000110111110010100000000000
 2⁻¹⁴ 0.100010100101000101000000000000000011111000011011100000000000
 2⁻¹⁵ 0.100010100101000101000000000000000011110101001100000000000000
 2⁻¹⁶ 0.1000101001010001010000000000000000111101010000000000000000
 2⁻¹⁷ 0.1000101001010001010000000000000000111110110101000000000000
 2⁻¹⁸ 0.1000101001010001010000000000000000111110110100000000000000
 2⁻¹⁹ 0.1000101001010001010000000000000000111101101000000000000000
 2⁻²⁰ 0.1000101001010001010000000000000000111101110000000000000000
 2⁻²¹ 0.1000101001010001010001010000000000111101100000000000000000
 2⁻²² 0.1000101001010001010000000000000000111110100000000000000000
 2⁻²³ 0.1000101001010001010000000000000000111111000000000000000000
 2⁻²⁴ 0.100010100101000101000000000000000011111000000000000000000000
 2⁻²⁵ 0.100010100101000101000000000000000010000000000000000000000000
 2⁻²⁶ 0.10001010010100010100010100000000000000000000000000000000000
 2⁻²⁷ 0.10001010010100010100000100000000000000000000000000000000000
 2⁻²⁸ 0.10001010010100010100000100000000000000000000000000000000000
 2⁻²⁹ 0.100010100101000101000
 2⁻³⁰ 0.100010100101000101000
 2⁻³¹ 0.10001010010100010100010100000000000000000000000000000000000
 2⁻³² 0.100010100101000101000
 2⁻³³ 0.100010100101000101000
 2⁻³⁴ 0.100010100101000101000
 2⁻³⁵ 0.100010100101000101000
 2⁻³⁶ 0.1000101001010001100
 2⁻³⁷ 0.1000101001010001000
 2⁻³⁸ 0.10001010010100100
 2⁻³⁹ 0.100010100101000
 2⁻⁴⁰ 0.100010100101000
 2⁻⁴¹ 0.100010100101000
 2⁻⁴² 0.1000101001100
 2⁻⁴³ 0.1000101001000
 2⁻⁴⁴ 0.100010100
 2⁻⁴⁵ 0.100010100
 2⁻⁴⁶ 0.100010100
 2⁻⁴⁷ 0.100011000
 2⁻⁴⁸ 0.1000100
 2⁻⁴⁹ 0.1001000
 2⁻⁵⁰ 0.100
 2⁻⁵¹ 0.100
 2⁻⁵² 0.100
 2⁻⁵³ 0.100



Tools for numerical algorithms

- This is a phenomena known as *subtractive cancellation*

- If we subtract two similar floating-point numbers, the result will have less precision than either operand

$$\frac{\sin(1.001) - \sin(1)}{0.001} = \frac{0.8420108663 - 0.8414709848}{0.001}$$

$$= \frac{0.0005398815}{0.001} = 0.5398815 \quad \cos(1) \approx 0.5403$$

0.07788 % relative error

- But we are only store four digits:

$$\frac{\sin(1.001) - \sin(1)}{0.001} = \frac{0.8420 - 0.8415}{0.001}$$

$$= \frac{0.0005}{0.001} = 0.5000$$

7.459 % relative error

- To get 0.0005399, $\sin(1.001)$ and $\sin(1)$ must be calculated to 7 digits



Tools for numerical algorithms

- How do we find numerical algorithms that help us get better approximations?
 - We will introduce, review and explore seven tools that we will choose from:
 1. Weighted averages
 2. Iteration
 3. Linear algebra
 4. Interpolation
 5. Taylor series
 6. Bracketing
 7. Intermediate-value theorem



Summary

- Following this topic, you now
 - Have seen how the binary representation appears to cause issues with what should be accurate formula
 - Know the seven tools we will use in this course
 - Understand what we will cover in the next seven topics



References

- [1] <https://en.wikipedia.org/wiki/Derivative>



Acknowledgments

Kevin Lee for noting an error on Slide 5.

Hassaan Ali Qazi for noting errors on Slides 3 and 4.



Colophon

These slides were prepared using the Cambria typeface. Mathematical equations use Times New Roman, and source code is presented using Consolas. Mathematical equations are prepared in MathType by Design Science, Inc. Examples may be formulated and checked using Maple by Maplesoft, Inc.

The photographs of flowers and a monarch butter appearing on the title slide and accenting the top of each other slide were taken at the Royal Botanical Gardens in October of 2017 by Douglas Wilhelm Harder. Please see

<https://www.rbg.ca/>

for more information.





Disclaimer

These slides are provided for the ECE 204 *Numerical methods* course taught at the University of Waterloo. The material in it reflects the author's best judgment in light of the information available to them at the time of preparation. Any reliance on these course slides by any party for any other purpose are the responsibility of such parties. The authors accept no responsibility for damages, if any, suffered by any party as a result of decisions made or actions based on these course slides for any other purpose than that for which it was intended.