Finite precision arithmetic in programming: Tips and pitfalls

Goal of this lecture

- · See some examples of what can go wrong with flooting point orithmetic.
- · Learn how to avoid some common mistakes.

Why do we wary about all of this?

- 1) Consequences of round-off error

 - · O.1+O.1+O.1 · · · · O.1 # I why

 10 times

 6 New compare real numbers

 With egudity
 - · Choosing tolerances for stopping
- (2) Cata strophic concellations

$$\frac{(1+x)-1}{x} + 1 - why not?$$

- · Quadratic formula
- 3) Floating point anthmetic is not?

$$6.2 + (0.3 + 0.1) \neq (0.2 + 0.3) + 0.1$$

1) Consequences of Round-off emor

Recall: Round-off error is the error that is made when we represent real numbers using finite precision anothmetic.

Example: X=0.1 cannot be represented exactly on the computer.

 $X = 0.1 = 1.10011001100110011 \times 2^{-4}$

How do we see the effects of this?

Python

x = 0.100000000000000555

Round-off error = 10-17

```
1 print("x+x
                           = \{:30.20f\}".format(x+x))
In [4]:
        2 print("x+x+x
                           = \{:30.20f\}".format(x+x+x))
        print("x+x+x+x
                           = \{:30.20f\}".format(x+x+x+x))
        ^{4} print("x+x+x+x+x = {:30.20f}".format(x+x+x+x+x))
                                                   flo-bruon
                           0.20000000000000001110
       x+x
                           0.30000000000000004441
       x+x+x
                                                    CION.
                           0.40000000000000002220
       x+x+x+x
                           0.50000000000000000000
       x+x+x+x+x =
```

```
In [8]: 

1 def roundoff():

Python
```

take home message: Avoid equality comparisons between real numbers

Example: Algorithm approximate the Square root Ja.

Using matplotlib backend: nbAgg
Populating the interactive namespace from numpy and matplotlib

```
In [23]:
         1 def squareroot(a):
              x = 1
              tol = 1e-8
             for k in range(20):
                  xp = (x + a/x)/2
                 if (abs(xp-x) < tol):</pre>
                       x = xp
                       return x
                  x = xp
        10
              return x
        ^{11} a = 5
        12 x = squareroot(a)
        print("Computed : {:.20f}".format(x))
        14 print("True : {:.20f}".format(sqrt(a)))
```

Computed: 2.23606797749978980505 True: 2.23606797749978980505

Use tolerances rather than exact equality.

Exemple Use Taylor Series to approximate Cos(x): $({}^{6}OS(\times))$: $COS(\times) = 1 - \frac{\chi^{2}}{2!} + \frac{\chi^{4}}{4!} - \frac{\chi^{6}}{6!} + \cdots = \frac{2}{n=0} (-1) \frac{\chi^{2}}{(2n)!}$ def taylor(x): T = 0for k in range(15): $T \leftarrow (-1)^{**}(k)^{*}x^{**}(2*k)/math.factorial(2*k)$ err = abs(cos(x)-T)print("{:5d} {:30.20f} {:12.4e}".format(k,T,err)) if cos(x) == T: printf('{:d} terms are needed'.format(k)) return T 10 print("No equality!") 11 return T 13 taylor(1.1) 0 1.000000000000000000000 5.4640e-01 1 0.3949999999999999674 5.8596e-02 2.4080e-03 0.45600416666666659937 0.45354366527777772999 5.2456e-05 0.45359682968278763893 7.0826e-07 0.45359611491689805218 6.5087e-09 0.45359612146891870044 4.3341e-11 0.45359612142535854495 2.1877e-13 0.45359612142557814707 8.3267e-16 Ovestion: 9 0.45359612142557725889 5.5511e-17 10 0.45359612142557725889 5.5511e-17 11 0.45359612142557725889 5.5511e-17 12 0.45359612142557725889 5.5511e-17 13 0.45359612142557725889 5.5511e-17 Stop here? 14 0.45359612142557725889 5.5511e-17 No equality! 0.45359612142557726 This loop would continue to

Taylor Series: Better:

0.4535961214689187

```
def taylor_better(x):
      tol = 1e-8
      T = 0
      for k in range(20):
          term = (-1)**(k)*x**(2*k)/math.factorial(2*k)
          T += term
          print("{:5d} {:30.20f} {:12.4e}".format(k,T,abs(term)))
          if abs(term) < tol:</pre>
               print('{:d} terms are needed'.format(k))
               return T
      print("No equality!")
      return T
14 taylor_better(1.1)
              1.0000000000000000000000
                                        1.0000e+00
    1
              0.3949999999999990674
                                        6.0500e-01
              0.45600416666666659937
                                        6.1004e-02
              0.45354366527777772999
                                        2.4605e-03
              0.45359682968278763893
                                        5.3164e-05
              0.45359611491689805218
                                        7.1477e-07
              0.45359612146891870044
                                        6.5520e-09
6 terms are needed
```

Stop when a reasonable tolerance has been achieved.

```
Example: Geometric Series

Compute S_n = \sum_{K=0}^{K} X

= \{ + X + X^2 + \cdots + X = \frac{|-X|}{|-X|} \}
```

```
def geoseries(x,n):
    strue = (1-x**n)/(1-x)
    S = 0
    for i in range(n):
        S = S*x + 1
    if (S == strue):
        print("Exact sum! S = {:24.20f}".format(S))
    else:
        print("Close : S = {:24.20f}".format(S))
        print("Close : error = {:24.4e}".format(abs(S-strue)))
```

```
geoseries(7.3,5)
```

```
Close : S = 3290.43110000000024228939
Close : error = 4.5475e-13
```

geoseries(7.1,5)

Exact sum! S = 2957.58909999999968931661

Since we don't know if we can
comparte a value exactly, we should
never use exact equality compansons.
with real numbers

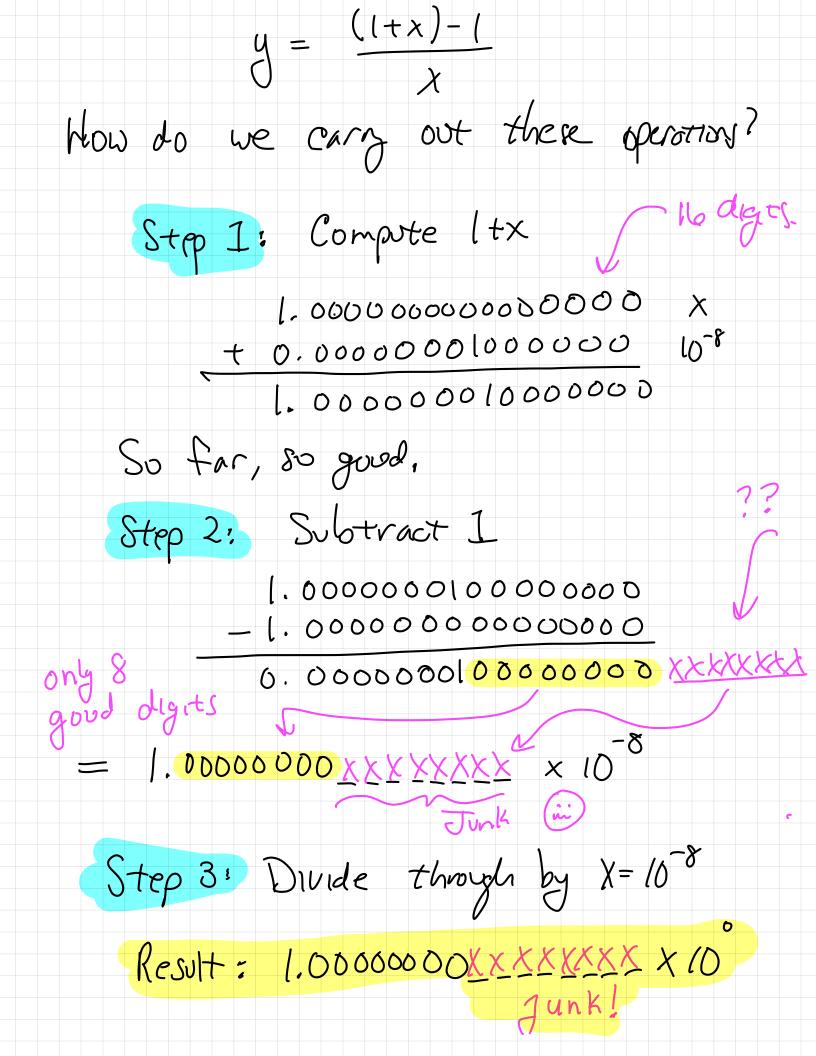
2 Catastrophic Cancellation

try the following:
$$y = \frac{1+x}{-1} \text{ for } x = \frac{-8}{10}$$
We expect to get $y = 1$. But instead, we get:

```
1 x = 1e-8
2 y = ((1 + x) - 1)/x
3 print('y = {:24.20f}'.format(y))
```

y = 0.9999999392252902908

What goes wrong?



Problem occurs because we try to Subtract two values that che very close in magnitude 1+x & 1, x <<1 This results in a "catastrophic" loss of accoracy. Matters one only made worse by multiplying by a lorge number. Example: Use a finite difference
formula to compute the derivative of f(x) = x: $f'(x) \approx \frac{f(x+h)-f(x)}{h}$ h <<1 his very small. From Calculus.

We know that f'(x) = 1.
But due to catastrophic cancellation
for small values of h, we get

or on error of about 10.

This loss of accuracy limits the usefulness of finite difference formulais for very small values of h.

Example: The quadratic formula: Solve 6.002x2 - 47.91x + 6 =0

$$X = -b - \int b^2 - 4ac$$
 $X = -b + \int b^2 - 4ac$

Let b^2 is much large than $4ac$,

If byo, we can expect a loss of accuracy when computing -b+161

It b to, we expect to lose accuracy computing -b-lb|

Quadretic Formula, continued: a = 0.002 b = -47.91 c = 6Since 620, the computation of does not lose any accoracy. But the computation of -b- Jb2-40c (mi), 161-161 will lose several digits of accuracy. A better way to compute X:
The roots satisfy $X, X_2 = \frac{C}{a}$ Comparte X, Using the guadrotic formula; Compute X2 USing $X_2 = \frac{C}{a_{X_1}}$

Quadrotic Formula continued.

Better way to compute roots:

```
def quadratic_formula(a,b,c):
                                                        right way to compte
      if b < 0:
           x1 = (-b + sqrt(b**2 - 4*a*c))/(2*a)
           x2 = c/(a*x1)
                                                        roots
      else:
           x2 = (-b - sqrt(b**2 - 4*a*c))/(2*a)
           x1 = c/(a*x2)
      return x1, x2
a = 0.002
^{11} b = -47.91
12 c = 6
x1,x2 = quadratic_formula(a,b,c)
14 \times 2_bad = (-b - sqrt(b**2 - 4*a*c))/(2*a)
15 print('x1
                  = {:24.16f}'.format(x1))
= {:24.16f}'.format(x2))
16 print('x2
^{17} print('x2 (bad) = {:24.16f}'.format(x2_bad))
print('Difference = {:24.4e}'.format((abs(x2-x2_bad))))
                 23954.8747645299954456
x1
                     0.1252354700030452 good.
x2
x2 (bad)
                     0.1252354700032043
Difference
                              1.5918e-13
```

Difference between good and bad.

3 Floating point arithmetic w not associative, and not always commutative!

Associative property: (atb) +c = at (btc) operations carried out in

Commutative:

Try this 1 x,= .2+ (.3+.1) g are there the $X_{2} = (.2 + .3) + .1)$ samo?

Associativity - cont.

```
1 x1 = 0.2 + (0.3 + 0.1)
2 x2 = (0.2 + 0.3) + 0.1
3 if x1 == x2:
4    print("x1 == x2")
5 else:
6    print("What happened? x1 != x2")
```

What happened? x1 != x2

0.6000000000000008882 0.5999999999999997780

order in which numbers one added can matter!

Rule: We can only over expect two real numbers, when represented on the computer, to be approximately equal

the orall