

舍入误差

1. 什么是舍入误差

例如计算一个圆的面积，有公式： $A = \pi r^2$ 。其中 π 是无限不循环小数3.1415926535.....。假设半径 r 的值没有误差的条件下，当 π 取无限不循环小数的时候，计算得到的圆的面积才是准确的。但由于计算机内存有限，不可能存储无限位数据，所有数据的存储空间和位数都是有限且固定的，如在C语言中单精度float的有效位数为8位，双精度double的有效位数为16位，如此 π 只能取到3.1415926或3.141592653589793，这与 π 的准确值相比是有误差的，这样的误差叫做舍入误差，那么利用具有舍入误差的 π 计算得到的圆的面积 A 也是具有舍入误差的。

1.1 计算机精度验证

Python的默认浮点数是双精度，Python浮点数采用IEEE 754标准来表示。在Python中可以通过 `sys.float_info.epsilon` 方法来获取计算机浮点运算的精度。

```
In [ ]: import sys
        print(sys.float_info.epsilon)
```

2.220446049250313e-16

以上结果表明计算机浮点运算的精度在2.220446049250313e-16，实际上这个值是 2^{-52} ，64位双精度浮点型数的二进制位的第0-52位为2进制小数尾值，决定了它的精度在2的负52次方。意味着小数点后第16位之后是不可信的，同时如果当两个浮点数的数据相差 10^{16} 量级，则较小的数的贡献就会引入较为明显的舍入误差，甚至被忽略。

```
In [ ]: print(1e-10+8e-24)
        print(1e-10+8e-25)
        print(1e-10+8e-26)
        print(1e-10+8e-27) # 相差16个数量级，计算结果已经不准确了
        print(1e-10+1e-27) # 相差17个数量级，较小的数的贡献被忽略
```

```
1.000000000000008e-10
1.000000000000008e-10
1.000000000000008e-10
1.000000000000002e-10
1e-10
```

再举一个简单的计算的例子

```
In [ ]: a=4/3 # a是一个无限循环小数，a=1.333333...
        b=a-1 # b=0.33333333...
        c=1/3
        print(b) # 输出结果是0.3333333333333326，误差为4e-17
        print(b+b+b) # 输出结果是0.999999999999998，误差为2e-16
        print(c+c+c) # 输出结果是1.0，说明舍入误差发生在b=a-1浮点运算时
```

```
0.3333333333333326
0.9999999999999998
1.0
```

以上例子说明在计算机中有限位数的浮点运算是必然会产生舍入误差的。

2. 舍入误差能造成多大的影响

数值计算课程上老师告诉我们数值计算中应该注意的几个问题：

1. 避免相近的数相减
2. 避免大数吃小数
3. 简化运算步骤、减少运算次数

那么到底这三条会产生多大的影响，我们逐条举个例子。

2.1 避免相近的数相减——会严重损失有效位数

比如计算 $1.234 - 1.233$ ，原本各有4位有效数字，最后的理论计算结果为0.001，有效位数只剩下1位，而计算机实际输出的是0.000999999999763531，误差为 10^{-16} 。如此再进行后面的计算，如： $10.333 \times (1234.567 - 1234.566)$ ，计算的理论结果为0.010333，而计算机实际输出的是0.010332999999755657，此时误差为 10^{-15} 。如果在此基础上再乘10.333，理论结果为0.106770889，而计算机实际输出为0.10677088899747521，误差变为 10^{-14} 。即如果乘数 >10 ，每计算一次，都会损失一位有效位数，误差也会逐渐累积提升一个数量级。长此以往，后果不堪设想。

```
In [ ]: a=1.234
b=1.233
c=a-b
print(c)
d1=10.333*c
print(d1)
e1=10.333*d1
print(e1)
```

```
0.0009999999999998899
0.010332999999998862
0.10677088899998824
```

2.2 避免大数吃小数

举一个Google的首席科学家Vincent Vanhoucke曾举过的例子，考虑在 10^9 的基础上，加上 10^{-6} ，重复 10^6 次，再减去 10^9 ，即 $10^9 + 10^6 * 10^{-6} - 10^9$ ，理论值应该为1。可是计算机最后的计算结果为0.95367431640625，误差高达0.04632568359375。

```
In [ ]: summ = 1000000000

for indx in range(1000000):
    summ += 0.000001

summ -= 1000000000
```

```
print(summ)
```

0.95367431640625

以上例子告诉我们，计算过程越多，累计误差越大，所以简化运算步骤、减少运算次数便很容易理解了。

如今解决大数吃小数的问题有一个Kahan算法具有比较好的效果。Kahan算法的核心是记录了小数的负的补全部分compensation，随着这个补全部分的不断积累，当这些截断误差积累到一定量级，它们在求和的时候也就不会被截断了，从而能够相对好地控制整个求和过程的精度。

```
In [ ]: def KahanSum(input):
    var sum = 0.0
    var c = 0.0
    for i = 1 to input.length do
        var y = input[i] - c      # Initially, c is zero; then it compensates prev
        var t = sum + y           # Low-order digits of y are lost
        c = (t - sum) - y         # recover the low-order digits of y, with negati
        sum = t
    next i

    return sum
```

比如，用 $10000.0 + \pi + e$ 来说明，我们依旧假设浮点型变量只能保存6位数值。此时，具体写出求和算式应该是： $10000.0 + 3.14159 + 2.71828$ ，它们的理论结果应该是10005.85987，约等于10005.9。但由于截断误差，第一次求和 $10000.0 + 3.14159$ 只能得到结果10003.1；这个结果再与2.71828相加，得到10005.81828，被截断为10005.8。此时结果就相差了0.1。运用Kahan求和法，我们的运行过程是（记住，我们的浮点型变量保存6位数值）

```
y = 3.14159 - 0.00000

t = 10000.0 + 3.14159
  = 10003.14159
  = 10003.1                // low-order digits have lost

c = (10003.1 - 10000.0) - 3.14159
  = 3.10000 - 3.14159
  = - (.0415900)           // recover the negative parts
of compensation errors

sum = 1003.1
```

下面我们将Kahan算法应用到Vincent Vanhoucke的例子中，可以发现能够获得一个准确的结果。

```
In [ ]: summ = 1000000000
c = 0.0

for indx in range(1000000):
    y = 0.000001 - c
```

```

t = summ + y
c = (t - summ) - y
summ = t

summ -= 1000000000

print(summ)

```

1.0

2.3 再举一些我们平常经常会用到的场景的例子

考虑分别利用向前差分（一阶精度）、中心差分（二阶精度）和五点差分（三阶精度）计算函数 $f(x) = e^x$ 在 $x = 1$ 处的一阶导数值。

```

In [ ]: import matplotlib.pyplot as plt
import numpy as np

def ForwardDiff(fx, x, h=0.001):
    # Forward difference
    return (fx(x+h) - fx(x))/h

def CentralDiff(fx, x, h=0.001):
    # Central difference
    return (fx(x+h) - fx(x-h))/h*0.5

def FivePointsDiff(fx, x, h=0.001):
    # Five points difference
    return (-fx(x+2*h) + 8*fx(x+h) - 8*fx(x-h) + fx(x-2*h)) / (12.0*h)

# choose h from 0.1 to 10^-t, t>=2
t = 15
hx = 10**np.linspace(0,-t, 30)

# The exact derivative at x=1
x0 = 1
fprimExact = np.exp(1)

# Numerical derivative using the three methods
fprimF = ForwardDiff(np.exp, x0, hx)
fprimC = CentralDiff(np.exp, x0, hx)
fprim5 = FivePointsDiff(np.exp, x0, hx)

# Relative error
felF = abs(fprimExact - fprimF)/abs(fprimExact)
felC = abs(fprimExact - fprimC)/abs(fprimExact)
fel5 = abs(fprimExact - fprim5)/abs(fprimExact)

# Plot
fig, ax = plt.subplots(1)
ax.loglog(hx, felF)
ax.loglog(hx, felC)
ax.loglog(hx, fel5)
ax.autoscale(enable=True, axis='x', tight=True)
ax.set_xlabel(r'Step length $h$')
ax.set_ylabel('Relative error')
ax.legend(['Forward difference', 'Central difference', 'Five points difference'])

```

Out[]: <matplotlib.legend.Legend at 0x1eb2b41fa00>

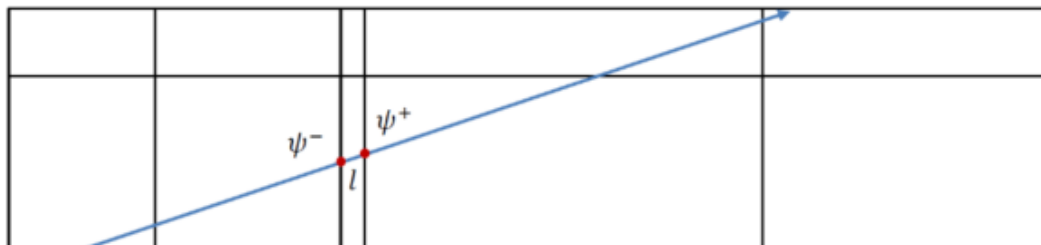
同时，我们也观察到，对于越是高阶的差分（如五点差分），它的离散化误差随 h 的下降速率越大，但也越早到达舍入误差的区域。因此，当我们遇到类似问题上，应选择合适阶

数的有限差分方法，并根据它的特性选择适合的 h 值。并不一定是越高阶的方法越好， h 越小越好。

MOC精细结构计算

在特征线方法中，核心的公式就是在每一段segment上利用入射角通量的指数衰减计算segment上的通量变化量，然后再将所有segment的通量变化量求积，即：

$$\delta\phi = \varphi^- - \varphi^+ = (\varphi^- - Q/\Sigma_T)(1 - e^{-l\Sigma_T})$$



考虑图中这一精细结构的计算，假设线段长度 l 为1cm，则光学距离约为0.106，那么 $(1 - e^{-l\Sigma_T}) \approx 0.1$ ，如果 $\delta\phi$ 恰好也在0.1附近，则 $(\varphi^- - Q/\Sigma_T)$ 便会遇到相近的数相减的问题，从而损失有效位数引入误差，而这种误差会在标通量和源的计算中不断累积，且网格细化并不能解决这个问题。

另一方面，如果结构太过精细，也可能会踩到大数吃小数的坑。

这或许也是MOC谜一样的敏感性分析的一项可能原因。