

1.3 Real numbers in Python

Decimal module

The module decimal implements the arithmetic we usually do by hand. All the standard arithmetic can be used and, if necessary, the results will be rounded to a number of digits which can be changed by setting the context parameter prec. However, the available functions are limited and the operations may take longer than with the built in data types.

Note that the numbers have to be input as strings (as otherwise Python would automatically round to floating point numbers).

```
import decimal as dec
```

```
x = dec.Decimal('0.6')
```

```
y = dec.Decimal('0.5999999999') # 10 significant digits
```

```
z = x - y
```

~~x = 0.6~~

```
print("x = {0}, y = {1}, x-y = {2}\n".format(x,y,z))  
print("representation of x: {!r}\n".format(x))  
print("type of x: {}\n".format(type(x)))
```

```
print("timing")  
%timeit(x-y)
```

```
x = 0.6, y = 0.5999999999, x-y = 1E-10
```

```
representation of x: Decimal('0.6')
```

```
type of x: <class 'decimal.Decimal'>
```

```
timing
```

```
101 ns ± 0.729 ns per loop (mean ± std. dev. of 7 runs, 100
```

Python's default real numbers

Python uses 64 bit floating point numbers with 53 binary significant digits per default for real numbers. The arithmetic with these numbers is implemented in hardware and is very fast. However, even simple numbers like 0.6 cannot be exactly represented as floating point numbers of this type and thus one gets rounding errors. The effect is even worse when one subtracts two numbers which are close and loses a significant amount of digits. This is called *cancellation*.

The usage of floating point numbers is illustrated below:

```
## Default floating point numbers
```

```
x = 0.6
```

```
y = 0.5999999999 # 10 significant decimal digits
```

```
z = x - y
```

```
print("x = {0}, y = {1}, x-y = {2}\n".format(x,y,z))
```

```
print("representation of x: {!r}\n".format(x))
```

```
print("type of x: {}\n".format(type(x)))
```

```
x = 0.6, y = 0.5999999999, x-y = 1.000000082740371e-10
```

```
representation of x: 0.6
```

```
type of x: <class 'float'>
```

```
print("printing the result with some more digits:\n")
print("    x = {0:3.100g}\n    y = {1:3.100g}\n    x-y = {2:3.100g}\n")
```

```
xstring = "{:3.100g}".format(x)
xmystring = "{:3.100g}".format(x-y)
print("significant decimal digits:    x: {0},    x-y: {1}\n")
```

```
print("timing")
%timeit(x-y)
```

printing the result with some more digits:

```
x = 0.59999999999999997779553950749686919152736663818359
y = 0.599999999899999996952150240758783183991909027099609
x-y = 1.000000082740370999090373516082763671875e-10
```

significant decimal digits: x: 53, x-y: 39

timing

40.7 ns \pm 0.389 ns per loop (mean \pm std. dev. of 7 runs, 1000000 loops each)

Real numbers in numpy

The numerical package numpy has three major floating point number systems: float16, float32 and float64. They are useful for the development of resource critical applications using lower accuracy arithmetic which is supported by some hardware including some graphics boards. In particular using these types is in principle a first way to save energy as the costliest operations are data transfers and using the types below on can thus cut costs by a factor of up to four.

The significant (binary) digits of the three numpy types are

- ▶ float16 has 11
- ▶ float32 has 24
- ▶ float64 has 53

Note that only 10, 23 and 52 of these bits are stored as it is assumed that the first significant bit is always 1. The accuracy and timing is illustrated below. On my laptop the timings were the same, actually, the 64 bit version was the fastest! The reason for this might be that the actual arithmetic is still done in the hardware floating point unit (using 64 or even higher accuracy) and the result is then rounded. As a direct conversion from float16 and float32 to decimal is not supported we first convert these two types to float64 (which can be done without error). Also note that the difference is less than the error for float16 and one gets then 100 percent error.

NB As we are computing with lower accuracy, we have changed the problem a bit compared to above!

numpy floating point numbers -- print out using the dec

```
import decimal as dec  
import numpy as np
```

three numpy types

```
x16 = np.float16(0.6)  
y16 = np.float16(0.59999)  
z16 = x16 - y16
```

```
x32 = np.float32(0.6)  
y32 = np.float32(0.59999)  
z32 = x32 - y32
```

```
x64 = np.float64(0.6)  
y64 = np.float64(0.59999)  
z64 = x64 - y64
```

```
# exact values
```

```
xex = dec.Decimal("0.6")
```

```
yex = dec.Decimal("0.59999")
```

```
zex = xex - yex
```

```
print("    x16 = {0}\n    x32 = {1}\n    x64 = {2}\n".format(xex, yex, zex))
```

```
    x16 = 0.60009765625
```

```
    x32 = 0.60000002384185791015625
```

```
    x64 = 0.599999999999999977795539507496869191527366638183
```

```

print("relative rounding errors:\n")
e16 = (dec.Decimal(float(x16)) - xex)/xex
e32 = (dec.Decimal(float(x32)) - xex)/xex
e64 = (dec.Decimal(x64) - xex)/xex
print("  for x:      float16 :   {:3.2g},    float32 :   {:3.2g},    float64 :   {:3.2g},    float128 :   {:3.2g},    float256 :   {:3.2g},    float512 :   {:3.2g},    float1024 :   {:3.2g},    float2048 :   {:3.2g},    float4096 :   {:3.2g},    float8192 :   {:3.2g},    float16384 :   {:3.2g},    float32768 :   {:3.2g},    float65536 :   {:3.2g},    float131072 :   {:3.2g},    float262144 :   {:3.2g},    float524288 :   {:3.2g},    float1048576 :   {:3.2g},    float2097152 :   {:3.2g},    float4194304 :   {:3.2g},    float8388608 :   {:3.2g},    float16777216 :   {:3.2g},    float33554432 :   {:3.2g},    float67108864 :   {:3.2g},    float134217728 :   {:3.2g},    float268435456 :   {:3.2g},    float536870912 :   {:3.2g},    float1073741824 :   {:3.2g},    float2147483648 :   {:3.2g},    float4294967296 :   {:3.2g},    float8589934592 :   {:3.2g},    float17179869184 :   {:3.2g},    float34359738368 :   {:3.2g},    float68719476736 :   {:3.2g},    float137438953472 :   {:3.2g},    float274877906944 :   {:3.2g},    float549755813888 :   {:3.2g},    float1099511627776 :   {:3.2g},    float2199023255552 :   {:3.2g},    float4398046511104 :   {:3.2g},    float8796093022208 :   {:3.2g},    float17592186044416 :   {:3.2g},    float35184372088832 :   {:3.2g},    float70368744177664 :   {:3.2g},    float140737488355328 :   {:3.2g},    float281474976710656 :   {:3.2g},    float562949953421312 :   {:3.2g},    float1125899906842624 :   {:3.2g},    float2251799813685248 :   {:3.2g},    float4503599627370496 :   {:3.2g},    float9007199254740992 :   {:3.2g},    float18014398509481984 :   {:3.2g},    float36028797018963968 :   {:3.2g},    float72057594037927936 :   {:3.2g},    float144115188075855872 :   {:3.2g},    float288230376151711744 :   {:3.2g},    float576460752303423488 :   {:3.2g},    float1152921504606846976 :   {:3.2g},    float2305843009213693952 :   {:3.2g},    float4611686018427387904 :   {:3.2g},    float9223372036854775808 :   {:3.2g},    float18446744073709551616 :   {:3.2g},    float36893488147419103232 :   {:3.2g},    float73786976294838206464 :   {:3.2g},    float147573952589676412928 :   {:3.2g},    float295147905179352825856 :   {:3.2g},    float590295810358705651712 :   {:3.2g},    float1180591620717411303424 :   {:3.2g},    float2361183241434822606848 :   {:3.2g},    float4722366482869645213696 :   {:3.2g},    float9444732965739290427392 :   {:3.2g},    float18889465931478580854784 :   {:3.2g},    float37778931862957161709568 :   {:3.2g},    float75557863725914323419136 :   {:3.2g},    float151115727451828646838272 :   {:3.2g},    float302231454903657293676544 :   {:3.2g},    float604462909807314587353088 :   {:3.2g},    float1208925819614629174706176 :   {:3.2g},    float2417851639229258349412352 :   {:3.2g},    float4835703278458516698824704 :   {:3.2g},    float9671406556917033397649408 :   {:3.2g},    float19342813113834066795298816 :   {:3.2g},    float38685626227668133590597632 :   {:3.2g},    float77371252455336267181195264 :   {:3.2g},    float154742504910672534362390528 :   {:3.2g},    float309485009821345068724781056 :   {:3.2g},    float618970019642690137449562112 :   {:3.2g},    float1237940039285380274899124224 :   {:3.2g},    float2475880078570760549798248448 :   {:3.2g},    float4951760157141521099596496896 :   {:3.2g},    float9903520314283042199192993792 :   {:3.2g},    float19807040628566084398385987584 :   {:3.2g},    float39614081257132168796771975168 :   {:3.2g},    float79228162514264337593543950336 :   {:3.2g},    float158456325028528675187087900672 :   {:3.2g},    float316912650057057350374175801344 :   {:3.2g},    float633825300114114700748351602688 :   {:3.2g},    float1267650600228229401496703205376 :   {:3.2g},    float2535301200456458802993406410752 :   {:3.2g},    float5070602400912917605986812821504 :   {:3.2g},    float10141204801825835211973625643008 :   {:3.2g},    float20282409603651670423947251286016 :   {:3.2g},    float40564819207303340847894502572032 :   {:3.2g},    float81129638414606681695789005144064 :   {:3.2g},    float162259276829213363391578010288128 :   {:3.2g},    float324518553658426726783156020576256 :   {:3.2g},    float649037107316853453566312041152512 :   {:3.2g},    float1298074214633706907132624082305024 :   {:3.2g},    float2596148429267413814265248164610048 :   {:3.2g},    float5192296858534827628530496329220096 :   {:3.2g},    float10384593717069655257060992658440192 :   {:3.2g},    float20769187434139310514121985316880384 :   {:3.2g},    float41538374868278621028243970633760768 :   {:3.2g},    float83076749736557242056487941267521536 :   {:3.2g},    float166153499473114484112975882535043072 :   {:3.2g},    float332306998946228968225951765070086144 :   {:3.2g},    float664613997892457936451903530140172288 :   {:3.2g},    float1329227995784915872903807060280344576 :   {:3.2g},    float2658455991569831745807614120560689152 :   {:3.2g},    float5316911983139663491615228241121378304 :   {:3.2g},    float10633823966279326983230456482242756608 :   {:3.2g},    float21267647932558653966460912964485513216 :   {:3.2g},    float42535295865117307932921825928971026432 :   {:3.2g},    float85070591730234615865843651857942052864 :   {:3.2g},    float170141183460469231731687303715884105728 :   {:3.2g},    float340282366920938463463374607431768211456 :   {:3.2g},    float680564733841876926926749214863536422912 :   {:3.2g},    float1361129467683753853853498429727072845824 :   {:3.2g},    float2722258935367507707706996859454145691648 :   {:3.2g},    float5444517870735015415413993718908291383296 :   {:3.2g},    float10889035741470030830827987437816582766592 :   {:3.2g},    float21778071482940061661655974875633165533184 :   {:3.2g},    float43556142965880123323311949751266331066368 :   {:3.2g},    float87112285931760246646623899502532662132736 :   {:3.2g},    float174224571863520493293247799005065324265472 :   {:3.2g},    float348449143727040986586495598010130648530944 :   {:3.2g},    float696898287454081973172991196020261297061888 :   {:3.2g},    float1393796574908163946345982392040522594123776 :   {:3.2g},    float2787593149816327892691964784081045188247552 :   {:3.2g},    float5575186299632655785383929568162090376495104 :   {:3.2g},    float11150372599265311570767859136324180752990208 :   {:3.2g},    float22300745198530623141535718272648361505980416 :   {:3.2g},    float44601490397061246283071436545296723011960832 :   {:3.2g},    float89202980794122492566142873090593446023921664 :   {:3.2g},    float178405961588244985132285746181186892047843328 :   {:3.2g},    float356811923176489970264571492362373784095686656 :   {:3.2g},    float713623846352979940529142984724747568191373312 :   {:3.2g},    float1427247692705959881058285969449495136382746624 :   {:3.2g},    float2854495385411919762116571938898990272765493248 :   {:3.2g},    float5708990770823839524233143877797980545530986496 :   {:3.2g},    float11417981541647679048466287755595961091061972992 :   {:3.2g},    float22835963083295358096932575511191922182123945984 :   {:3.2g},    float45671926166590716193865151022383844364247891968 :   {:3.2g},    float91343852333181432387730302044767688728495783936 :   {:3.2g},    float182687704666362864775460604089535377456991567872 :   {:3.2g},    float365375409332725729550921208179070754913983135744 :   {:3.2g},    float730750818665451459101842416358141509827966271488 :   {:3.2g},    float1461501637330902918203684832716283019655932542976 :   {:3.2g},    float2923003274661805836407369665432566039311865085952 :   {:3.2g},    float5846006549323611672814739330865132078623730171904 :   {:3.2g},    float11692013098647223345629478661730264157247460343808 :   {:3.2g},    float23384026197294446691258957323460528314494920687616 :   {:3.2g},    float46768052394588893382517914646921056628989841375232 :   {:3.2g},    float93536104789177786765035829293842113257979682750464 :   {:3.2g},    float187072209578355573530071658587684226515959365500928 :   {:3.2g},    float374144419156711147060143317175368453031918731001856 :   {:3.2g},    float748288838313422294120286634350736906063837462003712 :   {:3.2g},    float1496577676626844588240573268701473812127674924007424 :   {:3.2g},    float2993155353253689176481146537402947624255349848014848 :   {:3.2g},    float5986310706507378352962293074805895248510699696029696 :   {:3.2g},    float11972621413014756705924586149611790497021399392059392 :   {:3.2g},    float23945242826029513411849172299223580994042798784118784 :   {:3.2g},    float47890485652059026823698344598447161988085597568237568 :   {:3.2g},    float95780971304118053647396689196894323976171195136475136 :   {:3.2g},    float191561942608236107294793378393788647952342390272950272 :   {:3.2g},    float383123885216472214589586756787577295904684780545900544 :   {:3.2g},    float766247770432944429179173513575154591809369561091801088 :   {:3.2g},    float1532495540865888858358347027150309183618739122183602176 :   {:3.2g},    float3064991081731777716716694054300618367237478244367204352 :   {:3.2g},    float6129982163463555433433388108601236734474956488734408704 :   {:3.2g},    float12259964326927110866866776217202473468949912977468817408 :   {:3.2g},    float24519928653854221733733552434404946937899825954937634816 :   {:3.2g},    float49039857307708443467467104868809893875799651909875269632 :   {:3.2g},    float98079714615416886934934209737619787751599303819750539264 :   {:3.2g},    float196159429230833773869868419475239575503198607639501078528 :   {:3.2g},    float392318858461667547739736838950479151006397215279002157056 :   {:3.2g},    float784637716923335095479473677900958302012794430558004314112 :   {:3.2g},    float1569275433846670190958947355801916604025588861116008628224 :   {:3.2g},    float3138550867693340381917894711603833208051177722232017256448 :   {:3.2g},    float6277101735386680763835789423207666416102355444464034512896 :   {:3.2g},    float12554203470773361527671578846415332832204710888928069025792 :   {:3.2g},    float25108406941546723055343157692830665664409421777856138051584 :   {:3.2g},    float50216813883093446110686315385661331328818843555712276103168 :   {:3.2g},    float100433627766186892221372630771322662657637687111424552206336 :   {:3.2g},    float200867255532373784442745261542645325315275374222849104412672 :   {:3.2g},    float401734511064747568885490523085290650630550748445698208825344 :   {:3.2g},    float803469022129495137770981046170581301261101496891396417650688 :   {:3.2g},    float1606938044258990275541962092341162602522202993782792835301376 :   {:3.2g},    float3213876088517980551083924184682325205044405987565585670602752 :   {:3.2g},    float6427752177035961102167848369364650410088811975131171341205504 :   {:3.2g},    float12855504354071922204335696738729300820177623950262342682411008 :   {:3.2g},    float25711008708143844408671393477458601640355247900524685364822016 :   {:3.2g},    float51422017416287688817342786954917203280710495801049370729644032 :   {:3.2g},    float102844034832575377634685573909834406561420991602098741459288064 :   {:3.2g},    float205688069665150755269371147819668813122841983204197482918576128 :   {:3.2g},    float411376139330301510538742295639337626245683966408394965837152256 :   {:3.2g},    float822752278660603021077484591278675252491367932816789931674304512 :   {:3.2g},    float1645504557321206042154969182557350504982735865633579863348609024 :   {:3.2g},    float3291009114642412084309938365114701009965471731267159726697218048 :   {:3.2g},    float6582018229284824168619876730229402019930943462534319453394436096 :   {:3.2g},    float13164036458569648337239753460458804039861886925068638906788872192 :   {:3.2g},    float26328072917139296674479506920917608079723773850137277813577744384 :   {:3.2g},    float52656145834278593348959013841835216159447547700274555627155488768 :   {:3.2g},    float105312291668557186697918027683670432318895095400549111254310977536 :   {:3.2g},    float210624583337114373395836055367340864637790190801098222508621955072 :   {:3.2g},    float421249166674228746791672110734681729275580381602196445017243910144 :   {:3.2g},    float842498333348457493583344221469363458551160763204392890034487820288 :   {:3.2g},    float1684996666696914987166688442938726917102321526408785780068975640576 :   {:3.2g},    float3369993333393829974333376885877453834204643052817571560137951281152 :   {:3.2g},    float6739986666787659948666753771754907668409286105635143120275902562304 :   {:3.2g},    float13479973333575319897333507543509815336818572211270286240551805124608 :   {:3.2g},    float26959946667150639794667015087019630673637144422540572481103610249216 :   {:3.2g},    float53919893334301279589334030174039261347274288845081144962207220498432 :   {:3.2g},    float107839786668602559178668060348078522694548577690162289924414440996864 :   {:3.2g},    float215679573337205118357336120696157045389097155380324579848828881993728 :   {:3.2g},    float431359146674410236714672241392314090778194310760649159697657763987456 :   {:3.2g},    float862718293348820473429344482784628181556388621521298319395315527974912 :   {:3.2g},    float1725436586697640946858688965569256363112777243042596638790631055949824 :   {:3.2g},    float3450873173395281893717377931138512726225554486085193277581262111899648 :   {:3.2g},    float6901746346790563787434755862277025452451108972170386555162524223799296 :   {:3.2g},    float13803492693581127574869511724554050904902217944340773110325048447598592 :   {:3.2g},    float27606985387162255149739023449108101809804435888681546220650096895197184 :   {:3.2g},    float55213970774324510299478046898216203619608871777363092441300193790394368 :   {:3.2g},    float110427941548649020598956093796432407239217743554726184882600387580788736 :   {:3.2g},    float220855883097298041197912187592864814478435487109452369765200775161577472 :   {:3.2g},    float441711766194596082395824375185729628956870974218904739530401550323154944 :   {:3.2g},    float883423532389192164791648750371459257913741948437809479060803100646309888 :   {:3.2g},    float1766847064778384329583297500742918515827483896875618958121606201292619776 :   {:3.2g},    float3533694129556768659166595001485837031654967793751237916243212402585239552 :   {:3.2g},    float7067388259113537318333190002971674063309935587502475832486424805170479104 :   {:3.2g},    float14134776518227074636666380005943348126619871175004951664972849610340958208 :   {:3.2g},    float28269553036454149273332760011886696253239742350009903329945699220681916416 :   {:3.2g},    float56539106072908298546665520023773392506479484700019806659891398441363832832 :   {:3.2g},    float113078212145816597093331040047546785012958969400039613319782796882727665664 :   {:3.2g},    float226156424291633194186662080095093570025917938800079226639565593765455331328 :   {:3.2g},    float452312848583266388373324160190187140051835877600158453279131187530910662656 :   {:3.2g},    float904625697166532776746648320380374280103671755200316906558262375061821325312 :   {:3.2g},    float1809251394333065553493296640760748560207343510400633813116524750123642650624 :   {:3.2g},    float3618502788666131106986593281521497120414687020801267626233049500247285301248 :   {:3.2g},    float7237005577332262213973186563042994240829374041602535252466099000494570602496 :   {:3.2g},    float14474011154664524427946373126085988481658748083205070504932198000989141
```

```
print("timing - not much difference at this level")
print("* float16")
%timeit(x16-y16)
print("* float32")
%timeit(x32-y32)
print("* float64")
%timeit(x64-y64)
```

timing - not much difference at this level

* float16

87.7 ns \pm 0.776 ns per loop (mean \pm std. dev. of 7 runs, 10

* float32

75.5 ns \pm 0.608 ns per loop (mean \pm std. dev. of 7 runs, 10

* float64

76.7 ns \pm 0.836 ns per loop (mean \pm std. dev. of 7 runs, 10

npmath multiple precision

This module does provide floating point with choosable (binary) accuracy. We will choose 113 binary digits which is the standard for quadruple (128 bit) arithmetic.

In order to compute the exact numbers we first convert the npmath number to a string (note that we need to use 113 bit decimal precision to get the exact result). Then we convert this string to a Decimal. Note that multiple precision operations are substantially slower than the floating point ones.

Looking at the printout it seems that only the first 30 or so decimal digits are accurate and the later ones are wrong. Thus without losing much accuracy, one could set these later digits to zero.

However, we should remember, that the numerical approximation is binary number with 113 digits and the number is accurate to all the binary digits. If one now changes any of the later decimal digits and then rounds again to the nearest binary number one typically gets a larger error.

```
import mpmath as mpm
```

```
prec = 113
```

```
mpm.mp.prec = prec # set precision to quadruple
```

```
## Default floating point numbers -- print out using the d
```

```
x = mpm.mpf('0.6')
```

```
y = mpm.mpf('0.5999999999') # 10 significant decimal digits
```

```
z = x - y
```

```
xex = dec.Decimal("0.6") # exact values ..
```

```
yex = dec.Decimal("0.5999999999")
```

```
zex = xex - yex
```

```
xdec = dec.Decimal(mpm.nstr(x,prec)) # first convert to string
```

```
ydec = dec.Decimal(mpm.nstr(y,prec))
```

```
zdec = dec.Decimal(mpm.nstr(z,prec))
```



```
print("timing")  
%timeit(x-y)
```

timing

1.66 μ s \pm 9.08 ns per loop (mean \pm std. dev. of 7 runs, 100

Other options

One can use other C data types and can also get access to the 80 or 128 bit accuracy of the hardware processor. Especially for running on GPUs one may also use the data types these processors use natively. For our purposes the above methods will be sufficient.