**Problem 1**

$$cond(f) = \frac{|relative\ change\ in\ solution|}{|relative\ change\ in\ input\ data|} = \frac{[f(x+\Delta x, y+\Delta y) - f(x,y)]/f(x,y)}{[(x+\Delta x, y+\Delta y) - (x,y)]/(x,y)}$$

$$= \frac{|(x+\Delta x) - (y+\Delta y) - (x-y)|/(x-y)}{(\Delta x, \Delta y)/(x,y)} = \frac{(\Delta x - \Delta y)(|x|+|y|)}{(x-y)(|\Delta x|+|\Delta y|)}$$

$$\geq \frac{\Delta x - \Delta y}{(|\Delta x|+|\Delta y|)\varepsilon} \geq \frac{1}{\varepsilon}$$

we can conclude that subtraction is sensitive when $\varepsilon$ is small

**Problem 2**

Suppose a and b have the same sign. (ii) would be preferable since the result is guaranteed to lie in the interval [a; b], and no overflow would happen. The rounding process will do no harm to the problem, too.

For example,

a) $\beta = 10,\ t = 2,\ [L, U] = [-2, 2]$

b) $a = 5.7 \times 10^{-1},\ b = 5.9 \times 10^{-1}$

c) The intermediate results for i and ii are as follows

    i)     $(a+b) \approx 1.2 \times 10^{-1}$         (i1)

        $1.2 \times 10^{-1}/2.0 = 0.6 \times 10^{-1}$     (i2)

    ii)     $(b-a) = 0.2 \times 10^{-1}$         (ii1)

        $(b-a)/2 = 0.1 \times 10^{-1}$     (ii2)

        $a + 0.1 \times 10^{-1} = 5.8 \times 10^{-1}$     (ii3)

d) At step (i1) , according to the marks above, the problem will occur due to <u>rounding</u>.

**Problem 3: Bessel recurrence vs. floating point**

a) Please check problem3_a.py
   For z = 20, we have the following result
   n=2,   value=-0.160341       error=0
   n=3,   value=-0.0989014      error=1.82415e-15
   n=4,   value=0.130671        error=2.12408e-16
   n=5,   value=0.15117         error=3.67211e-16
   n=6,   value=-0.055086       error=7.55788e-16
   n=7,   value=-0.184221       error=3.01328e-16
   n=8,   value=-0.0738689      error=5.63611e-16
   n=9,   value=0.125126        error=2.21821e-16
   n=10,  value=0.186483        error=1.48837e-16
   n=11,  value=0.0613563       error=4.52367e-16
   n=12,  value=-0.118991       error=2.33259e-16
   n=13,  value=-0.204145       error=0
   n=14,  value=-0.146398       error=7.5836e-16
   n=15,  value=-0.000812069 error=1.10948e-13
   n=16,  value=0.14518         error=1.91181e-16
   n=17,  value=0.2331          error=1.19072e-16
   n=18,  value=0.25109         error=6.63242e-16
   n=19,  value=0.218862        error=1.26818e-15
   n=20,  value=0.164748        error=5.05419e-16
   n=21,  value=0.110634        error=2.2579e-15
   n=22,  value=0.0675829       error=3.69621e-15
   n=23,  value=0.0380487       error=1.82369e-16
   n=24,  value=0.0199291       error=3.13361e-15
   n=25,  value=0.00978117      error=1.95089e-15
   n=26,  value=0.00452381      error=1.91733e-15
   n=27,  value=0.00198074      error=6.56848e-16
   n=28,  value=0.000824178   error=1.31549e-16
   n=29,  value=0.000326963   error=1.65799e-15
   n=30,  value=0.000124015   error=2.40418e-15
   n=31,  value=4.50828e-05   error=1.05215e-15
   n=32,  value=1.57413e-05   error=5.81145e-15
   n=33,  value=5.28924e-06   error=6.50179e-14
   n=34,  value=1.71324e-06   error=3.26306e-14
   n=35,  value=5.35784e-07   error=4.86134e-14
   n=36,  value=1.62001e-07   error=4.06847e-14
   n=37,  value=4.74202e-08   error=2.76307e-14
   n=38,  value=1.34536e-08   error=4.67277e-14
   n=39,  value=3.70356e-09   error=4.55629e-14
   n=40,  value=9.90239e-10   error=5.53409e-14
   n=41,  value=2.57401e-10   error=3.9969e-14
   n=42,  value=6.51039e-11   error=1.07203e-14
   n=43,  value=1.60356e-11   error=5.21885e-14
   n=44,  value=3.84926e-12   error=2.64419e-14

n=45,  value=9.01145e-13  error=2.00571e-14
n=46,  value=2.05887e-13  error=3.53113e-14
n=47,  value=4.59366e-14  error=3.95662e-14
n=48,  value=1.00149e-14  error=4.25353e-14
n=49,  value=2.13469e-15  error=4.85951e-14
n=50,  value=4.45104e-16  error=2.37046e-14

b) Please check problem3_b.py
For z = 20, we have the following result
n=2,   value=-0.160341      error=0
n=3,   value=-0.0989014     error=1.82415e-15
n=4,   value=0.130671       error=2.12408e-16
n=5,   value=0.15117        error=7.34421e-16
n=6,   value=-0.055086      error=2.26736e-15
n=7,   value=-0.184221      error=4.51993e-16
n=8,   value=-0.0738689     error=1.8787e-15
n=9,   value=0.125126       error=4.43641e-16
n=10,  value=0.186483       error=5.9535e-16
n=11,  value=0.0613563      error=2.26184e-15
n=12,  value=-0.118991      error=5.83146e-16
n=13,  value=-0.204145      error=2.7192e-16
n=14,  value=-0.146398      error=1.51672e-15
n=15,  value=-0.000812069   error=1.96662e-13
n=16,  value=0.14518        error=3.82361e-16
n=17,  value=0.2331         error=2.38143e-16
n=18,  value=0.25109        error=2.21081e-16
n=19,  value=0.218862       error=5.07271e-16
n=20,  value=0.164748       error=1.17931e-15
n=21,  value=0.110634       error=2.50878e-16
n=22,  value=0.0675829      error=1.64276e-15
n=23,  value=0.0380487      error=5.65343e-15
n=24,  value=0.0199291      error=1.63644e-14
n=25,  value=0.00978117     error=6.01228e-14
n=26,  value=0.00452381     error=2.5117e-13
n=27,  value=0.00198074     error=1.19524e-12
n=28,  value=0.000824178    error=6.37739e-12
n=29,  value=0.000326963    error=3.77686e-11
n=30,  value=0.000124015    error=2.4639e-10
n=31,  value=4.50828e-05    error=1.75942e-09
n=32,  value=1.57413e-05    error=1.36796e-08
n=33,  value=5.28924e-06    error=1.15281e-07
n=34,  value=1.71324e-06    error=1.0488e-06
n=35,  value=5.3579e-07     error=1.02644e-05

```
n=36,  value=1.62019e-07   error=0.000107724
n=37,  value=4.74775e-08   error=0.00120889
n=38,  value=1.36483e-08   error=0.0144685
n=39,  value=4.38591e-09   error=0.184244
n=40,  value=3.45678e-09   error=2.49086
n=41,  value=9.44122e-09   error=35.6791
n=42,  value=3.52522e-08   error=540.476
n=43,  value=1.38618e-07   error=8643.39
n=44,  value=5.60806e-07   error=145691
n=45,  value=2.32893e-06   error=2.58441e+06
n=46,  value=9.91936e-06   error=4.81786e+07
n=47,  value=4.33001e-05   error=9.42606e+08
n=48,  value=0.000193591   error=1.93304e+10
n=49,  value=0.000885938   error=4.1502e+11
n=50,  value=0.00414751    error=9.31806e+12
```

c) the truncation error in (1) due to loss of precision caused by cancellation cannot be bounded anymore when n reaches round 30

d) Yes. When we calculate the result from 50, the number have less precision bits than it needs in the computer. Then we used the number with bits lost to do the calculation. The error will be accumulated during computing and then the precision would definitely lost. (Please check problem3_d.py)

**Problem 4: Gaussian elimination and partial pivoting**

a) Please check problem4_a.py
b) Please check problem4_b.py
c) Please check problem4_c{1,2,3}.py
   1) a 'random' matrix
      - condition number for matrix A :     2032.7
      - residual from un-pivoted solve:     9.57203e-11
      - error from un-pivoted solve:        4.5902e-11
      - residual from partially-pivoted solve:6.66282e-11
      - error from partially-pivoted solve:  2.56879e-10
      - residual from np.linalg.solve:       1.48289e-13
      - error from np.linalg.solve:          8.26606e-13

   This is not a well conditioned matrix, Gaussian elimination with partial pivoting is more accurate than Gaussian elimination without pivoting.

   2) the matrix given by
      - condition number for matrix A :     1.02
      - un-pivoted solves failed
      - residual from  partially-pivoted solve3.89423e-14
      - error from  partially-pivoted solve:  7.90039e-15
      - residual from np.linalg.solve:        3.92022e-14
      - error from np.linalg.solve:           7.95229e-15

   This is a well conditioned matrix. From the result, we can find that unpivoted case would be possible to fail solving the problem. However, partially-pivoted Gaussian elimination could solve the problem somewhat well.

   3) the matrix given by
      - condition number for matrix A :     1.30228
      - residual from un-pivoted solve:     6.38225e-08
      - error from un-pivoted solve:        6.42032e-08
      - residual from  partially-pivoted solve4.67187e-15
      - error from  partially-pivoted solve:  4.53856e-15
      - residual from np.linalg.solve:4.51273e-15
      - error from np.linalg.solve:    4.35569e-15

   This is a well conditioned matrix. Both Gaussian unpivoted and partially-pivoted Gaussian elimination could solve the problem. For this matrix, Gaussian elimination with partial pivoting is more accurate than Gaussian elimination without pivoting.

**Problem 5: Scaling a linear system**

    a) Please check problem5_a.py

        relative residuals: 2.01239e-15

        relative error: 9.28687e-14

        cond(A): 2032.7

    b) Please check problem5_b.py

        relative residuals: 1

        relative error: 9.28687e-14

        cond(DA): 2032.7

    c) Please check problem5_c.py

        relative residuals: 49.2039

        relative error: 1.54848e-13

        cond(DA): 11523.4

    d) Please check problem5_d.py

        relative residuals: 4967.92

        relative error: 8.31082e-14

        cond(DA): 633521

    e) Please check problem5_e.py

        relative residuals: 3.79728e+14

        relative error: 9.83454e-14

        cond(DA): 2.56716e+31

The scaling of case in (c) gives the worst accuracy.

In this case, relative residuals is still small comparing to case (d) and (e) .

Condition number is positively correlated with relative residual but not necessarily correlated with relative error.