

# Tutorial 2 solutions

1. From lecture, we know that a number,  $x$ , can be represented as follows in a floating point system:

$$x = \pm \sum_{k=0}^{p-1} \left( \frac{d_k}{\beta^k} \right) \beta^E$$

- Consider a floating point system with  $\beta=10$ ,  $p=4$ , and  $U=3$

The largest number we can represent is  $9.999 \times 10^3$   
(we want to make each digit in the mantissa as large as possible, and make the exponent as large as possible too)

- Therefore, for a general floating point system, the largest number,  $x_*$ , must satisfy the following:

- 1)  $d_k = \beta - 1 \quad \forall k$  (so each digit in the mantissa is as large as possible)
- 2)  $\text{sign} = +$
- 3)  $E = U$  (i.e. exponent is as large as possible)

Substituting these properties into our expression for  $x$ , we get:

$$x_* = \sum_{k=0}^{p-1} \left( \frac{\beta-1}{\beta^k} \right) \beta^U$$

- Now we can use the following equality:  $\sum_{n=0}^N \frac{a-1}{a^n} = a - a^{-N}$ , where  $N=p-1$ ,  $n=k$ ,  $a=\beta$

$$\therefore x_* = (\beta - \beta^{-(p-1)}) \beta^U$$

$$= (\beta - \beta^{-p+1}) \beta^U$$

$$= (\beta - \beta^{-p} \beta) \beta^U$$

$$= \beta(1 - \beta^{-p}) \beta^U$$

$\therefore$  the largest number we can represent is  $x_* = (1 - \beta^{-p}) \beta^{U+1}$

2. Recall that an  $n$ -bit binary variable can take on  $2^n$  unique values.

Now, if our floating point system can represent  $m$  unique values for the mantissa, and  $e$  unique values for the exponent, the total number of unique floating point numbers it can represent is:

$$m \cdot e$$

(ignoring the sign)

If we use 8 bits to store the exponent, we can have  $2^8$  unique values for the exponent. Using 23 bits for the mantissa means it can take on  $2^{23}$  unique values.

$\therefore$  we can represent  $2^8 \cdot 2^{23} = 2^{31}$  unique floating point numbers

If we use 9 bits for the exponent and 22 for the mantissa, we can represent  $2^9 \cdot 2^{22} = 2^{31}$  unique floating point numbers.

This means, both systems can represent the same number of floating point numbers. However, the question asks about normalized floats.

Recall that if the exponent is all zeroes, the number is de-normalized. Therefore, we should exclude this exponent from our calculations

$\therefore$  8-bit exponent & 23-bit mantissa can represent  $(2^8 - 1)(2^{23}) = 2^{31} - 2^{23}$  unique normalized floats

9-bit exponent & 22-bit mantissa can represent  $(2^9 - 1)(2^{22}) = 2^{31} - 2^{22}$  normalized floats

$\therefore$  this system has more normalized floating point numbers



$$a = 5.659 \times 10^4$$

$$c = 9.337 \times 10^2$$

$$B=10 \quad p=4 \quad L=-8 \quad U=8$$

$$b = 5.629 \times 10^4$$

$$d = 7.529 \times 10^{-1}$$

3.

1.  $a+b = 11.288 \times 10^4$  which gets rounded to  $1.129 \times 10^5$

Using  $\nearrow$   
"round to even"

2.  $a-b = 0.03 \times 10^4$  which is represented as  $3.000 \times 10^2$

3.  $c/d = 1240.138133 \xrightarrow{\text{rounds to}} 1.240 \times 10^3$

4.  $b \times d = 42380.741 \longrightarrow 4.238 \times 10^3$

5.  $a+c = 57523.7 \longrightarrow 57520$

$$b+d = 56290.7529 \longrightarrow 56290$$

$$(a+c) - (b+d) = 57520 - 56290 = 1230 \longrightarrow 1.230 \times 10^3$$

Better way of doing calculation:

$$(a-b) + (c-d)$$

$$a-b \longrightarrow 3.000 \times 10^2 \text{ (see part 2)}$$

$$b+d = 932.9471 \longrightarrow 9.329 \times 10^2$$

$$\therefore (a-b) + (c-d) = 1232.9 \longrightarrow 1.233 \times 10^3$$

notice that this is closer to the true value of  $1.2329471 \times 10^3$  than this

difference of  
squares

4. Rather than doing  $x^2 - y^2$ , do  $(x+y)(x-y)$

Why is this better?

Consider the following case, where  $x$  and  $y$  are <sup>very</sup> similar in size:

$y = x + \epsilon$ , where  $\epsilon \ll x$ , but there is still enough precision to represent  $y$  exactly

i.e. the difference between  $x$  and  $y$  ( $\epsilon$ ) is much smaller than  $x$  or  $y$ , but the floating point system has just enough precision to represent  $y$  exactly (so  $x + \epsilon \neq x$ )

If we calculate  $y^2$ , we get  $y^2 = (x + \epsilon)^2 = x^2 + 2x\epsilon + \epsilon^2$

since  $\epsilon \ll x$ ,  $\epsilon^2 \ll x^2$  so there is not enough precision to include it, so  $y^2$  gets represented as  $x^2 + 2x\epsilon$

Aside: (For example, if  $x = 1.0$ ,  $\epsilon = 1.0 \times 10^{-1}$  and  $p = 2$ :  
 $x + \epsilon = 1.1$  which can be represented exactly  
 $x^2 + \epsilon^2 = 1.01$  which is represented as  $1.0$  since  $p = 2$ )

Back to our question:

$y^2$  is represented as  $x^2 + 2x\epsilon$ , so  $x^2 - y^2 = -2x\epsilon$

Now let's do  $(x+y)(x-y)$  instead:

$$\begin{aligned} x+y &= x + (x+\epsilon) & x-y &= x - (x+\epsilon) \\ &= 2x + \epsilon & &= -\epsilon \end{aligned}$$

$$\therefore (x+y)(x-y) = -2x\epsilon - \epsilon^2$$

$\epsilon^2$  and  $2x\epsilon$  are much closer in size than  $\epsilon^2$  and  $x^2$ , so there is (hopefully) enough precision to retain the  $\epsilon^2$  term

ex. going back to our example here:

$$2x\epsilon = 2.0 \times 10^{-1}$$

$$\epsilon^2 = 1.0 \times 10^{-2}$$

$$\therefore -2x\epsilon - \epsilon^2 = -2.1 \times 10^{-1} \text{ which can be represented exactly so the } \epsilon^2 \text{ information is not lost}$$