

CMPT 295

Name: Junchen Li

Student ID: 301385486

Date: 2020/1/29

## Assignment 2

### Objectives:

In this assignment, you will gain familiarity with:

- IEEE floating point representation
- 

### Submission:

- Submit your document called **Assignment\_2.pdf**, which must include your answers to all of the questions in Assignment 2.
    - Add your full name and student number at the top of the first page of your document **Assignment\_2.pdf**.
  - When creating your assignment, first include the question itself and its number then include your answer, keeping the questions in its original numerical order.
  - **If you hand-write your answers (as opposed to using a computer application to write them):** When putting your assignment together, do not take photos (no .jpg) of your assignment sheets! Scan them instead! Better quality -> easier to read -> easier to mark!
  - Submit your assignment electronically on CourSys
- 

### Due:

- Thursday, Jan. 30 at 3pm
  - Late assignments will receive a grade of 0, but they will be marked in order to provide the student with feedback.
- 

### Requirements:

- **Show your work** (as illustrated in lectures).

---

**Marking scheme:**

This assignment will be marked as follows:

- Questions 1 and 2 will be marked for correctness.

The amount of marks for each question is indicated as part of the question.

A solution will be posted after the due date.

---

1. [8 marks] Floating point conversion and Rounding.

- a. Represent the following numbers in IEEE floating point representation (single precision), clearly showing the effect of rounding on the “frac” (mantissa), if rounding occurs, and express your final answer in binary and in hexadecimal:

- I.  $0.001111111_2$        $1.111111_2$  E=-3 exp= $2^{8-1}-1-3=127-3=124_{10}$  & it is a positive number.  $124_{10} \rightarrow 01111100_2$   
The final answer is: 0 01111100 111111.
- II.  $3.1416015625_{10}$        $11.0010010001_2$  E=1  $1.10010010001$  exp: $1+127=128_{10} \rightarrow 10000000_2$  And it is a positive number. The final answer is 0 10000000 10010010001
- III.  $-0.9_{10}$        $-0.1110011001100^{23}_2 \dots$  E=-1 rounding occurs (1 01111110 110011001100 ...)  
(0.1110011001<sub>2</sub> (001 < 1/2 = down) -0.11100...1101<sup>23</sup><sub>2</sub> The final answer is  $\approx -0.111001100110011001100110_2$  and it round to -0.8999999762 the hexadecimal is 0XBF666666.
- IV.  $1/3_{10}$  (a third)       $0.0101010101 \dots$  E=-2 exp: $127-2=125_{10} \rightarrow 01111101_2$  (0 01111101 01010101010101010101011)  
rounding occurs 0.010(101 > 1/2 = up) 0.011<sub>2</sub> The final answer is  $\approx 0.01010101010101010101011_2$  and it rounding to 0.3333333433 the hexadecimal is 0X3EAAAAAB.

- b. Convert 0x4AEA4C1A from IEEE floating point representation (single precision) to a real number.

**0100 1010 1110 1010 0100 1100 0001 1010<sub>2</sub> → 0 10010101 11010100100110000011010 149-127=22 move to right 22 bits'**  
**location 1.11010100100110000011010 11101010010011000001101.0 The final answer is 7677453.**

- c. Round the following binary numbers (rounding position is bolded -  $2^{-4}$  position) following the rounding rules of the IEEE floating point representation.

I.	$1.0011111_2$	<b>111</b> >1/2—Rounding up	<b>1.0100</b>	<b>1.2421875-1.25=-0.0078125</b>
II.	$1.1001001_2$	<b>001</b> <1/2—Rounding down	<b>1.1001</b>	<b>1.5703125-1.5626=0.0077125</b>
III.	$1.0111100_2$	<b>100</b> =1/2— half way—1.0111+0.1=1.1000 even (up to even)	<b>1.1000</b>	<b>1.46875-1.5=-0.03125</b>
IV.	$1.0110100_2$	<b>100</b> =1/2— half way—1.0110+0.1=1.0111 odd (down to even)	<b>1.0110</b>	<b>1.40625-1.375=0.03125</b>

For each of the above rounded binary numbers, indicate what type of rounding you performed and compute the value that is either added to or subtracted from the original number (listed above) as a result of the rounding process. In other words, compute the error introduced by the rounding process.

2. [12 marks] Creating hypothetical smaller floating-point representations based on the IEEE floating point format allows us to investigate this encoding scheme more easily, since the numbers are easier to manipulate and compute.

Below is a table listing several real numbers represented as 6-bit floating-point numbers ( $w = 6$ ). The format of these 6-bit floating-point numbers is as follows: 1 bit is used to express for the sign, 3 bits are used to express “exp” ( $k = 3$ ) and 2 bits are used to represent “frac” ( $n = 2$ ), in the following order: sign exp frac.

Complete the table (the same way as in Figure 2.35 in our textbook) then answer the questions below the table.

Tip: Have a look at Figure 2.35 in our textbook, which illustrates a similar table for a hypothetical 8-bit floating-point format. This will give you an idea of how to complete the table. Also, Figure 2.34 displays the complete range of these 6-bit floating point numbers as well as their values between -1.0 and 1.0. This diagram may be helpful when you are checking your work.

		Exponent			Fraction		Value		
Description	Bit representation	exp	E	$2^E$	frac	M	$M \cdot 2^E$	V (分数)	Decimal
zero	0 000 00	<u>0</u>	<u>1-3=-2</u>	<u>1/4</u>	<u>0/4</u>	<u>0/4=0</u>	<u>0</u>	<u>0</u>	<u>0.0</u>
Smallest positive denormalized	0 000 01	<u>0</u>	<u>-2</u>	<u>1/4</u>	<u><math>2^{-2}=1/4</math></u>	<u>1/4</u>	<u>1/16</u>	<u><math>(-1)^0 \cdot 1/16</math></u> <u>=1/16</u>	<u>0.0625</u>
	0 000 10	<u>0</u>	<u>-2</u>	<u>1/4</u>	<u><math>2^{-1}=1/2</math></u>	<u>1/2</u>	<u>1/8</u>	<u><math>(-1)^0 \cdot 1/8</math></u> <u>=1/8</u>	<u>0.125</u>
Largest positive denormalized	0 000 11	<u>0</u>	<u>-2</u>	<u>1/4</u>	<u><math>2^{-1}+2^{-2}=3/4</math></u>	<u>3/4</u>	<u>3/16</u>	<u><math>(-1)^0 \cdot 3/16</math></u> <u>=3/16</u>	<u>0.1875</u>
Smallest positive normalized	0 001 00	<u>1</u>	<u>-2</u>	<u>1/4</u>	<u>0/4</u>	<u><math>1+0/4=4/4=1</math></u>	<u>1/4</u>	<u><math>(-1)^0 \cdot 1/4</math></u> <u>=1/4</u>	<u>0.25</u>
	0 001 01	<u>1</u>	<u>-2</u>	<u>1/4</u>	<u>1/4</u>	<u>5/4</u>	<u>5/16</u>	<u><math>(-1)^0 \cdot 5/16</math></u> <u>=5/16</u>	<u>0.3125</u>
	0 001 10	<u>1</u>	<u>-2</u>	<u>1/4</u>	<u>1/2</u>	<u>3/2</u>	<u>3/8</u>	<u><math>(-1)^0 \cdot 3/8</math></u> <u>=3/8</u>	<u>0.375</u>
	0 001 11	<u>1</u>	<u>-2</u>	<u>1/4</u>	<u>3/4</u>	<u>7/4</u>	<u>7/16</u>	<u><math>(-1)^0 \cdot 7/16</math></u> <u>=7/16</u>	<u>0.4375</u>

	0 010 00	<u>2</u>	<u>-1</u>	<u>1/2</u>	<u>0</u>	<u>1+0=1</u>	<u>1/2</u>	$(-1)^0 \cdot 1/2$ <u>=1/2</u>	<u>0.5</u>
	0 010 01	<u>2</u>	<u>-1</u>	<u>1/2</u>	<u>1/4</u>	<u>5/4</u>	<u>5/8</u>	$(-1)^0 \cdot 5/8$ <u>=5/8</u>	<u>0.625</u>
	0 010 10	<u>2</u>	<u>-1</u>	<u>1/2</u>	<u>1/2</u>	<u>3/2</u>	<u>3/4</u>	$(-1)^0 \cdot 3/4$ <u>=3/4</u>	<u>0.75</u>
	0 010 11	<u>2</u>	<u>-1</u>	<u>1/2</u>	<u>3/4</u>	<u>7/4</u>	<u>7/8</u>	$(-1)^0 \cdot 7/8$ <u>=7/8</u>	<u>0.875</u>
One	0 011 00	<u>3</u>	<u>0</u>	<u>1</u>	<u>0</u>	<u>1</u>	<u>1*1=1</u>	$(-1)^0 \cdot 1$ <u>=1</u>	<u>1.0</u>
	0 011 01	<u>3</u>	<u>0</u>	<u>1</u>	<u>1/4</u>	<u>5/4</u>	<u>5/4</u>	$(-1)^0 \cdot 5/4$ <u>=5/4</u>	<u>1.25</u>
	0 011 10	<u>3</u>	<u>0</u>	<u>1</u>	<u>1/2</u>	<u>3/2</u>	<u>3/2</u>	$(-1)^0 \cdot 3/2$ <u>=3/2</u>	<u>1.5</u>
	0 011 11	<u>3</u>	<u>0</u>	<u>1</u>	<u>3/4</u>	<u>7/4</u>	<u>7/4</u>	$(-1)^0 \cdot 7/4$ <u>=7/4</u>	<u>1.75</u>
	0 100 00	<u>4</u>	<u>1</u>	<u>2</u>	<u>0</u>	<u>1</u>	<u>2</u>	$(-1)^0 \cdot 2$ <u>=2</u>	<u>2</u>
	0 100 01	<u>4</u>	<u>1</u>	<u>2</u>	<u>1/4</u>	<u>5/4</u>	<u>5/2</u>	$(-1)^0 \cdot 5/2$ <u>=5/2</u>	<u>2.5</u>
	0 100 10	<u>4</u>	<u>1</u>	<u>2</u>	<u>1/2</u>	<u>3/2</u>	<u>3</u>	$(-1)^0 \cdot 3$ <u>=3</u>	<u>3</u>
	0 100 11	<u>4</u>	<u>1</u>	<u>2</u>	<u>3/4</u>	<u>7/4</u>	<u>7/2</u>	$(-1)^0 \cdot 7/2$	<u>3.5</u>

								<u>=3</u>	
	0 101 00	<u>5</u>	<u>2</u>	<u>4</u>	<u>0</u>	<u>1</u>	<u>4</u>	$(-1)^0 * 4 = 4$	<u>4</u>
	0 101 01	<u>5</u>	<u>2</u>	<u>4</u>	<u>1/4</u>	<u>5/4</u>	<u>5</u>	$(-1)^0 * 5$	<u>5</u>
	0 101 10	<u>5</u>	<u>2</u>	<u>4</u>	<u>1/2</u>	<u>3/2</u>	<u>6</u>	$(-1)^0 * 6$	<u>6</u>
	0 101 11	<u>5</u>	<u>2</u>	<u>4</u>	<u>3/4</u>	<u>7/4</u>	<u>7</u>	$(-1)^0 * 7$	<u>7</u>
	0 110 00	<u>6</u>	<u>3</u>	<u>8</u>	<u>0</u>	<u>1</u>	<u>8</u>	$(-1)^0 * 8 = 8$	<u>8</u>
	0 110 01	<u>6</u>	<u>3</u>	<u>8</u>	<u>1/4</u>	<u>5/4</u>	<u>10</u>	$(-1)^0 * 10$	<u>10</u>
	0 110 10	<u>6</u>	<u>3</u>	<u>8</u>	<u>1/2</u>	<u>3/2</u>	<u>12</u>	$(-1)^0 * 12$	<u>12</u>
Largest positive normalized	0 110 11	<u>6</u>	<u>3</u>	<u>8</u>	<u>3/4</u>	<u>7/4</u>	<u>14</u>	$(-1)^0 * 14$	<u>14</u>
+ Infinity	<u>0 111 00</u>	–	–	–	–	–	–	$\infty$	–
NaN	<u>0 111 10</u>  (frac part cannot be all zero)	–	–	–	–	–	–	NaN	–

a. What is the value of the bias? Answer: #exp=3, #frac=2 Bias=  $2^{3-1}-1=3$

b. Consider two adjacent denormalized numbers. How far apart are they? Express this difference (“delta”) as a decimal number.

Answer:

0 001 10      0.375<sub>10</sub>

0 001 11      0.4375<sub>10</sub>

difference=0.4375-0.375= 0.0625

c. Consider two adjacent normalized numbers ...

a. with the **exp** field set to 001. How far apart are they?

Answer: 0 001 10->0.375<sub>10</sub>    0 001 11-> 0.4375<sub>10</sub>    difference=0.4375-0.375=0.0625.  $2^{-4}$

b. with the **exp** field set to 010. How far apart are they?

Answer: 0 010 10->0.75<sub>10</sub>    0 010 11->0.875<sub>10</sub>    difference=0.875-0.75 =0.125.  $2^{-3}$

c. with the **exp** field set to 011. How far apart are they?

Answer: 0 011 10->1.5<sub>10</sub>    0 011 11->1.75<sub>10</sub>    difference=1.75-1.5=0.25.  $2^{-2}$

Expressed these differences (“delta”) as decimal numbers.

d. Without doing any calculations, can you guess how far apart are two adjacent normalized numbers ...

a. with the **exp** field set to 100?  $2^{-1}=1/2=0.5$

b. with the **exp** field set to 101?  $2^0=1$

c. with the **exp** field set to 110?  $2^1=2$

e. What is the “range” (not contiguous) of real numbers that can be represented using this 6-bit floating-point representation?

Answer: the range is between -14 and 14. Range: [-14~14].



- f. What is the range of the normalized exponent  $E$  ( $E$  found in the equation  $v = (-1)^s M 2^E$ ) which can be represented by this 6-bit floating-point representation?

**Answer: Exp: 0~6 E: [-2~3]**

- g. Give an example of a real number that cannot be represented using this 6-bit floating-point representation, but is within the “range” of representable values.

**Answer: for example, 0 101 00->4<sub>10</sub> and 0 101 01->5<sub>10</sub>. The number between 4 and 5 cannot be representation by 6-bit floating-point. Like 4.5.**

- h. Give an example of a real number that would overflow if we were trying to represent it using this 6-bit floating-point representation. The best way to answer this question is to convert this real number into a 6-bit IEEE floating-point representation and clearly indicate why it would overflow.

**Answer: When over the biggest decimal value 14, going beyond this overflows to  $+\infty$ . For example the decimal number 15<sub>10</sub>. The binary of 15 is 01111<sub>2</sub> which cannot put 111 into exp part because it will become special value  $+\infty/-\infty$ . Frac part=000...0 which means operation that overflows. So 15 would have overflow occurred.**

- i. How close is the value of the “frac” of the largest normalized number to 1? In other words, what is  $\epsilon$  (epsilon)? Expressed it as a decimal number.

**Answer: 1-3/4=1/4. The epsilon is 1/4.**