## CSC205 section 1 Spring 2015 Homework 2

Marcy

## **How to Submit:**

Please submit your solutions (parts A and B separately) through Blackboard. Remember to put your name and homework number on all the documents that you submit as attachments.

Total possible points in this homework: 3 for Part B (You receive 1 bonus point towards Part A if you get all Part A questions correct.)

## Part A

- 1. You are given the 32-bit sequence 10000000 00000000 11111111 11111111 (or 0x8000FFFF). Answer the questions below.
  - a. This code represents the integer -65535 in a sign-magnitude system. How to represent +65535 in the same system?
  - b. This code represents the integer -2147418112 in a 1's complement system. How to represent +2147418112? In the same system?
  - c. This code represents the integer -2147418113 in a 2's complement system. How to represent +2147418113 in the same system?
- 2. Perform the addition of the following numbers in 8-bit 2's complement representation and indicate whether overflow occurs.

For example, 72+28:

_	51 example, 72 <u>· 20</u> :										
		Carry				1	1				
		72		0	1	0	0	1	0	0	0
	+	28		0	0	0	1	1	1	0	0
•		100		0	1	1	0	0	1	0	0

No overflow

a. -72+28:

	Carry		1	1	1				
	-72	1	0	1		1	0	Q	0
+	28	B	0	6	1	1	1	٥	6
	-44	1	1	0	1	0	1	6	6

Overflow?

b. 72+70:

Carry	1							
72	0	1	0	0	1	0	0	0
+ 70	0	1	0	0	Ô	1		0
142	1	0	O	0	1	1	1	0

Overflow?

c. -72-70:

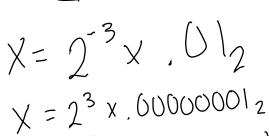
C	arry	1		7	1	١				
	-72		1	0		)	1	0	0	0
+	-70		l	$\bigcirc$	1	١	)	0	1	6
_	142		()	1	)	١	0	0	-	0

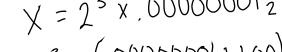
- Overflow?
- 3. In the IEEE-754 single precision floating point standard, what decimal number does the sequence 10111111 10000000 00000000 00000000 (or 0xBF800000) represent? You may leave it in the form  $(+/-)2^L \times K_{10}$ , where L and K are in decimal.

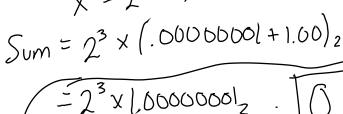
Hint: You may use the following table to help you analyze the number.

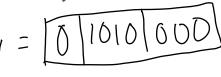
int. Tou may use the following table to help you unaryze the number.									
Field	Bit values	Meaning, or							
		decimal equivalent							
Sign =		-							
Biased exponent value =	127	20							
True exponent value =	0	20							
Significand value =	$\mathcal{O}$								
The number is:	1.0 × 2°	-1,0							

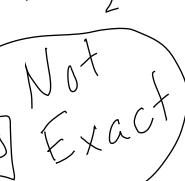
4. Consider a floating point number excess-7 system with 4-bit exponent (where values 0000 and 1111 are reserved), 3-bit significand, with implied binary point and implied 1. Perform the addition of the two numbers represented respectively by 0 0100 010 and 0 1010 000. Show how the answer is represented in this system. Is the representation exact?











## Part B

- 1. [Total 3 pts] This question examines the limitations of the *int* and the *float* data types in Java, in which an integer is 32-bit and a float uses the IEEE 754 single precision format.
  - a. [1pt] The following code segment computes the *exact* value of  $C_R^N = \frac{N!}{R!\cdot (N-R)!}$  for the special case of  $R = \lfloor N/2 \rfloor$  (meaning R is an integer that is half of N). Results are shown to the right of the code. Identify the type of error that occurs in the computation. Explain why it starts to occur at N=18 and not earlier or later.

```
Code segment 1 (in Java)
                                                           Output for N=1..30
public class Combination {
                                                           N=1
                                                                C(1,0)=1
  private static final int defaultNMax = 30;
                                                           N=2
                                                                C(2,1)=2
                                                           N=3 C(3,1)=3
  public Combination() {}
                                                           N=4 C(4,2)=6
                                                           N=5 C(5,2)=10
  int compute(int N, int R) {
                                                           N=6 C(6,3)=20
                                                           N=7
    int i;
                                                                C(7,3)=35
    int numerator;
                                                           N=8 C(8,4)=70
    int denominator;
                                                           N=9 C(9,4)=126
                                                           N=10 C(10,5)=252
                                                           N=11 C(11,5)=462
    numerator = 1;
    for (i=N-R+1; i<=N; i+=1)
                                                           N=12 C(12,6)=924
      numerator = numerator*i;
                                                           N=13 C(13,6)=1716
                                                           N=14 C(14,7)=3432
                                                           N=15 C(15,7)=6435
    denominator = 1;
    for (i=1; i \le R; i+=1)
                                                           N=16 C(16,8)=12870
      denominator = denominator*i;
                                                           N = 17
                                                                 C(17,8)=24310
                                                           N=18 C(18,9)=1276
                                                           N=19 C(19,9)=-2308
    return numerator / denominator;
                                                           N=20 C(20,10)=117
                                                           N=21 C(21,10)=9
  public static void main(String args[]) {
                                                           N=22 C(22,11)=19
    Combination c = new Combination();
                                                           N = 23
                                                                 C(23,11)=1
    int Nmax = defaultNMax;
                                                           N=24 C(24,12)=3
                                                           N=25 C(25,12)=0
    for (int N=1; N<=Nmax; N+=1) {</pre>
                                                           N=26 C(26,13)=0
      int result = c.compute(N, N/2);
System.out.println("N=" + N + "
                                                           N=27 C(27,13)=0
                                        C(" + N + "," +
                                                           N=28 C(28,14)=0
N/2 + ") = " + result);
                                                           N = 29
                                                                 C(29,14)=-1
                                                           N=30 C(30,15)=0
 }
}
```

- b. [1 pt] In the following code segment, the programmer has converted the data type from *int* to *float*, and the results are shown to the right of the code.
  - i. What is the first type of error that would occur? Explain.
  - ii. At what value of N would you start to loose confidence on the result when using *float* to compute exact integers?

(Hint: Consider what is the smallest positive integer that cannot be represented exactly using *float*.)

```
Output for N=1..30
Code segment 2 (in Java)
public class CombinationF {
                                                 N=1 C(1,0)=1.0
 private static final int defaultNMax = 30;
                                                 N=2
                                                      C(2,1)=2.0
                                                 N=3 C(3,1)=3.0
  public CombinationF() {}
                                                 N=4 C(4,2)=6.0
                                                 N=5 C(5,2)=10.0
                                                 N=6 C(6,3)=20.0
  float compute(float N, float R) {
    float i;
                                                 N=7 C(7,3)=35.0
                                                 N=8 C(8,4)=70.0
    float numerator;
    float denominator:
                                                 N=9 C(9,4)=126.0
                                                 N=10 C(10,5)=252.0
    numerator = 1;
                                                 N = 11
                                                       C(11,5)=462.0
                                                       C(12,6)=924.0
    for (i=N-R+1; i<=N; i+=1)
                                                 N = 12
      numerator = numerator*i;
                                                 N = 13
                                                       C(13,6)=1716.0
                                                 N = 14
                                                       C(14,7)=3432.0
    denominator = 1;
                                                 N = 15
                                                       C(15,7)=6435.0
    for (i=1; i<=R; i+=1)
                                                 N = 16
                                                       C(16,8)=12870.0
      denominator = denominator*i;
                                                 N = 17
                                                       C(17,8)=24310.0
                                                       C(18,9)=48620.0
                                                 N = 1.8
    return numerator / denominator;
                                                 N = 19
                                                       C(19,9) = 92378.0
  }
                                                 N = 20
                                                       C(20,10)=184756.0
                                                       C(21,10)=352716.0
                                                 N = 21
  public static void main(String args[]) {
                                                 N = 22
                                                       C(22,11) = 705432.0
    CombinationF c = new CombinationF();
                                                 N = 23
                                                       C(23,11)=1352078.0
    int Nmax = defaultNMax;
                                                 N = 24
                                                       C(24,12)=2704156.0
                                                 N=25
                                                       C(25,12)=5200300.5
                                                       C(26,13)=1.0400601E7
    for (int N=1; N<=Nmax; N+=1) {
                                                 N = 26
      float result = c.compute(N, N/2);
                                                 N = 27
                                                       C(27,13)=2.00583E7
      System.out.println("N=" + N + " C(" +
                                                 N=28 C(28,14)=4.0116604E7
N + "," + N/2 + ") = " + result);
                                                 N = 29
                                                       C(29,14)=7.7558776E7
    }
                                                 N=30 C(30,15)=1.55117552E8
 }
}
```

c. [1 pt] Suggest some possible ways that a programmer can protect the user from getting erroneous results when computing  $C_R^N = \frac{N!}{R!\cdot(N-R)!}$ .