

1. Como es bien sabido, la representación numérica binaria de tipo signo/exponente/mantisa, depende de la creatividad o precisión que desee el programador. Analicemos qué sucede con un sistema de numeración binario a $n > 0$ (entero) bits, para esto elige un valor par entre 10 y 20, crea el sistema de numeración entero similar al estudiado en clase, (la mantisa debe ser, al menos, la mitad del número elegido). Describe tu sistema de numeración, ¿Cuál es el menor y mayor de los números representados? Presenta un número que no se pueda representar en tu sistema seleccionado. Describe tu proceso. (4 pts) NOTA: punto extra si además se describe un sistema fraccionario que tenga sentido.

Demostración. Usaremos 16 bits en total. 1 para el signo, 3 para el exponente, 12 para la mantisa. Esto nos genera números en el rango

$$[-1 * 2^7 * (2^{12} - 1), 2^7 * (2^{12} - 1)] = [-524160, 524160]$$

Un número que no podemos representar es el $2^{12} + 1$, debido a que no es un múltiplo de dos y es mayor a $2^{12} - 1$, el mayor numero que podemos representar con la mantisa.

Listing 1: Generación de números posibles con el sistema elegido. Fraccionario y entero.

```
#include <stdio.h>
#define lsb 0.000244140625

void print_binary(FILE *file , int n) {
    for (int i = 11; i >= 0; i--) {
        int bit = (n >> i) & 1;
        fprintf(file , "%d" , bit);
    }
    fprintf(file , ",");
}

/*

16 bit machine:
1 bit for sign
3 bits for exponent
12 bits for mantissa

sgn * 2^exp * mantissa
max. number represented:
0111111111111111 = (2^(12) - 1) * 2^7 = 4095 * 2^7 = 524160

min. number represented:
1111111111111111 = -524160

impossible to represent the number 2^12 + 1 = 4097

*/

/*

print_representation() writes to a csv file "binary.csv"
all possible number representations in decimal format.
we then confirm that 4097 is not present.

*/
```

```

void print_representation() {
    FILE *file;
    file = fopen("binary.csv", "w+");
    /*
        header for the csv, representing
        the exponent of 2 which the mantissa
        is being multiplied by
    */
    fprintf(file, "binary, -
        2^{0},2^{1},2^{2},2^{3},2^{4},2^{5},2^{6},2^{7}\n");

    for (int i = 1; i < 1 << 12;
        i++) { // loop from i = 1 (we ignore 0), up to 2^12 - 1
        print_binary(file, i);
        for (int j = 0; j < 1 << 3; j++) { // loop from j = 0 up to 2^3 -
            1
            if (j == 7) {
                fprintf(file, "%d\n", i * (1 << j)); // mantissa * 2^j
            } else {
                fprintf(file, "%d,", i * (1 << j));
            }
        }
    }
}

void print_fractions() {
    FILE *file;
    file = fopen("fractions.csv", "w+");
    fprintf(file, "binary, -2^{-2},2^{-2},2^{-1},2^{0},2^{1},2^{2},2^{3}\n");
    int curr;
    double i;
    double mantissa;
    float exp;

    for (i = lsb, curr = 1; i < 1; i += lsb, curr++) {
        print_binary(file, curr);
        for (int j = 0; j < 7; j++) {
            if (j == 0) {
                exp = 1.0/4; // 2^{-2}
                mantissa = i; // implicit bit is 0
            }
            else {
                exp = (j >= 3) ? (1 << (j-3)) : 1.0 / (1 << (3 - j)); // using
                    2^{j-3} = 1/2^{3-j}
                mantissa = 1 + i; // implicit bit is 1
            }

            if (j == 6) {
                fprintf(file, "%f\n", mantissa * exp);
            } else {
                fprintf(file, "%f,", mantissa * exp);
            }
        }
    }
}

```

```

    }
  }
}

int main() {
    print_representation();
    print_fractions();
    return 0;
}

```

El código anterior genera dos archivos csv, uno con todos los valores (positivos) posibles de la representación elegida, y otro con los valores fraccionarios, que describimos más adelante. En la siguiente tabla, presentamos algunos valores enteros representables por nuestro sistema. El encabezado de la columna representa el exponente que multiplica la mantisa.

binary	2^0	2^1	2^2	2^3	2^4	2^5	2^6	2^7
000000000001	1	2	4	8	16	32	64	128
000000000010	2	4	8	16	32	64	128	256
000000000011	3	6	12	24	48	96	192	384
000000000100	4	8	16	32	64	128	256	512
000000000101	5	10	20	40	80	160	320	640
000000000110	6	12	24	48	96	192	384	768
000000000111	7	14	28	56	112	224	448	896
000000001000	8	16	32	64	128	256	512	1024
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
111111111111	4095	8190	16380	32760	65520	131040	262080	524160

Para el sistema fraccionario, definimos el bit implícito como

$$\begin{cases} b_{imp} = 1 & \text{if } 0 < \text{exp} < 7 \\ b_{imp} = 0 & \text{if } \text{exp} = 0 \text{ and } \text{mantissa} \neq 0 \end{cases}$$

y

$$\begin{cases} \text{exp} = \text{exp} - 3 & \text{if } 0 < \text{exp} < 7, \\ \text{exp} = -2 & \text{if } \text{exp} = 0, \end{cases}$$

y finalmente, cuando $\text{exp} = 7$,

$$\begin{cases} \pm\infty & \text{if } \text{mantissa} = 0, \\ \text{NaN} & \text{if } \text{mantissa} \neq 0. \end{cases}$$

De esta manera, obtenemos los valores (truncados para una mejor presentación)

binary	2^{-2}	2^{-2}	2^{-1}	2^0	2^1	2^2	2^3
000000000001	0.000061	0.250061	0.500122	1.000244	2.000488	4.000977	8.001953
000000000010	0.000122	0.250122	0.500244	1.000488	2.000977	4.001953	8.003906
000000000011	0.000183	0.250183	0.500366	1.000732	2.001465	4.002930	8.005859
000000000100	0.000244	0.250244	0.500488	1.000977	2.001953	4.003906	8.007812
000000000101	0.000305	0.250305	0.500610	1.001221	2.002441	4.004883	8.009766
000000000110	0.000366	0.250366	0.500732	1.001465	2.002930	4.005859	8.011719
000000000111	0.000427	0.250427	0.500854	1.001709	2.003418	4.006836	8.013672
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
111111111111	0.249939	0.499939	0.999878	1.999756	3.999512	7.999023	15.998047

De manera análoga, podemos obtener los valores negativos. Entonces, el número máximo (finito) de la representación fraccionaria es

$$2^3 \sum_{i=1}^6 2^{-i} = 15.998046875$$

mientras que el mínimo es el negativo del máximo. El más cercano a 0 es $\pm 2^{-12} = 0.000244140625$. \square

2. Consideremos la siguiente sucesión $\{x_n\}_{n \in \mathbb{N}} \subset \mathbb{R}$ definida de manera recurrente como: (4 pts)

$$x_2 = 2,$$

$$x_{n+1} = 2^{n-\frac{1}{2}} \sqrt{1 - \sqrt{1 - 4^{1-n} x_n^2}}$$

- (a) Escribe un programa para encontrar los primeros n términos de la sucesión.
(b) Aplica tu programa para encontrar los términos de la sucesión cuando

$$n = 6, 7, 8, 9, 10, 12, 18, 20$$

, ¿a qué valor tiende la sucesión? (ésta sucesión tiende a un valor distinto de cero).

- (c) Encuentra los términos 50 y 100, ¿qué sucede con la sucesión? Grafica los términos de la sucesión desde 2 hasta 100. Explica detalladamente el por qué del comportamiento que tiene la computadora al calcular estos términos.

Demostración. Primero mostraremos que utilizar directamente la secuencia resulta en los términos divergiendo a infinito debido a los errores de precisión de punto flotante. Podemos notar que la secuencia puede ser reescrita de tal manera que la cantidad de restas se puede reducir a 1.

$$\begin{aligned}
x_{n+1} &= 2^{n-1/2} \sqrt{1 - \sqrt{1 - 4^{1-n} x_n^2}} \\
&= 2^{n-1/2} \sqrt{1 - \sqrt{1 - 4^{1-n} x_n^2}} \cdot \frac{\sqrt{1 - \sqrt{1 + 4^{1-n} x_n^2}}}{\sqrt{1 + \sqrt{1 - 4^{1-n} x_n^2}}} \\
&= 2^{n-1/2} \frac{\sqrt{1 - (1 - 4^{1-n} x_n^2)}}{\sqrt{1 - \sqrt{1 + 4^{1-n} x_n^2}}} \\
&= 2^{n-1/2} \frac{2^{1-n} x_n}{\sqrt{1 - \sqrt{1 + 4^{1-n} x_n^2}}} \\
&= \frac{\sqrt{2} x_n}{\sqrt{1 - \sqrt{1 + 4^{1-n} x_n^2}}}.
\end{aligned}$$

Listing 2: Calcular secuencia de números que se acerca a π

```
#include <math.h>
#include <stdio.h>
#include <stdlib.h>
#define Sqrt2 1.4142135624

/*
this implementation rewrites the formula into


$$x_{n+1} = \frac{\sqrt{2}x_n}{\sqrt{1+\sqrt{1-4^{1-n}}x_n^2}}$$


although using division, this reduces the number of subtractions used,
which leads to less floating point precision errors
accumulating, resulting in the sequence staying at the
convergence point for longer.

*/

double non_naive(int n) {
    FILE *file;
    file = fopen("p2.csv", "w+");
    double x = 2.0; // first term
    double four_pow_one_minus_n = 1.0; // the term  $4^{1-n}$ , divided by 4
    each iteration
    for (int i = 3; i < n + 2; i++) {
        x = Sqrt2 * x / (sqrt(1 + sqrt(1 - (four_pow_one_minus_n /= 4) * x
            * x)));
        fprintf(file, "x_%d, -%.16lf\n", i, x);
    }
    return x;
}

/*
Naive implementation: this implementation uses the sequence formula


$$x_{n+1} = 2^{n-1/2} \sqrt{1-\sqrt{1-4^{1-n}}x_n^2}$$


This leads to the sequence diverging due to floating point precision
at the term 29, which then spirals into infinity.

*/
double naive(int n) {
    double x = 2.0; // first term
    double two_pow_n_minus_one = Sqrt2; // term  $2^{n-1/2}$ , multiplied by
    2 each iteration
    double four_pow_one_minus_n = 1.0; // term  $4^{1-n}$ , divided by 4
    each iteration
    FILE *file;
    file = fopen("naive_p2.csv", "w+");
    fprintf(file, "x_2, -%.16lf\n", x);
}
```

```

for (int i = 3; i < n + 2; i++) {
    x = (two_pow_n_minus_one *= 2) *
        sqrt(1 - sqrt(1 - (four_pow_one_minus_n /= 4) * x * x));
    // Note that the assignation of two_pow_n_minus_one is done before
    // being used by x
    // Similarly with four_pow_one_minus_n
    fprintf(file , "x-%d, -%.16lf\n", i, x);
}
return x;
}

int main(int argc, char *argv[]) {
    if (argc != 2) {
        fprintf(stderr, "Invalid arguments.-Usage: ./p2-(terms)");
        exit(1);
    }
    int terms = atoi(argv[1]);
    non_naive(terms);
    naive(terms);
    return 0;
}

```

Figure 1: Secuencia usando algoritmo ingenuo

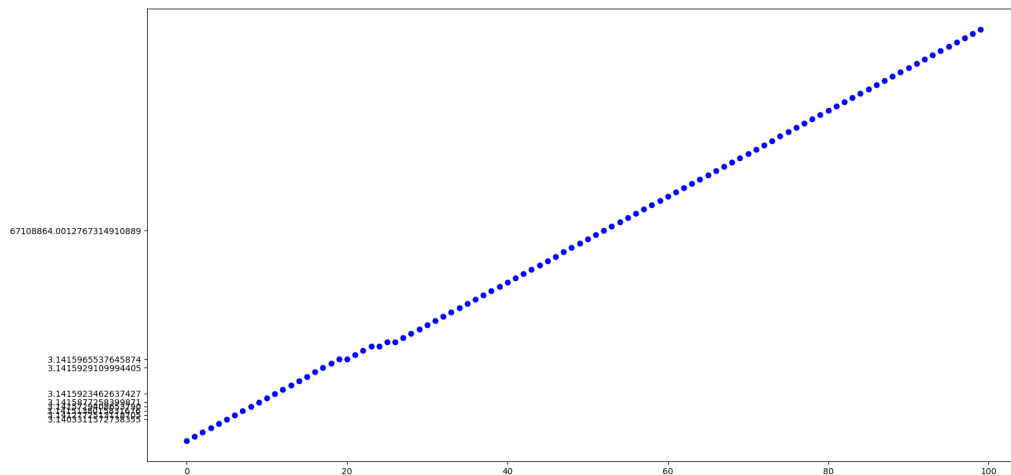
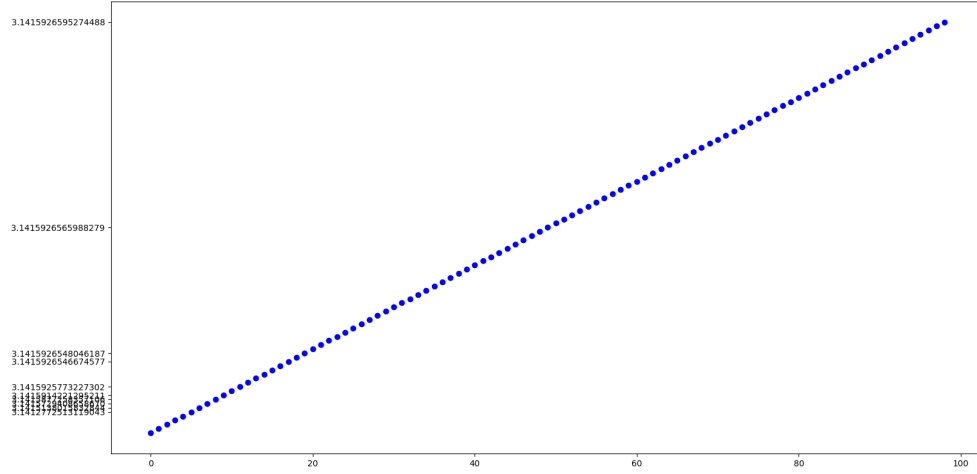


Figure 2: Secuencia usando algoritmo reduciendo restas



En el algoritmo ingenuo,

$$\begin{aligned}
 x_5 &= 3.1214451524539055 \\
 x_6 &= 3.1365484908043393 \\
 x_7 &= 3.1403311572738355 \\
 x_8 &= 3.1412772513118705 \\
 x_9 &= 3.1415138015831676 \\
 x_{10} &= 3.1415729408653790 \\
 x_{11} &= 3.1415877258399871 \\
 x_{12} &= 3.1415914221202459 \\
 x_{13} &= 3.1415923462637427 \\
 x_{14} &= 3.1415925771976707 \\
 x_{15} &= 3.1415926358946096 \\
 x_{16} &= 3.1415926548673574 \\
 x_{17} &= 3.1415926833264782 \\
 x_{18} &= 3.1415927592174677 \\
 x_{19} &= 3.1415929109994405 \\
 x_{20} &= 3.1415941252549588 \\
 x_{50} &= 4194304.0000797957181931 \\
 x_{100} &= 4722366482959487205376
 \end{aligned}$$

mientras que en el otro, tenemos

$$\begin{aligned}x_5 &= 3.1214451524539042 \\x_6 &= 3.1365484908043384 \\x_7 &= 3.1403311572738337 \\x_8 &= 3.1412772513119043 \\x_9 &= 3.1415138015832844 \\x_{10} &= 3.1415729408658670 \\x_{11} &= 3.1415877258357106 \\x_{12} &= 3.1415914221295211 \\x_{13} &= 3.1415923462482072 \\x_{14} &= 3.1415925773227302 \\x_{15} &= 3.1415926351361882 \\x_{16} &= 3.1415926496343785 \\x_{17} &= 3.1415926533037521 \\x_{18} &= 3.1415926542659216 \\x_{19} &= 3.1415926545512902 \\x_{20} &= 3.1415926546674577 \\x_{50} &= 3.1415926564792924 \\x_{100} &= 3.1415926594676811\end{aligned}$$

Como podemos observar, el segundo algoritmo (que reduce el numero de restas), se mantiene cerca de π , mientras que el ingenuo explota y diverge a infinito debido a las aproximaciones que se deben hacer con las restas de punto flotante. \square

3. Al número representable inmediatamente después de la unidad se le conoce como el ϵ de la computadora machine el cual nos permite calcular el número real representable inmediatamente posterior; ésto se obtiene al multiplicar el ϵ por el real y sumarlo a ese real. Utilizando el siguiente algoritmo, crea un código que calcule el machine de tu computadora, sólo hace falta indicar el tipo de variable que es "epsilon" (float o doble). El último valor impreso corresponde al epsilon de la computadora. (2 pts)

Demostración. La lista impresa es

1.00000000000000000000000000000000
0.50000000000000000000000000000000
0.25000000000000000000000000000000
0.12500000000000000000000000000000
0.06250000000000000000000000000000
0.03125000000000000000000000000000
0.01562500000000000000000000000000
0.00781250000000000000000000000000
0.00390625000000000000000000000000
0.00195312500000000000000000000000
0.00097656250000000000000000000000
0.00048828125000000000000000000000
0.00024414062500000000000000000000
0.00012207031250000000000000000000
0.00006103515625000000000000000000


```

0.0000305175781250000000
0.0000152587890625000000
0.0000076293945312500000
0.0000038146972656250000
0.0000019073486328125000
0.0000009536743164062500
0.0000004768371582031250
0.0000002384185791015625
0.0000001192092895507812
0.0000000596046447753906
0.0000000298023223876953
0.0000000149011611938477
0.0000000074505805969238
0.0000000037252902984619
0.0000000018626451492310
0.0000000009313225746155
0.0000000004656612873077
0.0000000002328306436539
0.0000000001164153218269
0.0000000000582076609135
0.0000000000291038304567
0.0000000000145519152284
0.0000000000072759576142
0.0000000000036379788071
0.0000000000018189894035
0.0000000000009094947018
0.0000000000004547473509
0.0000000000002273736754
0.0000000000001136868377
0.0000000000000568434189
0.0000000000000284217094
0.0000000000000142108547
0.0000000000000071054274
0.0000000000000035527137
0.0000000000000017763568
0.0000000000000008881784
0.0000000000000004440892
0.0000000000000002220446

```

Después de 52 iteraciones, obtenemos que

$$\epsilon = 0.0000000000000002220446 = 2.22046 \times 10^{-16} \approx 2^{-52}$$

.

Listing 3: Calcular epsilon de la máquina

```

#include <stdio.h>

int main() {
    double epsilon = 1;
    while (1+epsilon > 1) {
        printf("%.22f\n", epsilon);
        epsilon /= 2;
    }
}

```

