

# Tarea 1 Calculo Computacional

Victor Tortolero CI:24.569.609

## Respuesta 1

Tenemos que  $\frac{A+3}{13}$ , como  $A = 9$ , tendríamos  $\frac{9+3}{13} = \frac{12}{13}$ . Ahora procedemos a convertir a binario.

$\frac{12}{13} \times 2 = \frac{24}{13}, b_0 = 1$		$\frac{2}{13} \times 2 = \frac{4}{13}, b_7 = 0$
$\frac{11}{13} \times 2 = \frac{22}{13}, b_1 = 1$		$\frac{4}{13} \times 2 = \frac{8}{13}, b_8 = 0$
$\frac{9}{13} \times 2 = \frac{18}{13}, b_2 = 1$		$\frac{8}{13} \times 2 = \frac{16}{13}, b_9 = 1$
$\frac{5}{13} \times 2 = \frac{10}{13}, b_3 = 0$		$\frac{3}{13} \times 2 = \frac{6}{13}, b_{10} = 0$
$\frac{10}{13} \times 2 = \frac{20}{13}, b_4 = 1$		$\frac{6}{13} \times 2 = \frac{12}{13}, b_{11} = 0$
$\frac{7}{13} \times 2 = \frac{14}{13}, b_5 = 1$		
$\frac{1}{13} \times 2 = \frac{2}{13}, b_6 = 0$		

Por lo tanto tenemos que:

$$0,111011000100\overline{111011000100}1\dots$$

Observemos que el numero que vendria luego del bit 24 seria un 1. Entonces a la hora de redondear se suma 1. Por lo tanto, tenemos que  $Fl(\frac{12}{13})_{Truncado} = 0,111011000100111011000100$ , y que  $Fl(\frac{12}{13})_{Redondeado} = 0,111011000100111011000101$ .

**Por Truncamiento tenemos que:**

$$\begin{aligned}
 E_A &= |x - Fl(x)_{Truncado}| = 0,\underbrace{000\dots 000}_{24 \text{ Ceros}}111011000100\dots \\
 &= 0,\underbrace{111011000100111011000100}_{\text{Esto es } \frac{12}{13}}\dots \times 2^{-24} \\
 &= \frac{12}{13} \times 2^{-24} \approx 5,50196 \times 10^{-8} \\
 E_R &= \frac{E_A}{|x|} = \frac{\frac{12}{13} \times 2^{-24}}{\frac{12}{13}} \\
 &= 2^{-24} \approx 5,96046 \times 10^{-8}
 \end{aligned}$$

**Por Redondeo tenemos que:**

$$\begin{aligned}
 E_A &= |x - Fl(x)_{Redondeado}| = |x - (Fl(x)_{Truncado} + 1 \times 2^{-24})| \\
 &= |x - Fl(x)_{Truncado} - 1 \times 2^{-24}| \\
 &= \left| \frac{12}{13} \times 2^{-24} - 1 \times 2^{-24} \right| \\
 &= \left| \frac{12}{13} - 1 \right| \times 2^{-24} \\
 &= \frac{1}{13} \times 2^{-24} \approx 4,584 \times 10^{-9} \\
 E_R &= \frac{E_A}{|x|} = \frac{\frac{1}{13} \times 2^{-24}}{\frac{12}{13}} \\
 &= \frac{1}{12} \times 2^{-24} \approx 4,9670 \times 10^{-9}
 \end{aligned}$$

## Respuesta 2

Tenemos  $245696,09_{10}$ , procedemos a convertirlo a binario:

■ **Parte Entera:**

$\frac{245696}{2} = 122848, b_{17} = 0$	$\frac{239}{2} = 119, b_{07} = 1$
$\frac{122848}{2} = 61424, b_{16} = 0$	$\frac{119}{2} = 59, b_{06} = 1$
$\frac{61424}{2} = 30712, b_{15} = 0$	$\frac{59}{2} = 29, b_{05} = 1$
$\frac{30712}{2} = 15356, b_{14} = 0$	$\frac{29}{2} = 14, b_{04} = 1$
$\frac{15356}{2} = 7678, b_{13} = 0$	$\frac{14}{2} = 7, b_{03} = 0$
$\frac{7678}{2} = 3839, b_{12} = 0$	$\frac{7}{2} = 3, b_{02} = 1$
$\frac{3839}{2} = 1919, b_{11} = 1$	$\frac{3}{2} = 1, b_{01} = 1$
$\frac{1919}{2} = 959, b_{10} = 1$	$\frac{1}{2} = 0, b_{00} = 1$
$\frac{959}{2} = 479, b_{09} = 1$	
$\frac{479}{2} = 239, b_{08} = 1$	

Por lo tanto tenemos que  $245696_{10} = 11101111111000000_2$ .

■ **Parte Decimal:**

$\frac{9}{100} \times 2 = \frac{18}{100}, b_0 = 0$	$\frac{88}{100} \times 2 = \frac{176}{100}, b_5 = 1$
$\frac{18}{100} \times 2 = \frac{36}{100}, b_1 = 0$	$\frac{76}{100} \times 2 = \frac{152}{100}, b_6 = 1$
$\frac{36}{100} \times 2 = \frac{72}{100}, b_2 = 0$	$\frac{52}{100} \times 2 = \frac{104}{100}, b_7 = 1$
$\frac{72}{100} \times 2 = \frac{144}{100}, b_3 = 1$	
$\frac{44}{100} \times 2 = \frac{88}{100}, b_4 = 0$	

Por lo tanto tenemos que  $0,09_{10} = 00010111_2$ .

Entonces se tiene que  $245696,09_{10} \approx \overbrace{111011111111000000,000101}^{24bits} 11_2$ .

Si usamos **redondeo**:

$$Fl(245696,09)_{Redondeado} = 0,111011111111000000000110 \times 2^{18}$$

Si representamos este numero de vuelta en **decimal**:

$$111011111111000000,000110_2 = 245696,09375_{10}.$$

**Error absoluto y relativo:**

$$\begin{aligned} E_A &= |x - Fl(x)_{Redondeado}| = |245696,09 - 245696,09375| \\ &= -3,75 \times 10^{-3} \\ E_R &= \frac{E_A}{|x|} = \frac{-3,75 \times 10^{-3}}{245696,09} \approx -1,526275815 \times 10^{-8} \end{aligned}$$

## Respuesta 3

Después de correr el programa, se obtuvieron los siguientes datos

- **Para simple precisión:**  $\epsilon = 0,0000001192092895507812500000000000$

Iteracion	t	$\epsilon$
1	1.5000000000000000000000000000000000	0.5000000000000000000000000000000000
2	1.2500000000000000000000000000000000	0.2500000000000000000000000000000000
3	1.1250000000000000000000000000000000	0.1250000000000000000000000000000000
4	1.0625000000000000000000000000000000	0.0625000000000000000000000000000000
5	1.0312500000000000000000000000000000	0.0312500000000000000000000000000000
6	1.0156250000000000000000000000000000	0.0156250000000000000000000000000000
7	1.0078125000000000000000000000000000	0.0078125000000000000000000000000000
8	1.0039062500000000000000000000000000	0.0039062500000000000000000000000000
9	1.0019531250000000000000000000000000	0.0019531250000000000000000000000000
10	1.0009765625000000000000000000000000	0.0009765625000000000000000000000000
11	1.0004882812500000000000000000000000	0.0004882812500000000000000000000000
12	1.0002441406250000000000000000000000	0.0002441406250000000000000000000000
13	1.0001220703125000000000000000000000	0.0001220703125000000000000000000000
14	1.0000610351562500000000000000000000	0.0000610351562500000000000000000000
15	1.0000305175781250000000000000000000	0.0000305175781250000000000000000000
16	1.0000152587890625000000000000000000	0.0000152587890625000000000000000000

17	1.000007629394531200000000	0.0000076293945312500000000000000000
18	1.000003814697265600000000	0.0000038146972656250000000000000000
19	1.000001907348632800000000	0.0000019073486328125000000000000000
20	1.000000953674316400000000	0.0000009536743164062500000000000000
21	1.000000476837158200000000	0.0000004768371582031250000000000000
22	1.000000238418579100000000	0.0000002384185791015625000000000000
23	1.000000119209289600000000	0.0000001192092895507812500000000000
24	1.000000000000000000000000	0.0000000596046447753906250000000000

- Para doble precisión:  $\epsilon = 0,0000000000000002220446049250313100$

Iteracion	t	$\epsilon$
1	1.500000000000000000000000	0.5000000000000000000000000000000000
2	1.250000000000000000000000	0.2500000000000000000000000000000000
3	1.125000000000000000000000	0.1250000000000000000000000000000000
4	1.062500000000000000000000	0.0625000000000000000000000000000000
5	1.031250000000000000000000	0.0312500000000000000000000000000000
6	1.015625000000000000000000	0.0156250000000000000000000000000000
7	1.007812500000000000000000	0.0078125000000000000000000000000000
8	1.003906250000000000000000	0.0039062500000000000000000000000000
9	1.001953125000000000000000	0.0019531250000000000000000000000000
10	1.000976562500000000000000	0.0009765625000000000000000000000000
11	1.000488281250000000000000	0.0004882812500000000000000000000000
12	1.000244140625000000000000	0.0002441406250000000000000000000000
13	1.000122070312500000000000	0.0001220703125000000000000000000000
14	1.000061035156250000000000	0.0000610351562500000000000000000000
15	1.000030517578125000000000	0.0000305175781250000000000000000000
16	1.000015258789062500000000	0.0000152587890625000000000000000000
17	1.000007629394531200000000	0.0000076293945312500000000000000000
18	1.000003814697265600000000	0.0000038146972656250000000000000000
19	1.000001907348632800000000	0.0000019073486328125000000000000000
20	1.000000953674316400000000	0.0000009536743164062500000000000000
21	1.000000476837158200000000	0.0000004768371582031250000000000000
22	1.000000238418579100000000	0.0000002384185791015625000000000000
23	1.000000119209289600000000	0.0000001192092895507812500000000000
24	1.000000059604644800000000	0.0000000596046447753906250000000000
25	1.000000029802322400000000	0.0000000298023223876953130000000000
26	1.000000014901161200000000	0.0000000149011611938476560000000000
27	1.000000007450580600000000	0.0000000074505805969238281000000000

28	1.000000003725290300000000	0.0000000037252902984619141000000000
29	1.000000001862645100000000	0.0000000018626451492309570000000000
30	1.000000000931322600000000	0.0000000009313225746154785200000000
31	1.000000000465661300000000	0.0000000004656612873077392600000000
32	1.000000000232830600000000	0.0000000002328306436538696300000000
33	1.000000000116415300000000	0.0000000001164153218269348100000000
34	1.000000000058207700000000	0.0000000000582076609134674070000000
35	1.000000000029103800000000	0.0000000000291038304567337040000000
36	1.000000000014551900000000	0.0000000000145519152283668520000000
37	1.000000000007276000000000	0.0000000000072759576141834259000000
38	1.000000000003638000000000	0.0000000000036379788070917130000000
39	1.000000000001819000000000	0.0000000000018189894035458565000000
40	1.000000000000909500000000	0.0000000000009094947017729282400000
41	1.000000000000454700000000	0.0000000000004547473508864641200000
42	1.000000000000227400000000	0.0000000000002273736754432320600000
43	1.000000000000113700000000	0.0000000000001136868377216160300000
44	1.000000000000056800000000	0.0000000000000568434188608080150000
45	1.000000000000028400000000	0.0000000000000284217094304040070000
46	1.000000000000014200000000	0.0000000000000142108547152020040000
47	1.000000000000007100000000	0.0000000000000071054273576010019000
48	1.000000000000003600000000	0.0000000000000035527136788005009000
49	1.000000000000001800000000	0.0000000000000017763568394002505000
50	1.000000000000000900000000	0.0000000000000008881784197001252300
51	1.000000000000000400000000	0.0000000000000004440892098500626200
52	1.000000000000000200000000	0.0000000000000002220446049250313100
53	1.000000000000000000000000	0.0000000000000001110223024625156500

El valor de  $\epsilon$  es distinto de  $10^{-308}$ , porque como estamos continuamente sumando 1 con  $\epsilon$ , y  $\epsilon$  se vuelve mas pequeño con cada iteración, y su magnitud es muy pequeña comparada con la de 1 y la suma de  $1 + \epsilon$  deja de ser significativa.

Para precisión simple  $\delta = 0,00097656250$  y para precisión doble  $\delta = 0,0000000000181898940354585650$ .

Los valores de  $\epsilon$  y  $\delta$  son distintos ya que la magnitud de 10000 es mucho mayor a la de 1 y por lo tanto al sumarle números pequeños se llega de manera rápida a uno que no afecte la suma.

---

## Respuesta 4

- Ascendente precisión simple:
- Ascendente precisión doble:
- Descendente precisión simple:

- Descendente precisión doble:
  - Mayor a menor precisión simple:
  - Mayor a menor precisión doble:
  - Menor a mayor precisión simple:
  - Menor a mayor precisión doble:
- 

## Respuesta 5

El mayor valor que llego a tomar la sumatoria fue **15.4036827008740234**. Fueron sumados **2097152 términos** antes de que la computadora dejara de "sumar".

La computadora no llega infinito al realizar la sumatoria debido a la precisión decimal, llega a un punto en que la computadora al sumar dos números, las magnitudes entre ellos son muy distintas y por lo tanto se queda con el numero mas grande y es como si no se le sumara nada.

---

## Respuesta 6

Para  $x = 10$ , con precisión simple tenemos que  $e^{10} = \mathbf{22026,4667968750}$ , este resultado se obtuvo al sumar los términos desde un  $n = 0$ , y hasta que la suma dejara de "sumar", usando al final **32 iteraciones**.

Y para precisión doble tenemos  $e^{10} = \mathbf{22026,465760913433769019320607185363769531250}$ . que se obtuvo al cambiar el orden en que se suma y se empezó desde  $n = 32$  hasta  $n = 0$ . En este caso usamos **32 iteraciones**.

Si se empezara desde un  $n$  muy grande y hasta  $n = 0$ , se tendría un resultado mas preciso.

## Código Fuente

### repuesta3.c

---

```
1  #include <stdio.h>
2
3  void singlePrecision(float);
4  void doublePrecision(double);
5
6  int main(){
7      singlePrecision(1);
8      doublePrecision(1);
9      singlePrecision(10000);
10     doublePrecision(10000);
11 }
12
13 void singlePrecision(float x){
14     float t = 2 * x, epsilon=1;
15     int i=1;
16     while(t > x){
17         t = x + (epsilon *= 0.5);
18         printf("i=%d t=%.24f epsilon=%.34f\n", i, t, epsilon);
19         i++;
20     }
21 }
22
23 void doublePrecision(double x){
24     double t=2 * x, epsilon=1;
25     int i=1;
26     while(t > x){
27         t = x + (epsilon *= 0.5);
28         printf("i=%d t=%.24lf epsilon=%.34lf\n", i, t, epsilon);
29         i++;
30     }
31 }
```

---

## repuesta4.c

```
1  #include <stdio.h>
2
3  float productoEscalarAscendenteF(float[], float[]);
4  float productoEscalarDescendenteF(float[], float[]);
5
6  double productoEscalarAscendenteD(double[], double[]);
7  double productoEscalarDescendenteD(double[], double[]);
8
9  int main() {
10     float a_float[] = {2.718281828, -3.141592654, 1.414213562,
11                        0.5772156649, 0.3010299957};
12
13     float b_float[] = {1485.2497, 878366.9879, -22.37492,
14                        4773714.647, 0.000185049};
15
16     double a_double[] = {2.718281828, -3.141592654, 1.414213562,
17                           0.5772156649, 0.3010299957};
18
19     double b_double[] = {1485.2497, 878366.9879, -22.37492,
20                           4773714.647, 0.000185049};
21
22     printf("Ascendiente Precision Simple = %.34f\n", productoEscalarAscendenteF(a_float, b_float));
23     printf("Descendiente Precision Simple = %.34f\n", productoEscalarDescendenteF(a_float, b_float));
24     printf("----\n");
25     printf("Ascendiente Precision Doble = %.34lf\n", productoEscalarAscendenteD(a_double, b_double));
26     printf("Descendiente Precision Doble = %.34lf\n", productoEscalarDescendenteD(a_double, b_double));
27 }
28
29 // funciones float
30 float productoEscalarAscendenteF(float a[], float b[]){
31     int i;
32     float producto=0;
33     for(i=0; i < 5; i++){
34         //~ printf("%.30f x %.30f\n", a[i], b[i]);
35         producto += a[i] * b[i];
36     }
37     return producto;
38 }
39
40 float productoEscalarDescendenteF(float a[], float b[]){
41     int i;
42     float producto=0;
43     for(i=4; i >= 0; i--){
44         //~ printf("%.30f x %.30f\n", a[i], b[i]);
45         producto += a[i] * b[i];
46     }
47     return producto;
48 }
49
50
51 // funciones double
52 double productoEscalarAscendenteD(double a[], double b[]){
```



```
53     int i;
54     double producto=0;
55     for(i=0; i < 5; i++){
56         producto += a[i] * b[i];
57     }
58     return producto;
59 }
60
61 double productoEscalarDescendenteD(double a[], double b[]){
62     int i;
63     double producto=0;
64     for(i=4; i >= 0; i--){
65         producto += a[i] * b[i];
66     }
67     return producto;
68 }
```

---

## repuesta5.c

---

```
1  #include <stdio.h>
2  #include <math.h>
3
4  void serieArmonicaSingle();
5
6  int main(){
7      serieArmonicaSingle();
8  }
9
10 void serieArmonicaSingle(){
11     float serie=0, ant=1;
12     float k=1;
13
14     while(serie - ant != 0){
15         ant = serie;
16         serie += 1 / (k++);
17         //~ printf("serie: %.10f\n", serie);
18     }
19
20     printf("%f terminos.\n", k);
21     printf("Serie = %.24f\n", serie);
22 }
```

---

## repuesta6.c

---

```
1  #include <stdio.h>
2  #include <math.h>
3
4  float exponencialAdelante(float);
5  float exponencialAtras(float);
6  float fact(float);
7  double exponencialAdelanteD(double);
8  double exponencialAtrasD(double);
9  double factD(double);
10
11 int main(){
12     float x=10;
13     printf("Precision Simple:\n");
14     printf("\tHacia adelante = %.54f\n", exponencialAdelante(x));
15     printf("\tHacia atras    = %.54f\n", exponencialAtras(x));
16     printf("-----\n");
17     printf("Precision Doble:\n");
18     printf("\tHacia adelante = %.54lf\n", exponencialAdelanteD(x));
19     printf("\tHacia atras    = %.54lf\n", exponencialAtrasD(x));
20
21
22 }
23
24 // Funciones float
25 float exponencialAdelante(float x){
26     float e=0, factorial=0, potencia=0, ant=1, i=0;
27
28     while(e != ant){
29         factorial = factorial == 0 ? 1 : factorial * i;
30         potencia = (i++) == 0 ? 1 : potencia * x;
31         ant = e;
32         e += potencia / factorial;
33     }
34
35     return e;
36 }
37
38 float exponencialAtras(float x){
39     float e=0, factorial=0, potencia=0, n=33;
40
41     while(n >= 0){
42         factorial = fact(n);
43         potencia = pow(x, (n--));
44         e += potencia / factorial;
45     }
46
47     return e;
48 }
49
50 float fact(float n){
51     int i;
52     float f=1;
```

```

53         for(i=1; i <= n; i++){
54             f *= i;
55         }
56         return f;
57     }
58
59
60     // Funciones double
61     double exponencialAdelanteD(double x){
62         double e=0, factorial=0, potencia=0, ant=1, i=0;
63
64         while(e != ant){
65             factorial = (factorial == 0 ? 1 : factorial * i);
66             potencia = ((i++) == 0 ? 1 : potencia * x);
67             ant = e;
68             e += potencia / factorial;
69         }
70
71         return e;
72     }
73
74     double exponencialAtrasD(double x){
75         double e=0, factorial=0, potencia=0, n=33;
76
77         while(n >= 0){
78             factorial = fact(n);
79             potencia = pow(x, (n--));
80             e += potencia / factorial;
81         }
82
83         return e;
84     }
85
86     double factD(double n){
87         int i;
88         double f=1;
89         for(i=1; i <= n; i++){
90             f *= i;
91         }
92         return f;
93     }

```

---