

---

**NUMERICAL METHODS AND PROGRAMMING**

2024/2025

**Storage of numbers and algorithms**
**(PROBLEMS 1)**


---

1.— Write in fixed point and floating point, and in the denary, binary, octal and hexadecimal number bases the following numbers:

- a) 5.3125
- b) 0,1
- c) 1/3
- d)  $\pi$
- e)  $e$
- f) 52745916

**Solution 1.a)** We firstly change to binary base to make the conversion to octal and hexadecimal easier. The integer part and the decimal part are translated separately. For the integer part:

$$\begin{array}{l} 5 : 2 = 2 + \underbrace{1} / 2 \\ 2 : 2 = 1 + \underbrace{0} / 2 \\ 1 : 2 = 0 + \underbrace{1} / 2 \end{array} \Rightarrow 5 = (101)_2$$

To change to octal we take the digits in sets of 3 ( $8=2^3$ ) and for hexadecimal in sets of 4 ( $16=2^4$ ), es decir:

$$5 = (\underbrace{101}_3)_2 = (5)_8$$

$$5 = (\underbrace{0101}_4)_2 = (5)_{16}$$

For the decimal part:

$$\begin{array}{l} 0,3125 \times 2 = \underbrace{0},6250 \\ 0,6250 \times 2 = \underbrace{1},2500 \\ 0,2500 \times 2 = \underbrace{0},5000 \\ 0,5000 \times 2 = \underbrace{1},0000 \end{array} \Rightarrow 0,3125 = (0,0101)_2$$

The same way for the octal and hexadecimal bases:

$$0,3125 = (0.\underbrace{010}_3\underbrace{100}_3)_2 = (0,24)_8$$

$$0,3125 = (0.\underbrace{0101}_4)_2 = (0,5)_{16}$$

Thus:

$5,3125$	$=$	$(0,53125) 10^1$	$=$
$(101,0101)_2$	$=$	$(0,1010101)_2 2^3$	$=$
$(5,24)_8$	$=$	$(0,524)_8 8^1$	$=$
$(5,5)_{16}$	$=$	$(0,55)_{16} 16^1$	$=$

**Solution 1.b)**

$$\begin{array}{l}
 0,1 \times 2 = \underbrace{0}_{,2} \\
 0,2 \times 2 = \underbrace{0}_{,4} \\
 0,4 \times 2 = \underbrace{0}_{,8} \\
 0,8 \times 2 = \underbrace{1}_{,6} \\
 0,6 \times 2 = \underbrace{1}_{,2} \\
 0,2 \times 2 = \underbrace{0}_{,4}
 \end{array}
 \Rightarrow 0,1 = (0,0\overline{0011})_2$$

$$0,1 = (0,0\overline{0011})_2 = (0,\underbrace{000}\overline{110}\underbrace{011}\overline{001}\overbrace{100})_2 = (0,0\overline{6314})_8$$

$$0,1 = (0,\underbrace{0001}\overline{1001})_2 = (0,1\overline{9})_{16}$$

Then:

$0,1$	$=$	$(0,1) 10^0$	$=$
$(0,0\overline{0011})_2$	$=$	$(0,\overline{1100})_2 2^{-3}$	$=$
$(0,0\overline{6314})_8$	$=$	$(0,\overline{6314})_8 8^{-1}$	$=$
$(0,1\overline{9})_{16}$	$=$	$(0,1\overline{9})_{16} 16^0$	$=$

**Solution 1.c)**

$$1/3 = 0.\overline{3}$$

$$\begin{array}{l}
 0.\overline{3} \times 2 = \underbrace{0}_{,6} \\
 0.\overline{6} \times 2 = \underbrace{1}_{,3} \\
 0.\overline{3} \times 2 = \underbrace{0}_{,6}
 \end{array}
 \Rightarrow 0.\overline{3} = (0,\overline{01})_2$$

$$0.\overline{3} = (0,\overline{010}\overline{101})_2 = (0,\overline{25})_8$$

$$0.\overline{3} = (0,\overline{0101})_2 = (0,\overline{5})_{16}$$

Thus:

$0.\overline{3}$	$=$	$(0,\overline{3}) 10^0$	$=$
$(0,\overline{01})_2$	$=$	$(0,\overline{10})_2 2^{-1}$	$=$
$(0,\overline{25})_8$	$=$	$(0,\overline{25})_8 8^0$	$=$
$(0,5)_{16}$	$=$	$(0,5)_{16} 16^0$	$=$

**Solution 1.d)**

$$\pi \approx 3,141592653589793$$

$$\begin{array}{l} 3 : 2 = 1 + \underbrace{1} / 2 \\ 1 : 2 = 0 + \underbrace{1} / 2 \end{array} \Rightarrow 3 = (11)_2$$

$$3 = (\underline{011})_2 = (3)_8$$

$$3 = (\underline{0011})_2 = (3)_{16}$$

The conversion of the decimal part is shown in table 1.

$$0,141592653589793... = (0,001001000011111101101010100010001...)_2 =$$

$$(0,\underbrace{001\ 001\ 000\ 011\ 111\ 101\ 101\ 010\ 100\ 010\ 001\ \dots}_2) = (0,11037552421...)_8 =$$

$$(0,\underbrace{0010\ 0100\ 0011\ 1111\ 0110\ 1010\ 1000\ 1000\ 1\dots}_2) = (0,243F6A88...)_{16}$$

Then:

3,141592653589793...	=	(0,3141592653589793...) $10^1$
(11,001001000011111101101010100010001...) $_2$	=	
= (0,11001001000011111101101010100010001...) $_2 2^2$	=	
(3,11037552421...) $_8$	=	(0,311037552421...) $_8 8^1$
(3,243F6A99...) $_{16}$	=	(0,3243F6A99...) $_{16} 16^1$

**Solution 1.e)**

$$e = 2,718281828459045$$

$$\begin{array}{l} 2 : 2 = 1 + \underbrace{0} / 2 \\ 1 : 2 = 0 + \underbrace{1} / 2 \end{array} \Rightarrow 2 = (10)_2$$

$$2 = (\underline{010})_2 = (2)_8$$

$$2 = (\underline{0010})_2 = (2)_{16}$$

The conversion of the decimal part is shown in table 2.

$$0,718281828459045... = (0,101101111110000101010001011000101...)_2 =$$

$$(0,\underbrace{101\ 101\ 111\ 110\ 000\ 101\ 010\ 001\ 011\ 000\ 101\ \dots}_2) = (0,55760521305...)_8 =$$

$$(0,\underbrace{1011\ 0111\ 1110\ 0001\ 0101\ 0001\ 0110\ 0010\ 1\dots}_2) = (0,B7E15162...)_{16}$$

Tabla 1: Decimal part conversion of  $\pi$ .

1	0.141592653589793	0.283185307179586	0
2	0.283185307179586	0.566370614359172	0
3	0.566370614359172	1.132741228718340	1
4	0.132741228718344	0.265482457436688	0
5	0.265482457436688	0.530964914873376	0
6	0.530964914873376	1.061929829746750	1
7	0.061929829746752	0.123859659493505	0
8	0.123859659493505	0.247719318987009	0
9	0.247719318987009	0.495438637974019	0
10	0.495438637974019	0.990877275948037	0
11	0.990877275948037	1.981754551896070	1
12	0.981754551896074	1.963509103792150	1
13	0.963509103792148	1.927018207584300	1
14	0.927018207584297	1.854036415168590	1
15	0.854036415168594	1.708072830337190	1
16	0.708072830337187	1.416145660674370	1
17	0.416145660674374	0.832291321348748	0
18	0.832291321348748	1.664582642697500	1
19	0.664582642697496	1.329165285394990	1
20	0.329165285394993	0.658330570789985	0
21	0.658330570789985	1.316661141579970	1
22	0.316661141579971	0.633322283159941	0
23	0.633322283159941	1.266644566319880	1
24	0.266644566319883	0.533289132639765	0
25	0.533289132639765	1.066578265279530	1
26	0.066578265279532	0.133156530559063	0
27	0.133156530559063	0.266313061118126	0
28	0.266313061118126	0.532626122236251	0
29	0.532626122236251	1.065252244472500	1
30	0.065252244472504	0.130504488945007	0
31	0.130504488945007	0.261008977890015	0
32	0.261008977890015	0.522017955780029	0
33	0.522017955780029	1.044035911560050	1
34	0.044035911560059	0.088071823120117	0
35	0.088071823120117	0.176143646240234	0
36	0.176143646240234	0.352287292480469	0
37	0.352287292480469	0.704574584960937	0
38	0.704574584960937	1.409149169921870	1
39	0.409149169921875	0.818298339843750	0
40	0.818298339843750	1.636596679687500	1

Tabla 2: Decimal part conversion of  $e$ .

1	0.718281828459045	1.436563656918090	1
2	0.436563656918090	0.873127313836180	0
3	0.873127313836180	1.746254627672360	1
4	0.746254627672360	1.492509255344720	1
5	0.492509255344720	0.985018510689439	0
6	0.985018510689439	1.970037021378880	1
7	0.970037021378879	1.940074042757760	1
8	0.940074042757757	1.880148085515510	1
9	0.880148085515515	1.760296171031030	1
10	0.760296171031030	1.520592342062060	1
11	0.520592342062059	1.041184684124120	1
12	0.041184684124119	0.082369368248237	0
13	0.082369368248237	0.164738736496474	0
14	0.164738736496474	0.329477472992949	0
15	0.329477472992949	0.658954945985897	0
16	0.658954945985897	1.317909891971790	1
17	0.317909891971794	0.635819783943589	0
18	0.635819783943589	1.271639567887180	1
19	0.271639567887178	0.543279135774355	0
20	0.543279135774355	1.086558271548710	1
21	0.086558271548711	0.173116543097422	0
22	0.173116543097422	0.346233086194843	0
23	0.346233086194843	0.692466172389686	0
24	0.692466172389686	1.384932344779370	1
25	0.384932344779372	0.769864689558744	0
26	0.769864689558744	1.539729379117480	1
27	0.539729379117488	1.079458758234970	1
28	0.079458758234978	0.158917516469955	0
29	0.158917516469955	0.317835032939911	0
30	0.317835032939911	0.635670065879821	0
31	0.635670065879821	1.271340131759640	1
32	0.271340131759644	0.542680263519287	0
33	0.542680263519287	1.085360527038570	1
34	0.085360527038574	0.170721054077148	0
35	0.170721054077148	0.341442108154297	0
36	0.341442108154297	0.682884216308593	0
37	0.682884216308593	1.365768432617180	1
38	0.365768432617187	0.731536865234375	0
39	0.731536865234375	1.463073730468750	1
40	0.463073730468750	0.926147460937500	0

Thus:

2,718281828459045...	=	(0,2718281828459045...) 10 <sup>1</sup>
(10,101101111110000101010001011000101...)₂	=	
= (0,10101101111110000101010001011000101...)₂ 2 <sup>2</sup>		
(2,55760521305...)₈	=	(0,255760521305...)₈ 8 <sup>1</sup>
(2.B7E15162...)₁₆	=	(0,2B7E15162...)₁₆ 16 <sup>1</sup>

**Solution 1.f)** In this particular case is easier to convert the number to hexadecimal and then to octal and binary:

52745916	:	16	=	3296619	+	<u>12</u>	/	16	
3296619	:	16	=	206038	+	<u>11</u>	/	16	
206038	:	16	=	12877	+	<u>6</u>	/	16	
12877	:	16	=	804	+	<u>13</u>	/	16	⇒
804	:	16	=	50	+	<u>4</u>	/	16	52745916 = (324D6BC)₁₆
50	:	16	=	3	+	<u>2</u>	/	16	
3	:	16	=	0	+	<u>3</u>	/	16	

Converting now to binary and octal:

$$(324D6BC)_{16} = (\underbrace{0011}_B \underbrace{0010}_D \underbrace{0100}_4 \underbrace{1101}_{13} \underbrace{0110}_{10} \underbrace{1011}_{11} \underbrace{1100}_C)_2$$

$$(\underbrace{011}_3 \underbrace{001}_1 \underbrace{001}_1 \underbrace{001}_1 \underbrace{101}_{13} \underbrace{011}_{10} \underbrace{010}_{10} \underbrace{111}_{11} \underbrace{100}_4)_2 = (311153274)_8$$

Therefore:

52745916	=	(0,52745916) 10 <sup>8</sup>
(11001001001101011010111100)₂	=	
= (0,11001001001101011010111100)₂ 2 <sup>26</sup>		
(311153274)₈	=	(0,311153274)₈ 8 <sup>9</sup>
(324D6BC)₁₆	=	(0,324D6BC)₁₆ 16 <sup>7</sup>

2.— The above values are to be stored in the variables A, UD, UT, PI, E and GR of a FORTRAN program:

```

REAL * 4 A, UD, UT, PI, E, GR
A = 5.3125
UD = 0.1
UT = 0.333333
PI = 3.14159
E = 2.71828
GR = 52745916.
...
    
```

If for the mantissa of the numbers variables of type `REAL*4` are used, using 3 bytes (including the sign)

- a) Obtain the values that are actually stored in each case.
- b) Compare the stored values with the corresponding exact values, indicating the sources of error.
- c) Compare the rounding error produced in each case with its upper bound.
- d) Estimate how many significant figures should have been used to write the numbers  $1/3$ ,  $\pi$  and  $e$ . Repeat the estimation for the case in which the variables are to be stored as `REAL*8`.

**Solution 2.a)** We convert the constants to binary, in this case with the decimal digits the statement indicates. We express the numbers in floating point and group them in sets of 1 byte (8 digits):

$$\begin{aligned}
 5,3125 &= (0,10101010|00000000|00000000|0)_2 2^3 \\
 0,1 &= (0,11001100|11001100|11001100|1\dots)_2 2^{-3} \\
 0,333333 &= (0,10101010|10101010|10011111|1\dots)_2 2^{-1} \\
 3,14159 &= (0,11001001|00001111|11001111|1\dots)_2 2^2 \\
 2,71828 &= (0,10101101|11111000|01001100|1\dots)_2 2^2 \\
 52745916, &= (0,11001001|00110101|10101111|0\dots)_2 2^{26}
 \end{aligned}$$

Assuming 23 bits for the mantissa, with 1 for the sign and approximation by rounding, the actually stored values would be:

$$\begin{aligned}
 \text{A} &\rightarrow (0,10101010|00000000|00000000)_2 2^3 \\
 \text{UD} &\rightarrow (0,11001100|11001100|11001100)_2 2^{-3} \\
 \text{UT} &\rightarrow (0,10101010|10101010|10\underline{10000})_2 2^{-1} \\
 \text{PI} &\rightarrow (0,11001001|00001111|110\underline{1000})_2 2^2 \\
 \text{E} &\rightarrow (0,10101101|11111000|0100\underline{110})_2 2^2 \\
 \text{GR} &\rightarrow (0,11001001|00110101|101\underline{1000})_2 2^{26}
 \end{aligned}$$

NOTE: Actually, since the first bit is always a 1 it is not stored, we will not take this into account for the sake of simplicity. The underlined part is the one affected by rounding.

If we convert again to decimal we have the values actually stored:

A	≡	5,3125	=	5,3125
UD	≡	0,0999999940395	≠	0,1
UT	≡	0,333333015442	≠	0,333333
PI	≡	3,14159011841	≠	3,14159
E	≡	2,71827983856	≠	2,71828
GR	≡	52745920.	≠	52745916.

- Solution 2.b)**
- \* 5.3125 is exactly stored in A since its mantissa occupies less than 23 bits.
  - \* 0.1 and 52745916 are stored in UD and GR respectively with a rounding error due to their mantissas occupying more than 23 significant digits. Rounding with 23 bits losses information.

\* 0.333333, 3.14159 and 2.71828 are stored in UT, PI and E respectively with a rounding error, again because their mantissa requires more than 23 significant digits. Also this values have inherent errors since we are actually trying to use the numbers  $1/3$ ,  $\pi$  and  $e$ .

**Solution 2.c)** Let  $\bar{x}$  be the real number we want to represent,  $\hat{x}$  the value we save and  $x$  the value actually stored ( $\bar{x} = \pi$ ,  $\hat{x} = 3,14159$ ,  $x = 3,14159011841$ ). We define:

$$\left\{ \begin{array}{l} \text{error de redondeo} = r_x^A = \frac{\hat{x}-x}{\hat{x}} \\ \text{error inherente} = r_x^I = \frac{\bar{x}-\hat{x}}{\bar{x}} \\ \text{error total} = r_x = \frac{\bar{x}-x}{x} = \left(\frac{\bar{x}-\hat{x}}{\bar{x}}\right) + \left(\frac{\hat{x}-x}{\hat{x}}\right) \left(\frac{\hat{x}}{x}\right) \approx r_x^I + r_x^A \end{array} \right.$$

The upper bound for the storing error is the Machine Epsilon:

$$|r_x^A| \leq \frac{1}{2} 2^{-m+2} \Rightarrow |r_x^A| \leq r_M = 2^{-23} \approx 1,1909 \cdot 10^{-7}$$

Rounding errors:

$$\begin{aligned} r_A^A &= \frac{5,3125-5,3125}{5,3125} = 0 \\ r_{UD}^A &= \frac{0,1-0,0999999940395}{0,1} \approx 0,596050 \cdot 10^{-7} \rightarrow |r_{UD}^A| \approx 0,50 r_M \\ r_{UT}^A &= \frac{0,333333-0,333333015442}{0,333333} \approx -0,463260 \cdot 10^{-7} \rightarrow |r_{UT}^A| \approx 0,39 r_M \\ r_{PI}^A &= \frac{3,14159-3,14159011841}{3,14159} \approx -0,376911 \cdot 10^{-7} \rightarrow |r_{PI}^A| \approx 0,32 r_M \\ r_E^A &= \frac{2,71828-2,71827983856}{2,71828} \approx 0,593905 \cdot 10^{-7} \rightarrow |r_E^A| \approx 0,50 r_M \\ r_{GR}^A &= \frac{52745916.-52745920.}{52745916.} \approx -0,758353 \cdot 10^{-7} \rightarrow |r_{GR}^A| \approx 0,64 r_M \end{aligned}$$

Inherent errors:

$$\begin{aligned} r_A^I &= r_{UD}^I = r_{GR}^I = 0 \\ r_{UT}^I &= \frac{1/3-0,333333}{1/3} \approx 1,000000 \cdot 10^{-6} \rightarrow |r_{UT}^I| \approx 8,4 r_M \\ r_{PI}^I &= \frac{\pi-3,14159}{\pi} \approx 0,844664 \cdot 10^{-6} \rightarrow |r_{PI}^I| \approx 7,1 r_M \\ r_E^I &= \frac{e-2,71828}{e} \approx 0,672653 \cdot 10^{-6} \rightarrow |r_E^I| \approx 5,6 r_M \end{aligned}$$

Global errors:

$$\begin{aligned} r_A &= r_A^A + r_A^I = 0 \\ r_{UD} &= r_{UD}^A + r_{UD}^I \approx 0,596050 \cdot 10^{-7} \rightarrow |r_{UD}| \approx 0,50 r_M \\ r_{UT} &= r_{UT}^A + r_{UT}^I \approx 9,53674 \cdot 10^{-7} \rightarrow |r_{UT}| \approx 8 r_M \\ r_{PI} &= r_{PI}^A + r_{PI}^I \approx 8,069729 \cdot 10^{-7} \rightarrow |r_{PI}| \approx 6,8 r_M \\ r_E &= r_E^A + r_E^I \approx 7,320435 \cdot 10^{-7} \rightarrow |r_E| \approx 6,1 r_M \\ r_{GR} &= r_{GR}^A + r_{GR}^I \approx -0,758353 \cdot 10^{-7} \rightarrow |r_{GR}| \approx 0,64 r_M \end{aligned}$$



**Solution 2.d)** We can see that the inherent errors of UT, PI and E are very high. Thus, we observe that  $1/3$ ,  $\pi$  and  $e$  were introduced with high inherent errors in comparison with the Machine Epsilon  $r_M$ .

We determine the number of significant digits  $s$  we should have used in base 10 in order for the error to be the same order of magnitude of that of the Machine Epsilon:

$$\frac{1}{2} 10^{-(s+1)+2} \approx \frac{1}{2} 2^{-24+2} \Leftrightarrow 10^{-s+1} \approx 2^{-22} \Leftrightarrow s \geq 1 + 22 \log_{10} 2 = 7,62266 \Rightarrow \boxed{s \geq 8}$$

Thus, we should have written:

$$\begin{aligned} \text{UT} &= 0,33333333 \\ \text{PI} &= 3,1415927 \\ \text{E} &= 2,7182818 \end{aligned}$$

In case of using REAL\*8 variables:  $r_M \approx (1/2) 2^{-53+2}$

$$\frac{1}{2} 10^{-(s+1)+2} \approx \frac{1}{2} 2^{-53+2} \Leftrightarrow 10^{-s+1} \approx 2^{-51} \Leftrightarrow s \geq 1 + 51 \log_{10} 2 = 16,3525 \Rightarrow \boxed{s \geq 17}$$

and we should have introduced:

$$\begin{aligned} \text{UT} &= 0,333333333333333333 \\ \text{PI} &= 3,1415926535897932 \\ \text{E} &= 2,7182818284590452 \end{aligned}$$

**3.—** Compare the values stored in the variables GDIF and NGDIF in the following FORTRAN program:

```
REAL * 4 GDIF
INTEGER * 4 NGDIF
GDIF = 52745916. + 15.
GDIF = GDIF - 52745916.
...
NGDIF = 52745916 + 15
NGDIF = NGDIF - 52745916
...IF(GDIF.EQ.FLOAT(NGDIF)) STOP...
```

If for the mantissa of the variables REAL\*4, 3 bytes are used (including the sign), will the program stop when executing the instruction IF? What conclusions can be drawn from this example?

**Solution 3.** For each of the stored numbers in floating point, except for 0, the first bit after the decimal point is always a 1, therefore is not necessary to store it, thus we have an extra bit  $\Rightarrow m=25$ .

Considering this, the number 52745916. is stored as:

$$\text{GR} = (0,11001001|00110101|10101111)_2 2^{26}$$

and number 15. is stored as:

$$(0,11110000|00000000|00000000)_2 2^4$$

then, the statement  $\text{GDIF}=52745916.+15.$  will cast:

$$\begin{array}{r} 11001001|00110101|10101111|00 \\ \phantom{11001001|00110101|} + 11|11 \\ \hline 11001001|00110101|10110010|11 \end{array}$$

that will be stored as:

$$(0,11001001|00110101|10110011)_2 2^{26}$$

Then, the statement  $\text{GDIF}=\text{GDIF}-52745916.$  will cast:

$$\begin{array}{r} 11001001|00110101|10110011|00 \\ -11001001|00110101|10101111|00 \\ \hline 00000000|00000000|00000100|00 \end{array}$$

which will be stored as:

$$(0,10000000|00000000|00000000)_2 2^5 = (16.)_{10}$$

However, the operations with integer values will be more accurate because they can store up to 31 binary digits.

Therefore, the program will not stop since  $\text{GDIF}$  will contain the value 16. and  $\text{FLOAT}(\text{NGDIF})$  will give the value 15.

Conclusion: In a computer is not the same to operate with integer or real variables.

---

4.— A given computer takes a tenth of a second (at most) to:

- a) Obtain, through the equations of minimum least squares, the regression line that best approximates 100 data points.
- b) Calculate the modulus of a vector of 500 elements.
- c) Sort from highest to lowest a list of 50 numbers by the "bubble sort" method.
- d) Obtain the product of a square matrix of size 20 times a vector.
- e) Multiply two square matrices of size 8.
- f) Determine (by algebraic definition) the determinant of a square matrix of size 5.
- g) Find by sequential search a telephone number in a list (alphabetically sorted) of 100 subscribers.
- h) Find by bisection a telephone number in a list (alphabetically sorted) of 100 subscribers.

What is going to be the computational time if the size of the problem (size of the matrix, number of data, for each case) is ten, a hundred, a thousand, ten thousand, a hundred thousand or a million times larger?

---

**Solution 4.** We have to determine the computational complexity of each one of the items to analyze their dependency on the size of the problem.

**Solution 4.a)** Given  $\{x_k, y_k\}_{k=1, \dots, m}$ , we want to approximate  $y(x)$  by a first order polynomial,  $P_1(x) = ax + b$ , so that  $Q = \sum_{k=1}^m (y_k - P_1(x_k))^2$  verifies that  $\frac{\partial Q}{\partial a} = 0$  y  $\frac{\partial Q}{\partial b} = 0$ .

Thus:

$$\left\{ \begin{array}{l} \frac{\partial Q}{\partial a} = -2 \sum_{k=1}^m (y_k - P_1(x_k))x_k \\ \frac{\partial Q}{\partial b} = -2 \sum_{k=1}^m (y_k - P_1(x_k)) \end{array} \right\} \Rightarrow$$

$$\Rightarrow \left\{ \begin{array}{l} -2 \sum_{k=1}^m (y_k - (ax_k + b))x_k = 0 \\ -2 \sum_{k=1}^m (y_k - (ax_k + b)) = 0 \end{array} \right\} \Rightarrow \begin{bmatrix} m & \sum_{k=1}^m x_k \\ SIM & \sum_{k=1}^m x_k^2 \end{bmatrix} \begin{bmatrix} b \\ a \end{bmatrix} = \begin{bmatrix} \sum_{k=1}^m y_k \\ \sum_{k=1}^m x_k y_k \end{bmatrix}$$

Then:

$$D = \det \begin{bmatrix} m & \sum_{k=1}^m x_k \\ SIM & \sum_{k=1}^m x_k^2 \end{bmatrix} = m \sum_{k=1}^m x_k^2 - \left( \sum_{k=1}^m x_k \right)^2$$

$$b = \det \begin{bmatrix} \sum_{k=1}^m y_k & \sum_{k=1}^m x_k \\ \sum_{k=1}^m x_k y_k & \sum_{k=1}^m x_k^2 \end{bmatrix} / D = \left( \sum_{k=1}^m y_k \sum_{k=1}^m x_k^2 - \sum_{k=1}^m x_k y_k \sum_{k=1}^m x_k \right) / D$$

$$a = \det \begin{bmatrix} m & \sum_{k=1}^m y_k \\ \sum_{k=1}^m x_k & \sum_{k=1}^m x_k y_k \end{bmatrix} / D = \left( m \sum_{k=1}^m x_k y_k - \sum_{k=1}^m x_k \sum_{k=1}^m y_k \right) / D$$

If we operate as follows:

$$\left\{ \begin{array}{l}
 \mu_x = \left( \sum_{k=1}^m x_k \right) / m \quad \rightarrow \quad (m - 1) \text{ sums, 1 division} \\
 \mu_y = \left( \sum_{k=1}^m y_k \right) / m \quad \rightarrow \quad (m - 1) \text{ sums, 1 division} \\
 \sigma_x^2 = \left( \sum_{k=1}^m (x_k - \mu_x)^2 \right) / m \quad \rightarrow \quad (m) \text{ subtractions, } (m) \text{ products, } (m - 1) \text{ sums, 1 division} \\
 \sigma_{xy} = \left( \sum_{k=1}^m (x_k - \mu_x)(y_k - \mu_y) \right) / m \quad \rightarrow \quad (2m) \text{ subtractions, } (m) \text{ products, } (m - 1) \text{ sums, 1 division} \\
 \hline
 P_1(x) = \mu_y + \frac{\sigma_{xy}}{\sigma_x^2} (x - \mu_x) \quad \text{TOTAL} = (9m) \text{ operations}
 \end{array} \right.$$

Therefore the computational cost will be propotional to  $\boxed{\mathcal{O}(m)}$

**Solution 4.b)** Given a vector  $\mathbf{v} = \{v_i\}_{i=1,\dots,m}$ :

$$|\mathbf{v}| = \left( \sum_{i=1}^m v_i^2 \right)^{1/2} \Rightarrow \left\{ \begin{array}{l} m \text{ products} \\ (m - 1) \text{ sums} \\ 1 \text{ square root} \end{array} \right\} \Rightarrow \text{TOTAL} = 2m \text{ operations} \Rightarrow \boxed{\mathcal{O}(m)}$$

**Solution 4.c)**

$$\left\{ \begin{array}{l}
 \text{DO I} = 1, M - 1 \\
 \quad \text{DO J} = M, I + 1, -1 \\
 \quad \quad \text{IF(L(I).LT.L(J))THEN} \\
 \quad \quad \quad \text{LL} = \text{L(I)} \\
 \quad \quad \quad \text{L(I)} = \text{L(J)} \\
 \quad \quad \quad \text{L(J)} = \text{LL} \\
 \quad \quad \text{ENDIF} \\
 \quad \text{ENDDO} \\
 \text{ENDDO}
 \end{array} \right\} \Rightarrow \underbrace{[(m - 1) + (m - 2) + \dots + 1]}_{\frac{m(m - 1)}{2}} \text{ comparisons}$$

Not every operation requires a change because it depends on the initial order of the data. If they are already sorted they need 0 changes. If every single one would be necessary it would take  $\frac{m(m - 1)}{2}$  changes.

In any way, the computational complexity would be proportional to  $\frac{m(m - 1)}{2}$  in general, thus

$$\boxed{\mathcal{O}(m^2)}.$$

**Solution 4.d)**

$$\mathbf{w} = \mathbf{A}\mathbf{v} \Rightarrow w_i = \sum_{j=1}^m a_{i,j}v_j \quad ; \quad i = 1, \dots, m \quad \Rightarrow \quad m \text{ veces } \left\{ \begin{array}{l} m \text{ productos} \\ (m - 1) \text{ sumas} \end{array} \right.$$

$$\text{TOTAL} = m(2m - 1) \text{ operations} \quad \Rightarrow \quad \boxed{\mathcal{O}(m^2)}$$

**Solution 4.e)**

$$\mathbf{C} = \mathbf{AB} \Rightarrow c_{i,j} = \sum_{k=1}^m a_{i,k} b_{k,j} \quad ; \quad \begin{matrix} i = 1, \dots, m \\ j = 1, \dots, m \end{matrix} \Rightarrow m^2 \text{ times } \begin{cases} m \text{ products} \\ (m - 1) \text{ sums} \end{cases}$$

$$\text{TOTAL} = m^2(2m - 1) \text{ operations} \quad \Rightarrow \quad \boxed{\mathcal{O}(m^3)}$$

**Solution 4.f)**

$$\det(\mathbf{A}) = \sum_{j=1}^m a_{i,k} \mathbf{A}_{k,j} \Rightarrow 1 \text{ determinant of order } m \Rightarrow \begin{cases} m \text{ products} \\ (m - 1) \text{ sums} \\ m \text{ determinants of order } (m-1) \end{cases}$$

Let  $C_m$  be the number operations needed to obtain the determinant of order  $m$ . Then we know that:

$$C_m = m + (m - 1) + m C_{m-1} = -1 + m(2 + C_{m-1}) \quad \text{con} \quad C_1 = 0$$

Therefore:

$$\begin{aligned} C_1 &= -1 + 1 \\ C_2 &= -1 + 2(2 + C_1) = -1 + 2(1 + 1) \\ C_3 &= -1 + 3(2 + C_2) = -1 + 3(1 + 2(1 + 1)) \\ C_4 &= -1 + 4(2 + C_3) = -1 + 4(1 + 3(1 + 2(1 + 1))) \\ &\dots \\ C_m &= -1 + m(2 + C_{m-1}) = -1 + m(1 + (m - 1)(\dots(1 + 4(1 + 3(1 + 2(1 + 1)))))) \\ &= -1 + m + m(m - 1) + \dots + m(m - 1)\dots, 4 + m(m - 1)\dots, 4 3 + m(m - 1)\dots, 4 3 2 2 \\ &= -1 + \frac{m!}{(m - 1)!} + \frac{m!}{(m - 2)!} + \dots + \frac{m!}{3!} + \frac{m!}{2!} + \frac{m!}{1!} 2 \\ &= -1 + m! \left( \underbrace{1 + \frac{1}{1!} + \frac{1}{2!} + \dots + \frac{1}{(m - 1)!}}_{\rightarrow e \text{ when } m \rightarrow \infty} \right) \\ &\approx -1 + e m! \end{aligned}$$

Thus, the determinant of a matrix of order  $m$  has a cost of  $(-1 + e m!)$  operations, that's is

$$\boxed{\mathcal{O}(m!)}$$

**Solution 4.g)** We define:

$\mathbf{A}(\mathbf{M})$  as the array containing the surnames.

$\mathbf{T}(\mathbf{M})$  as the array containing the corresponding phone numbers.

$\mathbf{AA}$  as the surname whose phone number we are searching for.

Then:

$$\left\{ \begin{array}{l} \text{DO I = 1, M} \\ \quad \text{IF(A(I).EQ.AA) THEN} \\ \quad \quad \text{WRITE(6, *) 'THE PHONE NUMBER IS', T(I)} \\ \quad \quad \text{STOP} \\ \quad \text{ENDIF} \\ \text{ENDDO} \\ \text{WRITE(6, *) 'NON-EXISTENT PHONE NUMBER'} \end{array} \right.$$

The number of comparisons depends on the position of the number we are searching for in the file. Best case scenario only one comparison would be needed. Worst case,  $m$  comparisons. In the average case  $m/2$  comparisons.

Thus, the computational cost will be  $\boxed{= (m)}$

**Solution 4.h)** The program will be:

```

IF(A(1).EQ.AA) THEN
  WRITE(6, *) 'THE PHONE NUMBER IS', T(1)
  STOP
ELSE IF(A(M).EQ.AA) THEN
  WRITE(6, *) 'THE PHONE NUMBER IS', T(M)
  STOP
ELSE IF(A(1).LT.AA).AND.(AA.LT.A(M)) THEN
  IMIN = 1
  IMAX = M
  DO WHILE((IMAX - IMIN).GT.1)
    IMED = (IMIN + IMAX)/2
    IF(A(IMED).EQ.AA) THEN
      WRITE(6, *) 'THE PHONE NUMBER IS', T(IMED)
      STOP
    ELSE IF(A(IMED).LT.AA) THEN
      IMIN = IMED
    ELSE
      IMAX = IMED
    ENDIF
  ENDDO
ENDIF
WRITE(6, *) 'NON-EXISTENT PHONE NUMBER'

```

The number of comparisons to make will be, at maximum, approximately  $\log_2 m$ , thus the cost is  $\boxed{\mathcal{O}(\log_2 m)}$

With every computational cost determined and the data of the statement::

$$\left\{ \begin{array}{l} T_a \approx K_a m_a \quad \rightarrow \quad K_a \approx 0,1s/100 \\ T_b \approx K_b m_b \quad \rightarrow \quad K_b \approx 0,1s/500 \\ T_c \approx K_c m_c^2 \quad \rightarrow \quad K_c \approx 0,1s/50^2 \\ T_d \approx K_d m_d^2 \quad \rightarrow \quad K_d \approx 0,1s/20^2 \\ T_e \approx K_e m_e^3 \quad \rightarrow \quad K_e \approx 0,1s/8^3 \\ T_f \approx K_f m_f! \quad \rightarrow \quad K_f \approx 0,1s/5! \\ T_g \approx K_g m_g \quad \rightarrow \quad K_g \approx 0,1s/100 \\ T_h \approx K_h \log_2 m_h \quad \rightarrow \quad K_h \approx 0,1s/\log_2 100 \end{array} \right.$$

If the problem is  $n$  times larger:

$$\left\{ \begin{array}{l} m_a = 100 n \quad \rightarrow \quad T_a \approx n 0,1s \\ m_b = 500 n \quad \rightarrow \quad T_b \approx n 0,1s \\ m_c = 50 n \quad \rightarrow \quad T_c \approx n^2 0,1s \\ m_d = 20 n \quad \rightarrow \quad T_d \approx n^2 0,1s \\ m_e = 8 n \quad \rightarrow \quad T_e \approx n^3 0,1s \\ m_f = 5 n \quad \rightarrow \quad T_f \approx \frac{(5n)!}{5!} 0,1s \\ m_g = 100 n \quad \rightarrow \quad T_g \approx n 0,1s \\ m_h = 100 n \quad \rightarrow \quad T_h \approx \left( 1 + \frac{\log_2 n}{\log_2 100} \right) 0,1s \end{array} \right.$$

For  $n=10,100,1000,10000,100000,1000000$

$n$	$\frac{m_a}{T_a}$	$\frac{m_a}{T_a}$	$\frac{m_a}{T_a}$	$\frac{m_a}{T_a}$	$\frac{m_a}{T_a}$	$\frac{m_a}{T_a}$	$\frac{m_a}{T_a}$	$\frac{m_a}{T_a}$
1	$\frac{10^2}{0,1s}$	$\frac{5 \cdot 10^2}{0,1s}$	$\frac{5 \cdot 10^1}{0,1s}$	$\frac{2 \cdot 10^1}{0,1s}$	$\frac{8}{0,1s}$	$\frac{5}{0,1s}$	$\frac{10^2}{0,1s}$	$\frac{10^2}{0,1}$
10	$\frac{10^3}{1s}$	$\frac{5 \cdot 10^3}{1s}$	$\frac{5 \cdot 10^2}{10s}$	$\frac{2 \cdot 10^2}{10s}$	$\frac{8 \cdot 10^1}{1,67min}$	$\frac{5 \cdot 10^1}{8 \cdot 10^{53}years}$	$\frac{10^3}{1s}$	$\frac{10^2}{0,15s}$
$10^2$	$\frac{10^4}{10s}$	$\frac{5 \cdot 10^4}{10s}$	$\frac{5 \cdot 10^3}{16,7min}$	$\frac{2 \cdot 10^3}{16,7min}$	$\frac{8 \cdot 10^2}{1,16days}$	---	$\frac{10^4}{10s}$	$\frac{10^3}{10s}$
$10^3$	$\frac{10^5}{1,67min}$	$\frac{5 \cdot 10^5}{1,67min}$	$\frac{5 \cdot 10^4}{1,16días}$	$\frac{2 \cdot 10^4}{1,16days}$	$\frac{8 \cdot 10^3}{3,17years}$	---	$\frac{10^5}{1,67min}$	$\frac{10^4}{0,25s}$
$10^4$	$\frac{10^6}{16,7min}$	$\frac{5 \cdot 10^6}{16,7min}$	$\frac{5 \cdot 10^5}{0,317years}$	$\frac{2 \cdot 10^5}{0,317years}$	$\frac{8 \cdot 10^4}{3170years}$	---	$\frac{10^6}{16,7min}$	$\frac{10^5}{0,3s}$
$10^5$	$\frac{10^7}{2,78h}$	$\frac{5 \cdot 10^7}{2,78h}$	$\frac{5 \cdot 10^6}{31,7years}$	$\frac{2 \cdot 10^6}{31,7years}$	$\frac{8 \cdot 10^5}{3,17 \cdot 10^6years}$	---	$\frac{10^7}{2,78h}$	$\frac{10^6}{0,35s}$
$10^6$	$\frac{10^8}{1,16days}$	$\frac{5 \cdot 10^8}{1,16days}$	$\frac{5 \cdot 10^7}{3170years}$	$\frac{2 \cdot 10^7}{3170years}$	$\frac{8 \cdot 10^6}{3,17 \cdot 10^9years}$	---	$\frac{10^8}{1,16days}$	$\frac{10^7}{0,4s}$

5.— In order to obtain numerically the value of  $\alpha = \sqrt{2}$  two iterative algorithms are proposed:

$$1) x_{n+1} = \frac{2 + x_n(10 - x_n)}{10}; \quad 2) x_{n+1} = \frac{x_n}{2} + \frac{1}{x_n}$$

For both of them:

- Analyze the evolution of the absolute truncation error between two consecutive iterations.
- Determine for which initial values  $x_0$  does the algorithm converge and for which ones it does not. Provide examples.
- Simplificar el estudio anterior suponiendo que la aproximación inicial  $x_0$  es "suficientemente buena". Obtener el orden de convergencia.
- Sabiendo que el valor  $\sqrt{2}$  se encuentra comprendido entre 1.41 y 1.42, estimar cuantas iteraciones hay que realizar partiendo de la aproximación inicial  $x_0=1.415$  para obtener una mejor aproximación con 5, 10, 15, 20 y 100 cifras significativas exactas. Realizar algunas iteraciones para verificar el resultado.

### Solution 5.a)

$$1) \quad x_{n+1} = \frac{2 + x_n(10 - x_n)}{10}$$

Let  $e_n$  be the error at iteration  $n$ , we can write:

$$\left. \begin{aligned} e_n &= \alpha - x_n \rightarrow x_n = \alpha - e_n \\ e_{n+1} &= \alpha - x_{n+1} \rightarrow x_{n+1} = \alpha - e_{n+1} \end{aligned} \right\}$$

Taking it to the expression of the method:

$$\begin{aligned} \alpha - e_{n+1} &= \frac{2 + (\alpha - e_n)(10 - (\alpha - e_n))}{10} = \frac{2 + 10(\alpha - e_n) - (\alpha - e_n)^2}{10} \Rightarrow \\ &\Rightarrow 10\alpha - 10e_{n+1} = 2 + 10\alpha - 10e_n - \alpha^2 - e_n^2 + 2\alpha e_n \end{aligned}$$

We also know that  $\alpha^2 = 2$ , thus:

$$e_{n+1} = \left(1 - \frac{2\alpha}{10} + \frac{e_n}{10}\right) e_n = \underbrace{\left(1 - \frac{\sqrt{2}}{5}\right) e_n + \frac{1}{10} e_n^2}_{\text{Asintotically linear}}$$

$$2) \quad x_{n+1} = \frac{x_n}{2} + \frac{1}{x_n}$$

Using again the expressions:

$$\left. \begin{aligned} e_n &= \alpha - x_n \rightarrow x_n = \alpha - e_n \\ e_{n+1} &= \alpha - x_{n+1} \rightarrow x_{n+1} = \alpha - e_{n+1} \end{aligned} \right\}$$

and taking them to the expression of the method:



$$\begin{aligned} \alpha - e_{n+1} &= \frac{\alpha - e_n}{2} + \frac{1}{\alpha - e_n} \Rightarrow e_{n+1} = \alpha - \frac{\alpha}{2} + \frac{e_n}{2} - \frac{1}{\alpha - e_n} = \\ &= \frac{\alpha + e_n}{2} - \frac{1}{\alpha - e_n} = \frac{(\alpha^2 - e_n^2) - 2}{2(\alpha - e_n)} \end{aligned}$$

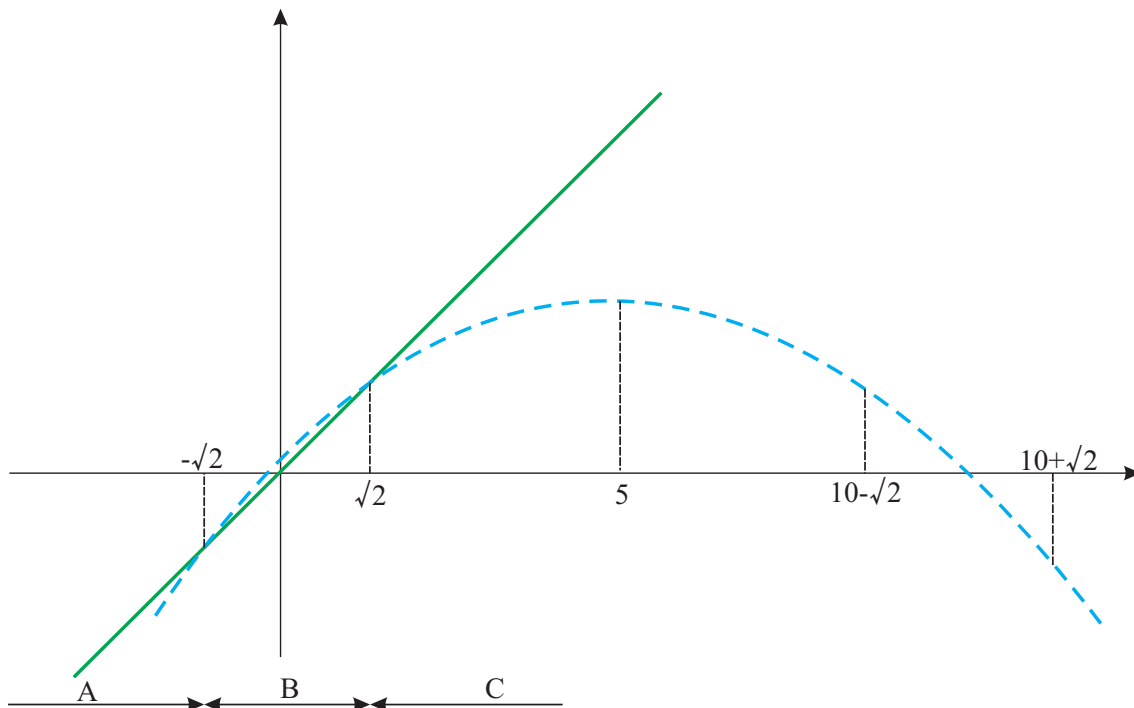
If we already know that  $\alpha^2 = 2$ :

$$e_{n+1} = -\frac{1}{2(\alpha - e_n)}e_n^2 = \underbrace{-\frac{1}{2(\sqrt{2} - e_n)}e_n^2}_{\text{Asintotically quadratic}}$$

**Solution 5.b)**

1)  $x_{n+1} = F(x_n)$  ,  $F(x) = \frac{2 + x(10 - x)}{10}$

Fixed points:  $x = F(x) \Leftrightarrow 10x = 2 + 10x - x^2 \Leftrightarrow x = \pm\sqrt{2}$

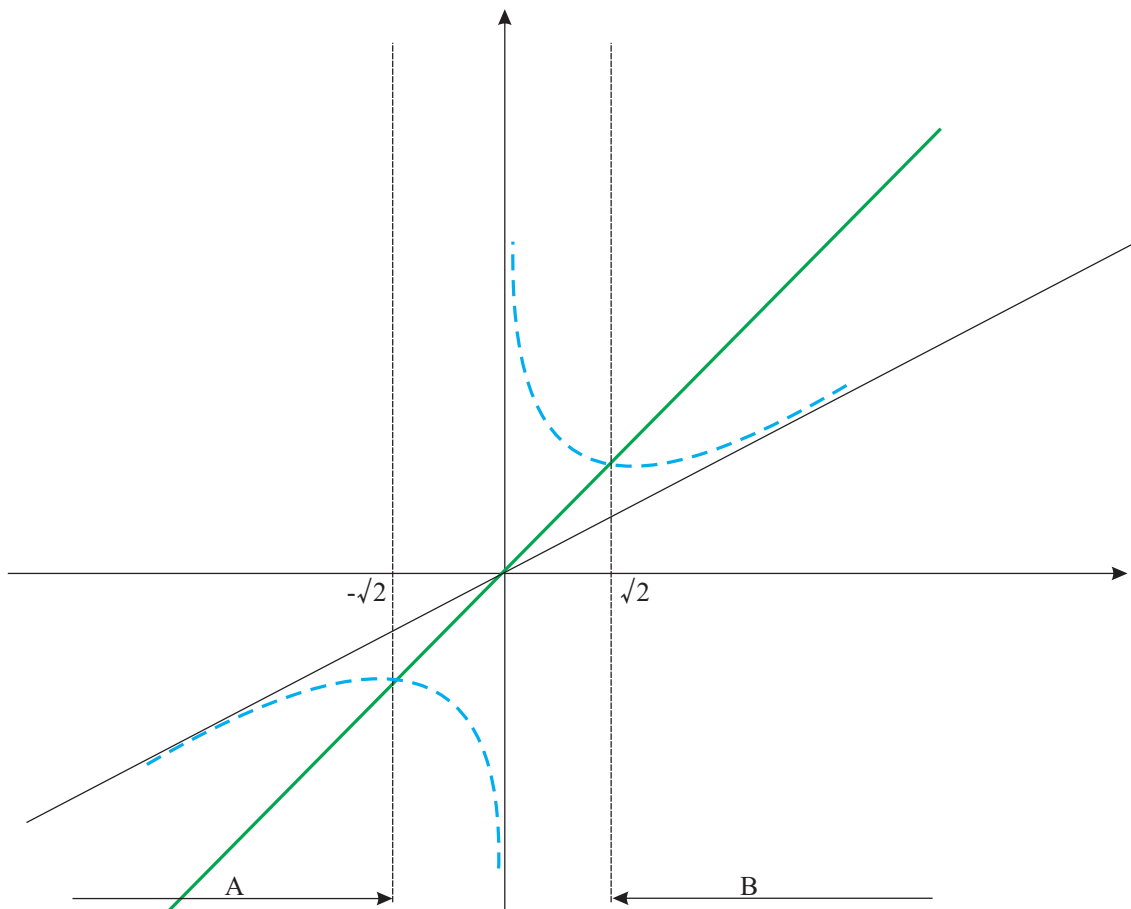


{	Area A : $x_0 < -\sqrt{2}$	$\rightarrow$ diverges to $-\infty$
	---- : $x_0 = -\sqrt{2}$	$\rightarrow x_n = -\sqrt{2} \quad \forall n$
	Area B : $x_0 \in (-\sqrt{2}, \sqrt{2})$	$\rightarrow$ converges to $+\sqrt{2}$ (increasing), F.A.C = $\left 1 - \frac{\sqrt{2}}{5}\right  \approx 0,717$
	---- : $x_0 = \sqrt{2}$	$\rightarrow x_n = \sqrt{2} \quad \forall n$
	Area C : $x_0 \in (\sqrt{2}, 5)$	$\rightarrow$ converges to $+\sqrt{2}$ (decreasing), F.A.C = $\left 1 - \frac{\sqrt{2}}{5}\right  \approx 0,717$
	---- : $x_0 \in [5, 10 - \sqrt{2})$	$\rightarrow x_1 \in \text{Area C}$
	---- : $x_0 = 10 - \sqrt{2}$	$\rightarrow x_n = \sqrt{2} \quad \forall n \geq 1$
	---- : $x_0 \in (10 - \sqrt{2}, 10 + \sqrt{2})$	$\rightarrow x_1 \in \text{Area B}$
	---- : $x_0 = 10 + \sqrt{2}$	$\rightarrow x_n = -\sqrt{2} \quad \forall n \geq 1$
	---- : $x_0 > 10 + \sqrt{2}$	$\rightarrow x_1 \in \text{Area A}$

Thus, it converges to  $\sqrt{2}$  for  $x_0 \in (-\sqrt{2}, 10 + \sqrt{2})$

2)  $x_{n+1} = F(x_n) \quad , \quad F(x) = \frac{x}{2} + \frac{1}{x}$

Fixed points:  $x = F(x) \Leftrightarrow x = \frac{x}{2} + \frac{1}{x} \Leftrightarrow x = \pm\sqrt{2}$



$$\left\{ \begin{array}{ll} \text{Area A} : x_0 < -\sqrt{2} & \rightarrow \text{converges to } -\sqrt{2} \text{ (increasing), quadratic} \\ \text{---} : x_0 = -\sqrt{2} & \rightarrow x_n = -\sqrt{2} \quad \forall n \\ \text{---} : x_0 \in (-\sqrt{2}, 0) & \rightarrow x_1 \in \text{Area A} \\ \text{---} : x_0 = 0 & \rightarrow x_1 \text{ undefined} \\ \text{---} : x_0 \in (0, \sqrt{2}) & \rightarrow x_1 \in \text{Area B} \\ \text{---} : x_0 = \sqrt{2} & \rightarrow x_n = \sqrt{2} \quad \forall n \\ \text{Area B} : x_0 > \sqrt{2} & \rightarrow \text{converges to } \sqrt{2} \text{ (increasing), quadratic} \end{array} \right.$$

Thus, converges to  $-\sqrt{2}$  for  $x_0 \in (-\infty, 0)$  and to  $\sqrt{2}$  for  $x_0 \in (0, +\infty)$ .

**Solution 5.c)** If the initial guess is "good enough":

$$|e_0| \leq \delta \ll 1 \quad \Rightarrow \quad |e_n| \leq \delta \ll 1 \quad \forall n$$

$$\begin{aligned} \mathbf{1)} \quad e_{n+1} &= \left(1 - \frac{\sqrt{2}}{5}\right) e_n + \frac{1}{10} e_n^2 \\ |e_{n+1}| &\leq \left[ \left|1 - \frac{\sqrt{2}}{5}\right| + \frac{1}{10} |e_n| \right] |e_n| \\ |e_{n+1}| &\leq \left[ \left|1 - \frac{\sqrt{2}}{5}\right| + \delta \right] |e_n| \end{aligned}$$

Since we know that  $|e_n| = |x_n - \alpha|$ :

$$|x_{n+1} - \alpha| \leq \underbrace{\left[ \left|1 - \frac{\sqrt{2}}{5}\right| + \delta \right]}_{\lambda_1 \leq 0,717 + \delta < 1} |x_n - \alpha| \Rightarrow$$

$\Rightarrow$  Linear convergence (at least)

$$\begin{aligned} \mathbf{2)} \quad e_{n+1} &= -\frac{1}{2(\sqrt{2} - e_n)} e_n^2 \\ |e_{n+1}| &\leq \left| \frac{1}{2(\sqrt{2} - \delta)} \right| |e_n|^2 \end{aligned}$$

Applying again  $|e_n| = |x_n - \alpha|$ :

$$|x_{n+1} - \alpha| \leq \underbrace{\left| \frac{1}{2(\sqrt{2} - \delta)} \right|}_{\lambda_2} |x_n - \alpha|^2 \Rightarrow$$

$\Rightarrow$  quadratic convergence subject to  $\lambda_2 |x_0 - \alpha|^{2-1} < 1 \Leftrightarrow |x_0 - \alpha| < 2(\sqrt{2} - \delta)$ , which is verified.

**Solution 5.d)**

$$\left. \begin{array}{l} \alpha \in [1,41, 1,42] \\ x_0 = 1,415 \end{array} \right\} \Rightarrow |x_0 - \alpha| \leq 0,005 \quad \text{con } \alpha = \sqrt{2}$$

We want  $x_n \approx \alpha$  with  $d$  exact digits, not including sign, thus:

$$\left| \frac{\alpha - x_n}{\alpha} \right| \leq \frac{1}{2} 10^{-(d+1)+2} = \frac{1}{2} 10^{-d+1}$$

$$\begin{aligned} 1) \quad e_{n+1} &\approx \left(1 - \frac{\sqrt{2}}{5}\right) e_n \Rightarrow e_n \approx \left(1 - \frac{\sqrt{2}}{5}\right)^n e_0 \\ &\Rightarrow |\alpha - x_n| \approx \left(1 - \frac{\sqrt{2}}{5}\right)^n |\alpha - x_0| \leq 0,005 \left(1 - \frac{\sqrt{2}}{5}\right)^n \\ &\Rightarrow \left| \frac{\alpha - x_n}{\alpha} \right| \leq \frac{0,005}{\sqrt{2}} \left(1 - \frac{\sqrt{2}}{5}\right)^n \end{aligned}$$

If  $\frac{0,005}{\sqrt{2}} \left(1 - \frac{\sqrt{2}}{5}\right)^n \leq \frac{1}{2} 10^{-d+1}$  then  $\left| \frac{\alpha - x_n}{\alpha} \right| \leq \frac{1}{2} 10^{-d+1}$

Thus:

$$\left(1 - \frac{\sqrt{2}}{5}\right)^n \leq \frac{\sqrt{2}}{0,005} \frac{1}{2} 10^{-d+1} \Leftrightarrow \underbrace{n \log_{10} \left(1 - \frac{\sqrt{2}}{5}\right)}_{<0} \leq \log_{10} \left(\frac{\sqrt{2}}{0,005} \frac{1}{2}\right) + (-d+1) \Leftrightarrow$$

$$\Leftrightarrow d - 1 - \log_{10} \left(\frac{\sqrt{2}}{0,005} \frac{1}{2}\right) \leq \underbrace{n \left[-\log_{10} \left(1 - \frac{\sqrt{2}}{5}\right)\right]}_{>0} \Leftrightarrow$$

$$\Leftrightarrow n \geq \frac{1}{-\log_{10} \left(1 - \frac{\sqrt{2}}{5}\right)} d + \frac{-1 - \log_{10} \left(\frac{\sqrt{2}}{0,005} \frac{1}{2}\right)}{-\log_{10} \left(1 - \frac{\sqrt{2}}{5}\right)} \approx 6,926 d - 21,820 \Leftrightarrow$$

$$\Leftrightarrow n_1 \geq 6,926 d - 21,820$$

$$\begin{aligned} 2) \quad e_{n+1} &\approx \left(-\frac{1}{2\sqrt{2}}\right) e_n^2 \Rightarrow e_n \approx -\left(\frac{1}{2\sqrt{2}}\right)^{1+2+2^2+\dots+2^{n-1}} e_0^{2^n} = -\left(\frac{1}{2\sqrt{2}}\right)^{2^n-1} e_0^{(2^n)} \\ &\Rightarrow |\alpha - x_n| \approx \left(\frac{1}{2\sqrt{2}}\right)^{2^n-1} 0,005^{(2^n)} \\ &\Rightarrow \left| \frac{\alpha - x_n}{\alpha} \right| \leq \frac{2\sqrt{2}}{\sqrt{2}} \left(\frac{0,005}{2\sqrt{2}}\right)^{2^n} \end{aligned}$$

If  $2 \left( \frac{0,005}{2\sqrt{2}} \right)^{2^n} \leq \frac{1}{2} 10^{-d+1}$  then  $\left| \frac{\alpha - x_n}{\alpha} \right| \leq \frac{1}{2} 10^{-d+1}$

Thus:

$$\begin{aligned} \left( \frac{0,005}{2\sqrt{2}} \right)^{2^n} \leq \frac{1}{4} 10^{-d+1} &\Leftrightarrow 2^n \underbrace{\log_{10} \left( \frac{0,005}{2\sqrt{2}} \right)}_{<0} \leq \log_{10} \left( \frac{1}{4} \right) + (-d + 1) \Leftrightarrow \\ &\Leftrightarrow d - 1 - \log_{10} \left( \frac{1}{4} \right) \leq 2^n \underbrace{\left[ -\log_{10} \left( \frac{0,005}{2\sqrt{2}} \right) \right]}_{>0} \Leftrightarrow \\ &\Leftrightarrow 2^n \geq \frac{d}{-\log_{10} \left( \frac{0,005}{2\sqrt{2}} \right)} + \frac{-1 - \log_{10} \left( \frac{1}{4} \right)}{-\log_{10} \left( \frac{0,005}{2\sqrt{2}} \right)} \approx 0,363d - 0,145 \Leftrightarrow \\ &\Leftrightarrow n_2 \geq \frac{\log_{10}(0,363d - 0,145)}{\log_{10}(2)} \end{aligned}$$

The table shows an estimation of the number of iterations needed for each method to reach the desired significant digits:

d	n <sub>1</sub>	n <sub>2</sub>
5	13	1
10	48	2
15	83	3
20	117	3
100	671	6

Example:

i	x <sub>i</sub> (1)	x <sub>i</sub> (2)
0	<u>1.415</u>	<u>1.415</u>
1	<u>1.4147775</u>	<u>1.41421378092</u>
2	<u>1.41461796255</u>	<u>1.41421356237</u>
3	<u>1.41450356455</u>	<u>1.41421356237</u>
4	<u>1.41442153114</u>	
5	<u>1.41436270436</u>	
6	<u>1.41432051841</u>	
7	<u>1.41429026553</u>	
8	<u>1.41426857001</u>	
9	<u>1.4142530112</u>	
10	<u>1.41424185323</u>	
11	<u>1.41423385129</u>	
12	<u>1.41422811268</u>	
13	<u>1.41422398721</u>	
14	<u>1.41422104578</u>	
15	<u>1.41421892915</u>	

- 
- 6.— An engineer working on road layout is creating a FORTRAN program. One part of the program has to print a list of the cut and fill areas along the road axis. By default the listing starts at kilometer point (PK) 0,000 and a data is printed every 0,1 Km. The program works apparently well, but users complain that sometimes strange values appear in the PK column, such as 70,999 instead of 71,000.
- Make a FORTRAN program that repeatedly adds the value 0,1 in a variable of type `REAL*4` (initialized to zero) and prints the results. Verify that after a certain number of operations the printed result is not a multiple of 0,1. Repeat the calculations with variables of type `REAL*8`, and check that the effect persists, although it takes longer to occur.
  - Check that the same effect occurs in a spreadsheet. To do this, start a column with the value zero in the first row and generate the values of each of the following rows by adding the value 0,1 to the result of the previous row.
  - Propose an alternative to solve the problem satisfactorily.
- 

**Solution 6.a)** The program with variables `REAL*4` is:

```

c ===== Program Iterativesum_R4
c =====
c This program sums iteratively in single precision the number a = 0.1
c until the rounding error affects the third decimal.
c =====

      implicit real * 4(a - h, o - z)
      character * 20 spk
      logical next

      npk = 0
      dpk = 0.1
      xpk = 0.0
      next = .TRUE.

      do while(next)
         npk = npk + 1
         xpk = xpk + dpk
         write(spik,'(f20.3)') xpk
         write(6,100) xpk
         next = (spk(20 : 20).eq.'0')
      enddo

      write(6,101) npk, spk
100 format('P.K. ',f20.3)
101 format(' Problem in term ',i20,' - - - > xpk =', a)

      end

```

The program with variables REAL\*8 is:

```

c ===== Programa Iterativesum.R8
c =====
c This program sums iteratively in double precission the numero a = 0.1
c until the rounding error affects the third decimal.
c =====

      implicit real * 8(a - h, o - z)
      character * 20 spk
      logical next

      npk = 0
      dpk = 0.1d + 00
      xpk = 0.0d + 00
      next = .TRUE.

      do while(seguir)
         npk = npk + 1
         xpk = xpk + dpk
         write(spk,'(f20.3)') xpk
         write(6,100) xpk
         next = (spk(20 : 20).eq.'0')
      enddo

      write(6,100) npk, spk
100 format('P.K. ',f20.3)
101 format(' Problem in term ',i20,' --- > xpk =',a)

      end

```

We can see that the operation  $xpk=xpk+dpk$  creates accumulation of storing errors.

The same effect will remain if the write a loop in which the index variable is of type real.

```

      do x = 0., 100., 0.1
         write(6,*) x
      enddo

```

**Solution 6.b)** The effect also appears when solving the problem in a spreadsheet since the operations being performed are exactly the same.

**Solution 6.b)** To solve the problem it would only be necessary to rewrite the loop changing the sum operation by a product of two real numbers, which has less propagation of errors. It means, changing  $xpk=xpk+dpk$  by  $xpk=0.1*FLOAT(npk)$