# PH20018 Coursework 1

## 1. The diffraction limit of a telescope

To do the calculations on this exercise I have assumed the value of $\lambda$ to be equal to 500 nm and the value of $I_o$ equal to 1. The Linux cc compiler has been used for the entirety of this coursework.

For both parts of the question, I have used the trapezium method to calculate the value of

$$J_m(x) = \frac{1}{\pi} \int_0^\pi \cos(m\theta - x\sin(\theta)) \, d\theta.$$

To make the code look cleaner and more readable I have defined the function

$$y_m(x,\theta) = \cos(m\theta - x\sin(\theta))$$

```
1.  float y(float t, float x, int m){
2.
3.          return cos(m*t - x*sin(t));
4.  }
```

Hence the expression for the approximation of $J_m(x)$ would be

$$J_m(x) = \frac{1}{\pi} * \frac{h}{2} * \left[ y_m(x,0) + 2 \sum_{i=1}^{N} y_m(x, h*i) + y_m(x,\pi) \right]$$

where $N = 10000$ is the number of subdivisions of the interval and $h = pi/N\text{-}1$ is the length of each interval on the x-axis.

```
1.  float J(int m, float x){
2.
3.          int N = 10000;
4.
5.          float h = pi / (float) (N-1);
6.
7.          int i = 1;
8.  //the varia a will store the sum of the 'middle value' of ym(x)
9.          float a = 0;
10.         while(i < N+1){
11.                 a = a + y(h*i ,x,m);
12.                 i++;
13.          }
14.
15.         return (1/pi)*0.5*h*(y(0.0, x, m) + 2*a + y(pi, x, m));
16.
17.  }
```

### Part a)

Once the functions needed are defined, we need to plot the values for $J_0(x)$, $J_1(x)$ and $J_2(x)$ within the interval from x = 0 to x = 20. To do so I have decided to write the values on a file in CSV format in order to make it easier to use the file with Python or Excel.

```
1.  void main(){
2.          //Open a .csv file to write the calculated values
3.          fp = fopen("data_1a.csv", "w");
4.
```

```
5.          //Table headers
6.          fprintf(fp, "x,J0(x),J1(x),J2(x)\n");
7.
8.          //Calculate the values of x between 0 and 20 for m = 0,1,2
9.          float x = 0.0;
10.           while(x <= 20){
11.               float a[3];
12.
13.               for(int m = 0; m < 3; m++){
14.                   a[m] = J(m,x);
15.               }
16.               //Print the values of each x for m=0,1,2
17.               fprintf(fp,                "%.2f,%f,%f,%f\n",x,
    a[0],a[1],a[2]);
18.
19.               x = x + 0.1;
20.           }
21.          fclose(fp);
22.
23.  }
```

This code will give us a file named data_1a.csv. Using Python, we can read this file and plot a graph for the different values of m and x. (Figure 1)

### Part b)

Adding a third function which gives us the value of the intensity as a function of r, the distance from the centre of the image, we can calculate the diffraction pattern for m = 1. We will use the formula
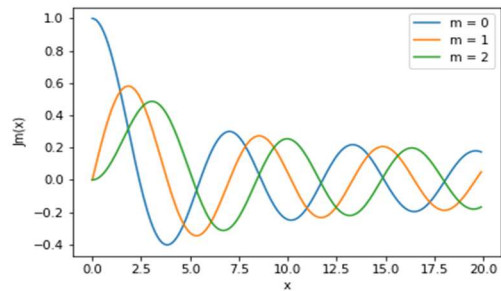


*Figure 1: The graph shows our approximation of the Bessel function for differnt values of m.*

$$I(x) = I_o \left( \frac{2 * J_1(x)}{x} \right)^2 \quad where \quad x = \frac{\pi * r}{\lambda} * 100$$

For the sake of simplicity, the function takes r in μm and λ in nm. Our function on C looks like this:

```
1. float Intensity(float r){
2.        //Simplifying the formula for x using r in μm and L in nm
3.        float x = pi * r * 100 / L;
4.        return Io*pow(2*J(1,x)/x,2);
5.
6. }
```

Following a similar procedure from the previous part of the exercise, the values of the intensity are written on a CSV and then plotted with Python. (Figure 2)

```
1. void main(){
2.        //open csv file
3.        fp = fopen("data_1b.csv", "w");
4.        // Table headers
5.        fprintf(fp, "r, I\n");
6.        float I;
```

```
7.
8.          int r = -25;
9.          while (r < 25 ){
10.             //Exception when x = 0 because it would be undefined
    with Bessel's formula
11.                 if (r == 0){
12.                     I = Io;
13.                 }else{
14.                     I = Intensity(r);
15.                 }
16.
17.                 fprintf(fp, "%d,%lf\n",r, I);
18.                 r++;
19.         }
20.         fclose(fp);
21.  }
```

In this exercise, we encounter a problem when r = 0 because it implies that x = 0 and therefore the value of the intensity at r = 0 is not well defined. However, because it is a division of floats the function in C returns a value, but it is a large number with does not make sense in this context. Therefore, I have decided to assign the value of I$_o$ to I(r) when r = 0.
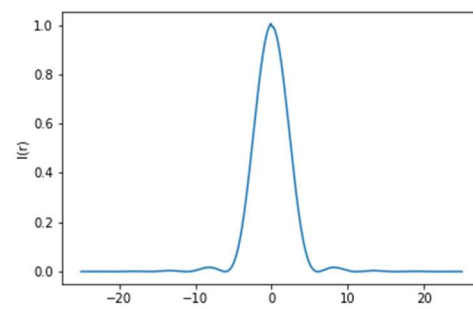


*Figure 2: This graph represents the diffraction pattern*

## 2.  Root-finding

Finding roots numerically with a computer although it is very useful has some limitations. In this case, we are using the **secant method**. This process is a variation of the Newton-Rapson method that is useful when the derivative of the function is unknown. The method works by plugging two initial guesses on the formula

$$x_{n+1} = x_n - f(x_n)\frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})}$$

and it will give you a better guess. If the process is repeated enough times the guess given by the formula might converge and thus approximate one of the roots. Depending on the two initial guesses the algorithm will either converge to one of the roots or diverge.

The main two problems with this method are the rounding errors when using floats and the fact that depending on which initial guesses one uses the solution might diverge. Another limitation worth pointing out is the difficulty this method has to find roots when f'(x) = 0, therefore the x-axis is tangent to the function.

The main source of rounding errors comes from divisions that have as a result a number with a very large number of decimal positions. In these cases, the computer approximates the value according to its capacity. This can lead to some inaccuracies. For example, when trying to check is two numbers are equal. Sometimes because of rounding errors two numbers that should be the same are not exactly the same but really close. In this scenario the solution I applied to fix it was to find the absolute value of the difference and check whether it was smaller than the required precision.

When the process diverges, we see two scenarios either the series of predictions go to plus or minus infinity or it gets on an infinite loop. The first scenario can be solved by checking on every iteration if the result of the secant method is too big and if it is, the program should restart everything with a different pair of initial guesses. On the other hand, to check if the program has fallen on an infinite loop the program can count the number of iterations and stop the loop when this number gets bigger than a pre-established limit.

Once these problems are addressed, we can write a program that finds all the roots of the function automatically with a level of precision specified by the user. This program would consist of a loop that changes the initial guesses and stores the solutions on an array. Afterwards, another set of nested loops would delete the repeated roots found and print on the screen all the solutions that the program has found. This algorithm is illustrated in the following flowcharts.
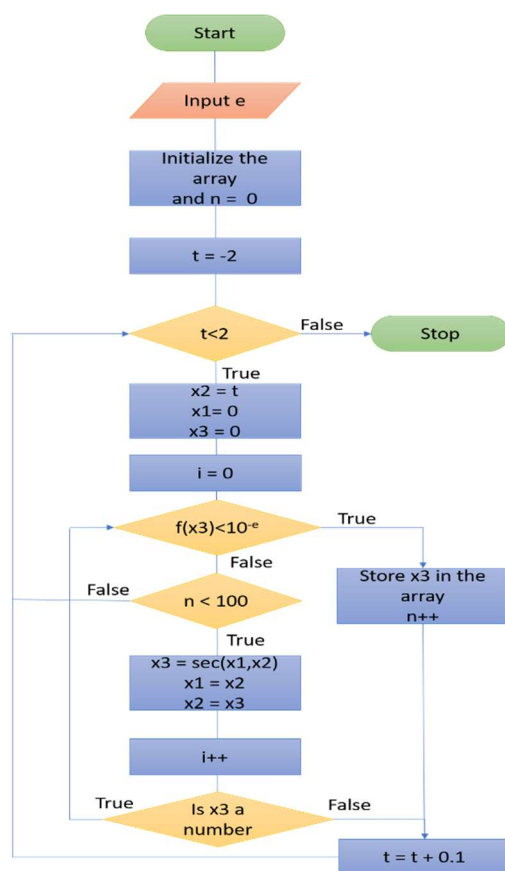


*Figure 4: This illustration shows the main process that the program flows in order to find the roots and append them to the array.*
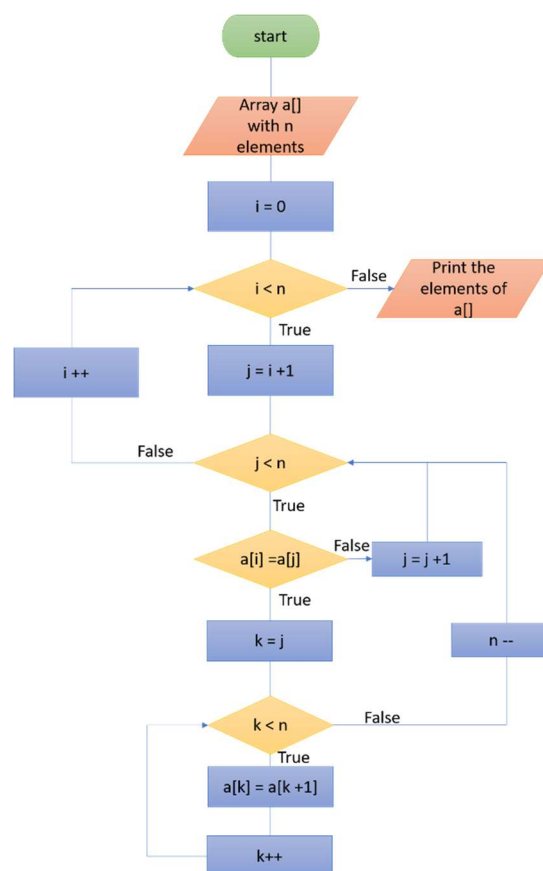
*Figure 3 The chart above represents the algorithm to delete the repeated roots from the array.*

Now that we know the structure the program should follow in order to find the roots of any function; we need to write a script that solves the equation

$$f(x) = \frac{1}{(x+1)^2} + \frac{1}{(x-1)^2} - 10$$

but that can also be easily generalised to solve any other function.

When it comes to writing the code first of all we have to define the functions that are going to be used throughout the program. In this case for the sake of readability of the code one function represents f(x) and the other the secant method formula.

```
1. //Define the function
2. float f(float x){
3.    return 1/pow(x+1,2) + 1/pow(x-1,2) - 10;
4. }
5.
6. //Secant method equation
7. float sec(float x1,float x2){
8.    return x2 - f(x2) * (x2-x1)/(f(x2) - f(x1));
9. }
```

The next step would be to define the algorithm that keeps applying the equation until a root is found and stops if the process diverges. To check if the process is diverging the number of iterations is stored on a variable and every iteration is checked to be less than 100, it could be much less with this equation but I believe it is good to give plenty of margins I case the program is used to solve a more complex equation. I decided to wrap this algorithm into a function so the main() is less full.

```
1. float findRoot(float x1,float x2, int e){
2.         float x3 = 0.0;
3.         int n_iterations = 0;
4.         //Execute the secant method
5.         while (fabs(f(x3))> pow(10,-e)){
6.
7.             if (n_iterations< 100){
8.                 x3 = sec(x1,x2);
9.                 x1 = x2;
10.                    x2 = x3;
11.                 n_iterations++;
12.                 }
13.             else{
14.                 return 0.0/0.0;
15.                 }
16.             }
17.         return x3;
18. }
19.
```

 It is worth noticing that when the number of iterations exceeds the limit the function is forced to return the data type *nan*, which stands for not a number. This rather weird action is justified because it is a way of turning two different problems into one. The first one is when *x2* goes to infinity, which at some point the computer assigns the value *inf* to the variable which leads to the *x3* variable being assigned *nan* on the following iteration. And the second one is when the maximum number of iterations is passed. By forcing the root to have value *nan* the program can just check is the value of the root is different to *nan* before appending it to the array, as shown on another code snippet below.

Inside the main function, we first find the lines of code that ask the user the number of decimal figures needed. In this case, because we are using floats the maximum number of decimal positions is 6 and the program is not able to find any roots if a number higher than 5

significant figures is required. If more precession was needed instead of using floats one could use doubles.

```
1. //Define the how many decimal positions of precision are requied
2.    int e;
3.    printf("Type the number (1 to 5) of significant figures:   ");
4.    scanf("%d", &e);
```

Secondly, we have the loop that changes the value of one of the initial guesses every time a root if found and fills the array with the valid ones.

```
1. float x1, x2;
2. //Array that will store the roots of f(x)
3. float a[100];
4. //Count the roots found
5. int n = 0;
6.
7. // Loop to change the initial value of x2
8. for (float t = -2.0; t < 2.1; t = t +0.1){
9.   // Assaing initial values to x1, x2, x3
10.          x2 = t;
11.          x1 = 0.0;
12.          float root = findRoot(x1,x2,e);
13.          if(!isnan(root)){
14.                a[n] = root;
15.                n++;
16.          }
17.    }
```

Once the array is filled with all the roots that the program has managed to find we need to remove the ones repeated. The array *a* has been initialized with a size of 100 to make sure it can fit all the roots. Otherwise, the program would yield a *core dumped error*. The variable *n* counts how many roots have been appended, which is useful for the algorithm that removes the duplicated roots. The main idea of this algorithm is to compare each element of the array with the ones after it, starting with the first one. If a repeated number is found it is deleted and every element is shifted one position to the "left".

```
1.    // Check for duplicated roots and remove them
2.    int i,j,k;
3.    for(i=0; i<n; i++){
4.          for(j=i+1; j<n; j++){
5.                // If any duplicate found
6.                if(fabs(a[i] - a[j]) < pow(10,-e)){
7.                      // Delete the current duplicate element
8.                      for(k=j; k<n; k++){
9.                            a[k] = a[k + 1];
10.                      }
11.                      //Decrement size after removing
   duplicate element
12.                      n--;
13.                      // If shifting of elements occur
   then don't increment j
14.                      j--;
15.                }
16.          }
17.    }
```

The final step is to show on the screen the unique roots calculated. This is easily done with a loop iterating over the array which now has only the unique roots.

```
1.  //Print array after deleting duplicated elements
2.    printf("\nThe program has found %d solutions : \n", n);
3.
4.            for(i=0; i<n; i++){
5.                    printf("x%d = %f\n",i+1, a[i]);
6.    }
```

The output on the console would look something like this:

```
ecf35@linux3:~/ccourse/coursework1$ ./a.out
Type the number (1 to 5) of significant figures:  4

The program has found 4 roots:
x1 = -0.678003
x2 = -1.319209
x3 = 0.678002
x4 = 1.319209
```

When the equation f(x) = 0 is solved analytically we get four solutions:

$$x_{1,2} = \pm\sqrt{\frac{11}{10} + \frac{\sqrt{40}}{10}} \approx \pm 1.319209014, \quad x_{3,4} = \pm\sqrt{\frac{11}{10} - \frac{\sqrt{40}}{10}} \approx \pm 0.678002637$$

## 3. Multivariable equations

Solving systems of equations numerically is extremely useful. One of the most important applications is solving non-linear differential equations. The algorithm used in this case is an adaptation of the Newton method.

If there are have *n* non-linear equations with *n* variables can be written like this:

$$F_0(x_0, \dots, x_n) = 0$$

$$\vdots$$

$$F_n(x_0, \dots, x_n) = 0$$

If we use vectorial notation the system can be written as **F(x) = 0**.

Following the logic from the Newton method, we need to approximate **F(x)** by a linear function with the form **F(x) = M x + c**, where M is an *n x n* matrix and c is a vector.

To approximate an **F(x)** with a linear function we can use the first two terms of its Taylor expansion. Given the value of **F** and its partial derivatives with respect to **x** at some point $x_i$, we can approximate the value at some point **x**$_{i+1}$ by the two first term in a Taylor series expansion around **x**$_i$ :

$$F(x_{i+1}) \approx F(x_i) + J(x_i)(x_{i+1} - x_i)$$

Where **J** is the Jacobian. The idea is to approximate **F(x**$_{i+1}$**) = 0**. Therefore, we can rewrite the equation as

$$F(x_i) + J(x_i)\delta = 0, \text{where } \delta = (x_{i+1} - x_i)$$

Now it is clearer what the program must solve:

1. Calculate $\delta = -J(x_i)^{-1} F(x_i)$.
2. $x_{i+1} = (x_i + \delta)$

## Implementation

In this case, we are attempting to solve the following system of equations:

$$(x + 1)^2 + (y + 1)^2 = 25$$

$$xy + y^2 = 5$$

First of all, we need to define some global variables that will be used throughout the program.

```
1. //Array of pointers of the functions
2. float(*f[2])(float,float);
3. //Ventor F(x)
4. float F[2];
5. //Jacobian matrix
6. float m[2][2];
7.
8. //Number of significant figures
9. int e;
```

We have defined the array that Will hold the pointers of the functions, this makes it easier to generalize the approximation of the derivative. Secondly, we see the vector that will save the values of $f_1(x, y)$ and $f_2(x, y)$. Then the matrix $m$ will store the values of $J(x_i)$. Lastly, $e$ will store the number of significant figures required.

Next, we need to define the functions that we will be using and the expressions for the approximations of the derivatives:

```
1. float f1 (float x, float y){
2.    return pow(x+1,2) + pow (y+1,2)-25;
3. }
4.
5. float f2 (float x, float y){
6.    return x*y + pow(y,2) - 5;
7. }
8.
9. float dfx(int i, float x, float y){
10.         return (f[i](x,y) - f[i](x +0.001,y))/(0.001);
11.  }
12.
13.  float dfy (int i,float x, float y){
14.         return (f[i](x,y) - f[i](x,y +0.001))/(0.001);
15.  }
```

The next function that we need to define is the one that calculates the Jacobian. In this case, I have decided to create a void function that fills up the matrix $m$ declared above, instead of a function that returns a matrix. Returning arrays and matrices with C ca be rather tricky because we would need to work with pointers.

```
1.  //Calculates de jacobian matrix
2.  void J(float r[2]){
3.
4.     //Aproximate de derivatives
5.     m[0][0] = dfx(0,r[0],r[1]);
6.     m[0][1] = dfy(0,r[0],r[1]);
7.     m[1][0] = dfx(1,r[0],r[1]);
8.     m[1][1] = dfy(1,r[0],r[1]);
9.
10.    // inverts the matix
11.    float det;
12.    det = m[0][0] * m[1][1] - m[1][0]*m[0][1];
13.    float inv[2][2];
14.    inv[0][0] =   m[1][1] / det;
15.    inv[1][0] = - m[1][0] / det;
16.    inv[0][1] = - m[0][1] / det;
17.    inv[1][1] =   m[0][0] / det;
18.
19.    // saves the found values on the matrix declared above
20.    for (int i = 0 ; i < 2 ; i++){
21.            for (int j = 0 ; j < 2 ; j++){
22.                    m[i][j] = inv[i][j];
23.            }
24.    }
25.  }
```

Once the Jacobian is defined, we need to build a function that calculates $\delta$. In this case, the function will return an array, since it is only a 1D array working with it is not complicated. This function will create a vector and fill it with the values obtained after multiplying the inverse of the Jacobian and **F(x)**.

```
1.  float *delta(float r[2]){
2.
3.     //Create the jacobian
4.     J(r);
5.
6.     static float d[2];
7.     //Find the delta
8.     for (int i = 0 ; i < 2; i++){
9.             d[i] = m[i][0] * f[0](r[0],r[1])+ m[i][1] *
    f[1](r[0],r[1]);
10.            }
11.
12.            return d;
13.  }
```

 The last function that needs to be defined before the main is the one that calculates the norm of **F(x)**. This is useful because on every iteration we need to check whether or not $f_1(x, y)$ = and $f_2(x, y) = 0$, and thus the norm of **F**. For the sake of simplicity, I have decided to calculate the norm as the simple sum of the absolute value of each component. This does not cause any problems because we just need to check if all the components are zero and not the real magnitude.

```
1. //Find the norm of F(x)
2. float norm(float f[], int size){
3.
4.    float n = 0;
5.    for(int i = 0; i < size; i++){
6.            n = n + fabs(f[i]);
7.    }
8.    return n;
9. }
```

Inside the *main* function, we first assign the value of *e* for the precision and introduce la pointers of the $f_1$ and $f_2$ inside the vector f. Secondly, we initialize the arrays: *r* that will store the values of x and y and will be updated on every iteration, the array *a* in a similar way to the previous exercise *a* will collect every solution found, and finally *p* that its use will be explained at the appendix.

```
1.    printf("Please type the number of significant figures required:
      \n");
2.    scanf("%d", &e);
3.
4.    //Put the functions in the array
5.    f[0] = f1;
6.    f[1] = f2;
7.
8.    float r[2];
9.    //Array to store the values to plot the graph
10.   float p[1001][4];
11.   //Array to store all the roots found
12.   float a[1001][2];
13.   //n number of elements in a , m number of elements in p
14.   int n = 0,m = 0;
15.
```

Following a similar logic to the previous exercise, the main algorithm consists of a for loop that gives the initial value of x and y, and a nested while loop that updated the initial values until it finds a solution. The program has the same way of checking for divergence that exercise 2 has. To give an initial value for x and y I have decided to use a spiral function. This way the initial values cover a uniform area around the origin that could be useful when trying to solve more complex systems.

```
1. for( float t = 0; t < 100 ; t= t + 0.1){
2.
3.            float *x = &r[0];
4.                 float *y = &r[1];
5.
6.            r[0] = 0.1*(cos(t) + t*sin(t));
7.            r[1] = 0.1*(sin(t) - t*cos(t));
8.
9.            p[n][0] = *x;
10.           p[n][1] = *y;
11.
12.           F[0] = f[0](*x,*y);
13.           F[1] = f[1](*x,*y);
14.
15.           float norm_F = norm(F, 2);
16.           int n_iterations = 0;
```

```
17.
18.         while (fabs(norm_F) > pow(10,-e) && n_iterations <
   100){
19.
20.                 float *d = delta(r);
21.
22.                 r[0] = r[0] + d[0];
23.                 r[1] = r[1] + d[1];
24.
25.                 F[0] = f[0](*x,*y);
26.                 F[1] = f[1](*x,*y);
27.                 norm_F = norm(F, 2);
28.                 n_iterations++;
29.         }
30.
31.         if (!isnan(r[0])&&!isnan(r[1])&&n_iterations < 100){
32.
33.                 a[n][0] = *x;
34.                 a[n][1] = *y;
35.
36.                 p[n][2] = *x;
37.
38.                 m++;
39.                 n++;
40.         }
41.  }
```

Once the array *a* is filed with all the solutions found we can adapt and apply the algorithm used in the previous exercise to remove duplicated solutions.

```
1.   //remove duplicates
2.   int i,j,k;
3.   for(i=0; i<n; i++){
4.         for(j=i+1; j<n; j++){
5.                 // If any duplicate found
6.                 if(fabs(a[i][1] - a[j][1]) < pow(10,-e) &&
   fabs(a[i][0]-a[j][0])< pow(10,-e)){
7.
8.                         // Delete the current duplicate element
9.                         for(k=j; k<n; k++){
10.                                a[k][0] = a[k + 1][0];
11.                                a[k][1] = a[k + 1][1];
12.                        }
13.                        //Decrement size after removing
   duplicate element
14.                        n--;
15.                        // If shifting of elements occur then
   don't increment j
16.                        j--;
17.                }
18.         }
19.  }
```

Finally, the same lines of code that generate the output of the second exercise can be used to produce this one.

```
1.    //Print array after deleting duplicate elements
2.    printf("\nThe program has found %d solutions : \n", n);
3.    for(i=0; i<n; i++){
4.        printf("x%d = %f, y = %f\n",i+1, a[i][0],a[i][1]);
5.    }
```
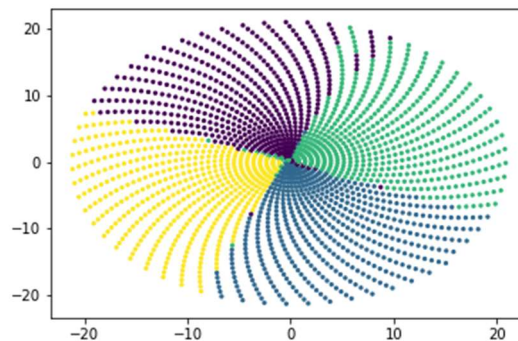
```
ecf35@linux1:~/ccourse/coursework1$ ./a.out
Please type the number of significant figures required:
5

The program has found 4 solutions :
x1 = -2.460051, y = 3.782076
x2 = 2.906849, y = -4.120341
x3 = 3.546557, y = 1.080581
x4 = -5.993355, y = -0.742316
```

The roots of uncertainties and limitations are the same as in exercise 2.

### Appendix

During the development of this piece of software, I encountered some problems. But without a doubt, the hardest to solve was that at some point the program would find every solution but one. To see what was happening I came up with the idea of plotting a graph showing the initial guess with a point in the 2D plane and the colour of the point would indicate which solution had that initial guess lead to.



Here is where the array *p* comes into play. This array is saving the initial values of x and y and the x coordinate of the root.

Using a series of loops the data on *p* is written on a CSV file and the graph is plotted with python.

This feature can be very useful when working with more complex functions.

```
1.  // Fill the last column with the indexes and write the .csv
2.    FILE *fp;
3.    fp = fopen("map.csv", "w");
4.    for(i = 0 ; i<m; i++){
5.        for( j = 0; j<n ; j++){
6.            if (fabs( p[i][2] - a[j][0]) < pow(10,-e)){
7.                p[i][3] = j;
8.                fprintf(fp,"%f, %f, %f\n",p[i][0],
   p[i][1],p[i][3]);
9.
10.            }
11.        }
12.    }
13.    fclose(fp);
```