

# Iterative search and lineal regression

Gradient descent iterative method and lineal regression

**Blanca Cano Camarero**

Department: DECSAI

University: ETSIIT, Granada university

Country: Spain

Date: April 1, 2021

Mathematics and computer engineering degrees,  
Doble Grado matemática e informática



# Contents

<b>1</b>	<b>Gradient descent</b>	<b>5</b>
1.1	Gradient descent 's algorithm . . . . .	5
1.1.1	Introduction . . . . .	5
1.1.2	Math . . . . .	5
1.1.3	Algorithm . . . . .	6
1.1.4	Problem 1 . . . . .	6
1.1.5	Problem 2 . . . . .	8
1.1.6	Final conclusion about finding global functions' minimum . . . . .	14
<b>2</b>	<b>Linear Regression</b>	<b>15</b>
2.1	Linear regresion . . . . .	15
2.1.1	Stochastic gradient descendent . . . . .	15
2.1.2	How to plot . . . . .	15
2.2	Experiment . . . . .	16
2.2.1	a) Generate a training sample . . . . .	16
2.2.2	b) Labels, noise and map . . . . .	16
2.2.3	Estimate the fitting error of $E_{in}$ using SDG . . . . .	18
	<b>Bibliography</b>	<b>19</b>



# Chapter 1

## Gradient descent

### 1.1 Gradient descent 's algorithm

#### 1.1.1 Introduction

Gradient descent is a general technique for minimizing twice-differentiable functions through its slope. [1] It is used to find local minimums. The start point is crucial in the search.

The basic idea is to update the weights using the gradients until it is not possible to continuous minimizing the error.

#### 1.1.2 Math

In order to understand the algorithm we are going to define:

Let  $w(0) \in \mathbb{R}^d$  be an arbitrary initial point,  $E : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  a class  $C^1$  function. The learning rate or step size  $\eta \in \mathbb{R}^+$  is a experimental coefficient about how much are we going to follow the slope to obtain the new weight. Let  $w(t) \in \mathbb{R}^d$   $t \in \mathbb{N}$  be the weight for  $t$  iteration which is defined as

$$w(t+1) = w(t) - \eta \nabla E_{in}(w(t))$$

#### Properties

- This algorithm gives local minimums.
- Convergence it is not assured, so it would be necessary some stop criteria.
- For a convex function it would be a unique global minimum.
- The learning rate  $\eta$  variable in time is important: fixes learning gradient descent algorithm.

### 1.1.3 Algorithm

The following code snippet implement the algorithm, where  $w(0)$  is the *initial\_point*,  $E$  is the error,  $\nabla E_{in}(w)$  is *gradient\_function* and finally  $\eta$  is *eta*. The value  $\eta = 0.1$  is a heuristic basic on purely practical observation [1].

In order to avoid an infinite search, our stop criteria are a limit in the number of iterations *max\_sitter* and an error tolerance.

```
def gradient_descent(initial_point, loss_function,
                    gradient_function, eta, max_iter, target_error):
    '''
    initicial point: w_0
    E: error function
    gradient_function
    eta: step size

    ### stop conditions ###
    max_iter
    target_error

    #### return ####
    (w, iterations)
    w: the coordenates that minimize E
    it: the numbers of iterations needed to obtain w

    '''

    iterations = 0
    error = E( initial_point[0], initial_point[1])
    w = initial_point

    while ( (iterations < max_iter) and (error > target_error)):

        w = w - eta * gradient_function(w[0], w[1])

        iterations += 1
        error = loss_function(w[0], w[1])

    return w, iterations
```

### 1.1.4 Problem 1

We want to solve the following problem:

Use gradient descent's algorithm to find a minimum for the function

$$E(u, v) = (u^3 e^{(v-2)} - 2 * v^2 e^{-u})^2.$$

Set  $(u, v) = (1, 1)$  as initial point and use learning rate  $\eta = 0.1$ .

**Compute analytically the gradient of  $E(u, v)$**

$$\begin{aligned} \nabla E(u, v) &= \left( \frac{\partial}{\partial u} (u^3 e^{(v-2)} - 2 * v^2 e^{-u})^2, \frac{\partial}{\partial v} (u^3 e^{(v-2)} - 2 * v^2 e^{-u})^2 \right) = \\ &= \left( 2(u^3 e^{(v-2)} - 2 * v^2 e^{-u})(3u^2 e^{(v-2)} + 2v^2 e^{-u}), 2(u^3 e^{(v-2)} - 2 * v^2 e^{-u})(u^3 e^{(v-2)} - 4v e^{-u}) \right) \end{aligned}$$

**Number of iterations and final coordinates.**

Firstable we need to use 64-bits float, so we are going to use the data type `float64` of numpy library [2].

The functions' declaration are:

```
def dEu(u,v):
    '''
    Partial derivate of E with respect to the variable u
    '''
    return np.float64(
        2
        *( 3* u**2 * np.e**(v-2) + 2*v**2 * np.e**(-u) )
        *( u**3 * np.e**(v-2) - 2*v**2 * np.e**(-u))
    )

def dEv(u,v):
    '''
    Partial derivate of E with respect to the variable v
    '''
    return np.float64(
        2*
        ( u**3 * np.e**(v-2) - 2*v**2 * np.e**(-u) )
        *( u**3 * np.e**(v-2) - 4*v * np.e**(-u))
    )

def gradE(u,v):
    '''
    gradient of E
    '''
    return np.array([dEu(u,v), dEv(u,v)])
```

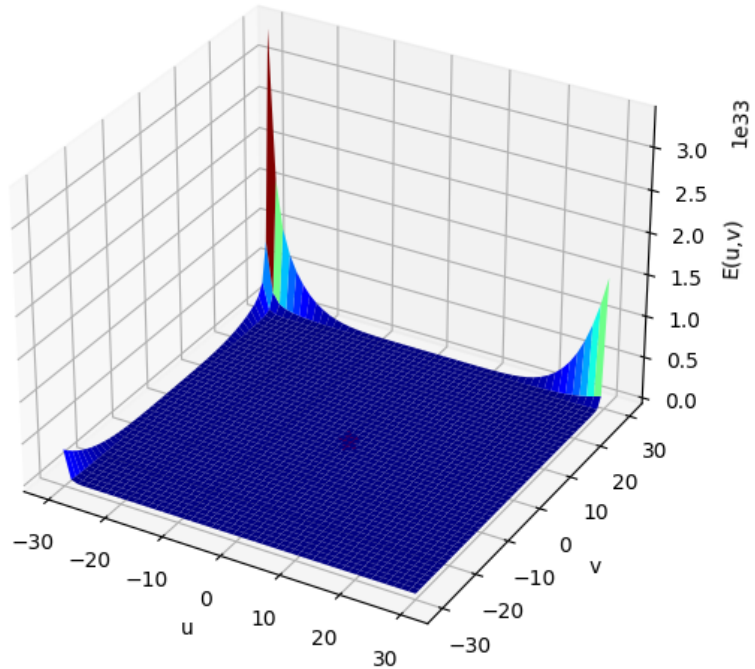
To obtain the number of iterations and the final coordinates, the only thing we need to do is to call *gradient\_descent* function with the initial conditions:

```
eta = 0.01
max_iter = 10000000000
target_error = 1e-14
initial_point = np.array([1.0,1.0])
w, it = gradient_descent( initial_point,
                           E,
                           gradE,
                           eta,
                           max_iter,
                           target_error )
```

The result are:

- Numbers of iterations: 178.
- Final coordinates: (1.162, 0.924).

A 3d graph with the result is



### 1.1.5 Problem 2

For function  $f(x, y) = (x + 2)^2 + 2(y - 2)^2 + 2 \sin(2\pi x) \sin(2\pi y)$



**Use gradient descent to minimize  $f$** 

The initial point is  $(x_0 = -1, y_0 = 1)$ , learning rate is  $\eta = 0.01$  and the maximum number of iterations must be 50. Plot the result and repeat the experiment with  $\eta = 0.1$ .

Firstly we are going to calculate partial derivatives and gradient of  $f$ .

$$\frac{\partial}{\partial x} f = 2(x + 2) + 2 \sin(2\pi y) \cos(2\pi x) 2\pi = 2(x + 2) + 4\pi \sin(2\pi y) \cos(2\pi x)$$

$$\frac{\partial}{\partial y} f = 2(y - 2) + 4\pi \sin(2\pi x) \cos(2\pi y)$$

It is important to realise that  $f(x, y) < 0$  for some values in  $\mathbb{R}^3$  so the error target has been omitted in this algorithm.

Now the new algorithm is

```
def gradient_descent_trace(initial_point, loss_function,
    gradient_function, eta, max_iter):
    '''
        inicial point: w_0
        loss_function: error function
        gradient_function
        eta: step size

        ### stop conditions ###
        max_iter

        #### return ####
        (w, iterations)
        w: the coordenates that minimize loss_function
        it: the numbers of iterations needed to obtain w

    '''

    iterations = 0
    error = loss_function( initial_point[0], initial_point[1])
    w = [initial_point]

    while iterations < max_iter:

        new_w = w[-1] - eta * gradient_function(w[-1][0], w[-1][1])

        iterations += 1
```

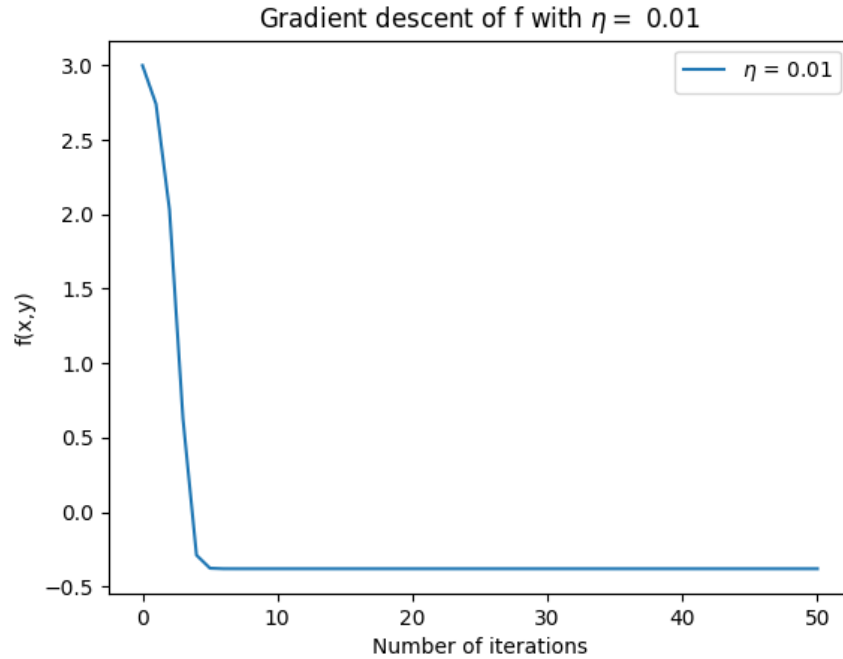
```

    error = loss_function(new_w[0], new_w[1])
    w.append( new_w )

    return w, iterations

```

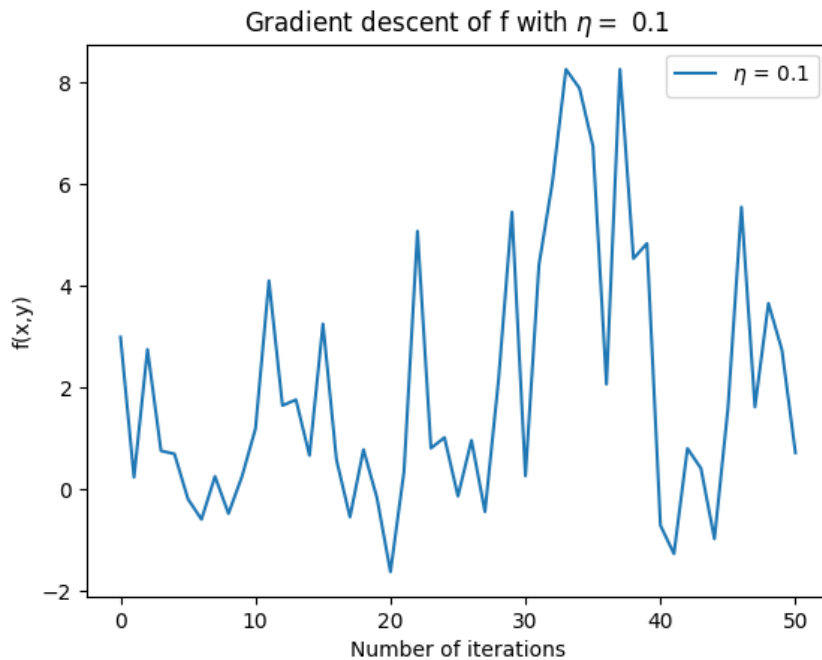
After 50 iterations for  $\eta = 0.01$ , the final coordinates are  $(-1.269, 1.287)$  and their value is  $-0.381$ . The graph, which shows the relation between iterations and the function minimization is



As far as we have seen, before the 10<sup>th</sup> iteration we are really close to the minimum and stay there without fluctuate.

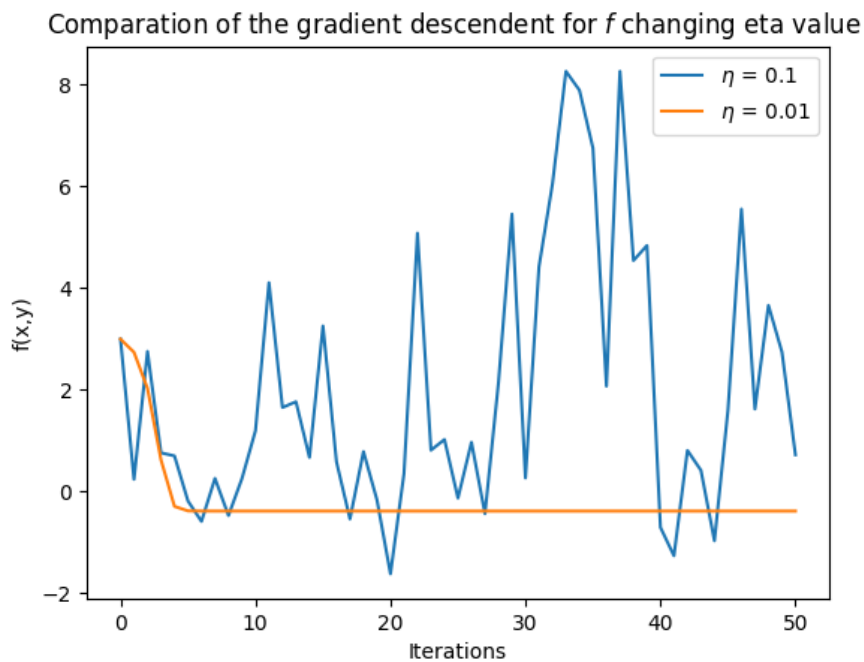
On the other hand, after 50 iterations for  $\eta = 0.1$  the final coordinates are  $(-2.939, 1.608)$  and their value is 0.724, so as we can see this result is worse than the last one.

In the following graph we can see the evolution fluctuation of the images iteration by iteration



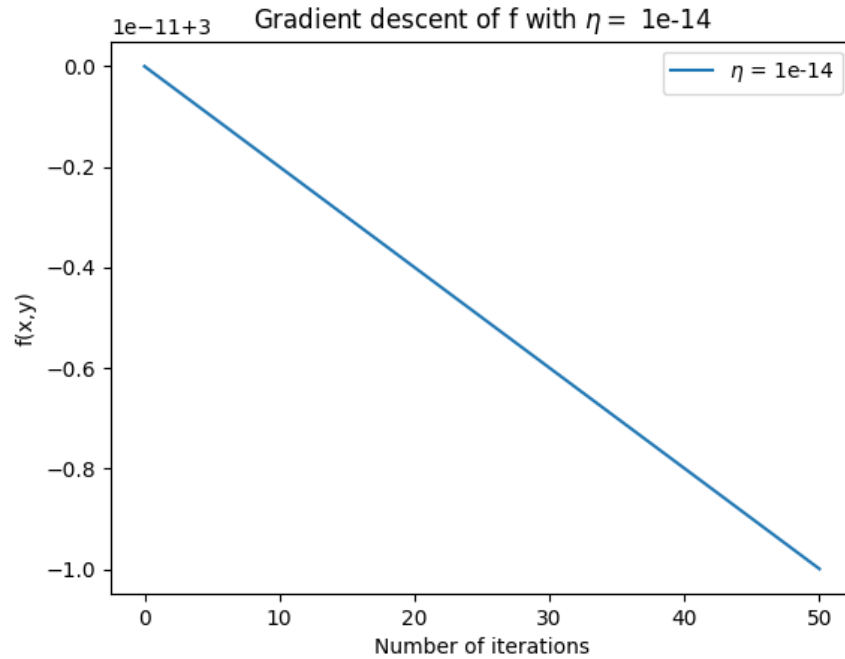
The reason for this irregularity is that the step size is too big, so it skips the minimum.

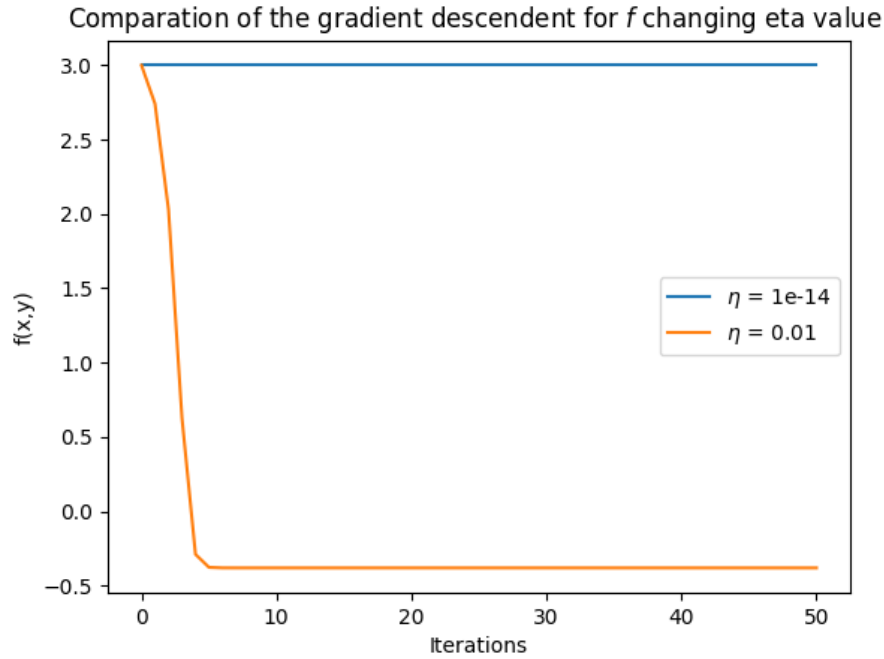
We can also compare the two experiment in the following graph.



Moreover, basic on its mathematics proof, which use Taylor's series, we know that it should be small, but if it is too small the algorithm will never reach the minimum in a right time.

Let see a new example, now  $\eta = 10^{-14}$ , after 50 iteration the final coordinates are  $(-1, 1)$  and the value is 3, so this new selection is even worse that the one with the bigger step size, although it goes without oscillate.





As a conclusion, a priory, it is difficult to select a step size value, each problem should have an appropriate one, and the selection must be empirical. Even though some heuristic [1] tell that  $\eta = 0.01$  its a good try.

### Minimum value

Before running the algorithm is important to think about a good value for the learning rate  $\eta$ . Based on the last section  $\eta = 0.1$  is a good one.

After run the program the results are

Initial point	Final coordinates	Final value
(-0.5 -0.5)	(-0.793 -0.126)	9.125
(1 1)	(0.677 1.29)	6.437
( 2.1 -2.1)	( 0.149 -0.096 )	12.491
(-3 3)	(-2.7315 2.713)	-0.381
(-2 2)	(-2. 2.)	0

This example gives the idea that the local minimums found depend on the start point and a priory, unless we know some properties of the function such as convexity or monotony we are not able to assure that the minimum found is global.

For a mathematical study we need differentiable functions, and a technique to find where their zeros are, so we need to solve their equations. Some time this is not possible and we use numeric methods.

### 1.1.6 Final conclusion about finding global functions' minimum

- Properties differentiable, arbitrary initial point, eta, computational cost...  
TO-DO

## Chapter 2

# Linear Regression

### 2.1 Linear regresion

#### 2.1.1 Stochastic gradient descent

Stochastic gradient descent (SGD) is a sequential version of the gradient descent. Instead of considering the full batch gradient on all  $N$  training data points, we consider a stochastic version of the gradient. First, pick a training data point  $(x_n, y_n)$  uniformly random (hence the name 'stochastic') and consider only the error on that data point.

The gradient of this single data point's error is used for the weight update in exactly the same way that the gradient was used in batch gradient descent.

Motivo de 1, -1 para compensar la etiqueta

Hacer 200

W no importa en la práctica (podemos manipularlo)

VISUALIZAR MODELO CON SCATTER PLOT (DIAPPOSITIVA 22, CON LOS DATOS ENTRENAMIENTO Y TEST)

E w deberá de ser la misma

Incluir valor de  $E_{in}$  y error de clasificación (lo normal es que el error en test sea un poquito peor que los de entrenamiento).

Gradiente descendiente minimizar iterativamente:

En el segundo no tiene por qué darse que una iteración sea mejor (otra cosa es que mejore la tendencia mejor)

Si se utilizan todos sí que debería de ir mejorando de manera global.

#### 2.1.2 How to plot

As we know  $w = (cte, intensity - coefficient, symmetry - coefficient) = (c, i, s)$  have three parameters. And we are going to plot in a 2D graph where  $c + xi + ys = 0$ , it is a line, so we are going to compute two points, the one that x

$$(x = 0) : y = \frac{-c}{y} \quad (x = 1) : y = \frac{-c - i}{y}$$

## 2.2 Experiment

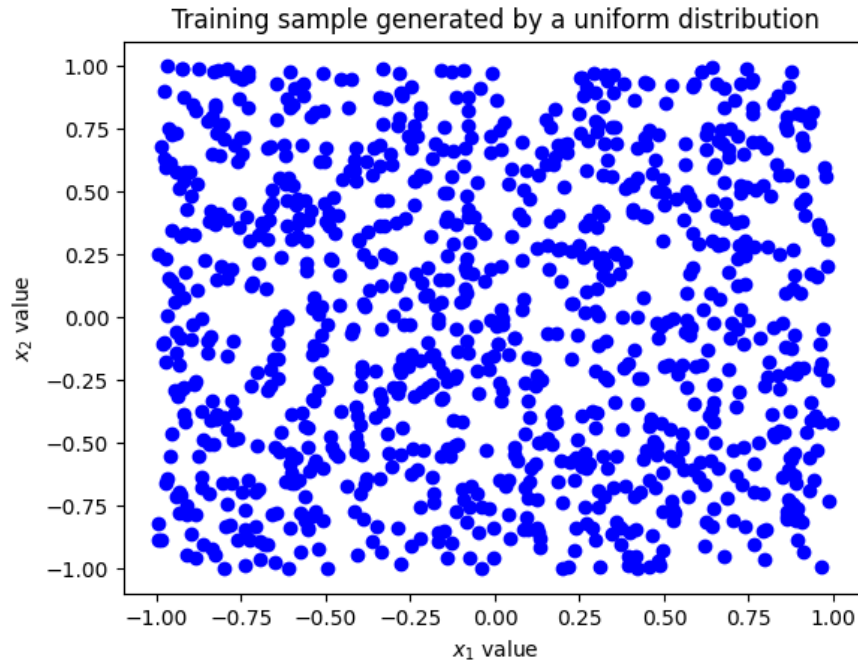
### 2.2.1 a) Generate a training sample

We are going to use a uniform generation:

```
def simula_unif(N, d, size):
    ''' generate a training sample of N points
    in the square [-size,size]x[-size,size]
    '''
    return np.random.uniform(-size,size,(N,d))
```

After fixed random seed to 1.

The final 2D map is



### 2.2.2 b) Labels, noise and map

b) Let's consider the function  $f(x_1, x_2) = \text{sign}((x_1 - 0.2)^2 + x_2^2 - 0.6)$  that we will use to assign a label to each point of the previous sample. We introduce noise on the labels, randomly changing the sign of 10 % of them. Draw the obtained labels map.

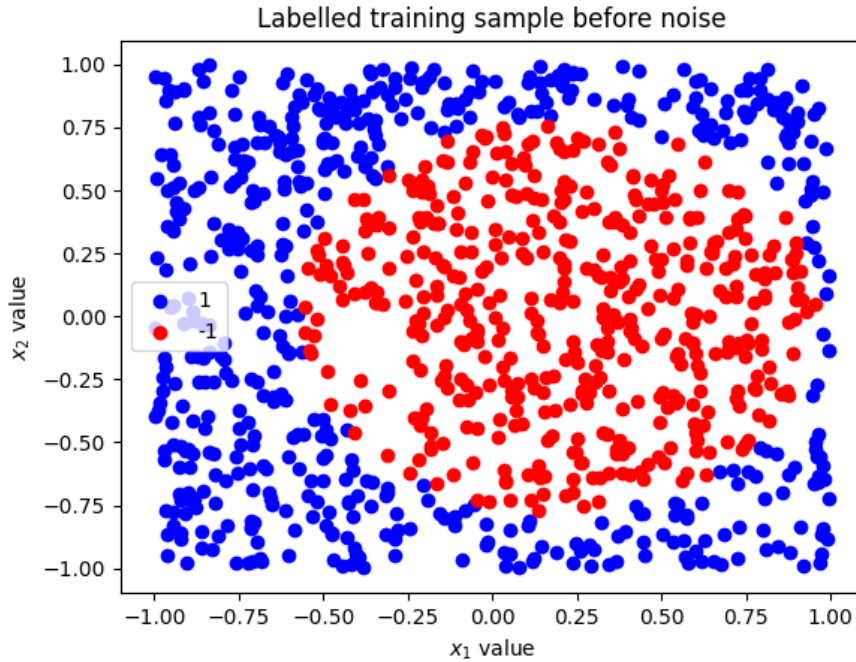
Before plotting, while analysing the function is important to have in mind that a circumference with radius  $r$  and center  $(c_1, c_2) \in \mathbb{R}^2$  are the points  $(x_1, x_2) \in \mathbb{R}^2$  that verify



$$(x_1 - c_1)^2 + (x_2 - c_2)^2 = r^2$$

Therefore looking at  $f$  it is easy to think that we are going to see a circle of radius  $\sqrt{0.6}$  and center  $(0.2, 0)$ .

The plotting is



In order to introduce noise on the label we are going to change randomly the sign of the 10% of the labels obtained by  $b$ .

```
#labels
y = np.array( [f(x[0],x[1]) for x in training_sample ] )

index = list(range(size_training_example))
np.random.shuffle(index)

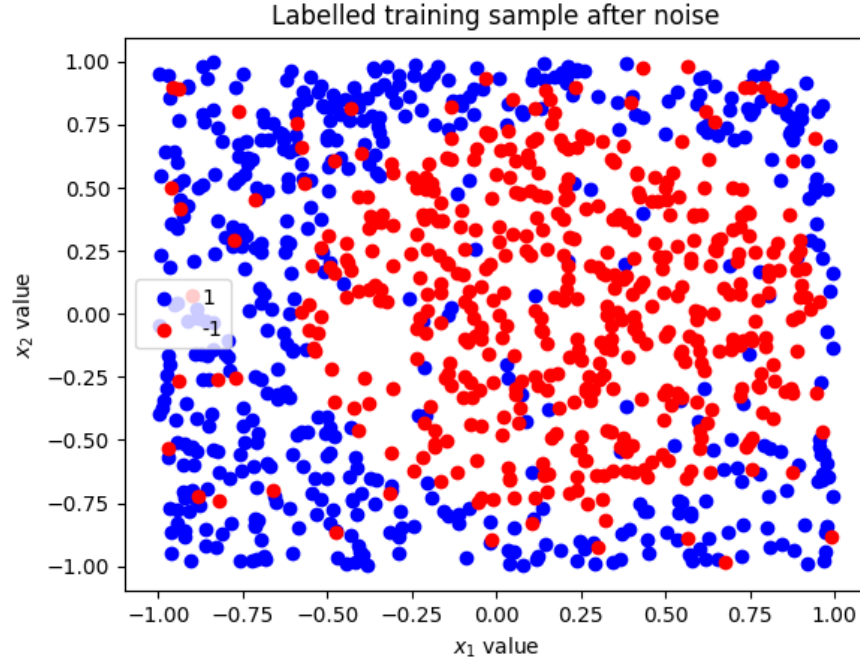
percent_noisy_data = 10.0
size_noisy_data = int((size_training_example *percent_noisy_data)/ 100 )

noisy_y = np.copy(y)
for i in index[:size_noisy_data]:
    noisy_y[i] *= -1
```

As we can see the idea behind the snippet is simple: The noised labels would be a copy of the original one and 10% of the data would change their

sign.

The final map is :



### 2.2.3 Estimate the fitting error of $E_{in}$ using SDG

c) Using  $(1, x_1, x_2)$  as feature vector, fit a linear regression model to the generated dataset and estimate the weights  $w$ . Estimate the fitting error of  $E_{in}$  using Stochastic Gradient Descent (SGD).

# Bibliography

- [1] Hsuan-Tien Lin Yaser S. Abu-Mostafa, Malik Magdon-Ismail. *Learning From Data. A Short Course*. AMLbook, 2012.
- [2] Numpy documentation. Numpy basic data types documentation, 2021.