

```
In [1]: #help function
def set_spines():
    ax = plt.gca() # ax stands for 'get current axis'
    ax.spines['right'].set_color('none')
    ax.spines['top'].set_color('none')
    ax.axis.set_ticks_position('bottom')
    ax.spines['bottom'].set_position(('data',0))
    ax.yaxis.set_ticks_position('left')
    ax.spines['left'].set_position(('data',0))

In [40]: %matplotlib inline
import matplotlib.pyplot as plt
import numpy as np
import math
from mpl_toolkits.mplot3d import Axes3D
from random import random
from mpl_toolkits.mplot3d import Axes3D
import matplotlib.pyplot as plt
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LinearRegression
```

# Linear Regression

- 1 Simple Linear Regression
  - 1.1 Fitting the regression line
  - 1.2 Gradient Descent over simple linear regression
  - 1.3 Effect of different values for learning rate
- 2 Multiple Linear Regression
  - Regularization of gradient descent by learning rate and max iterations
- Conclusion

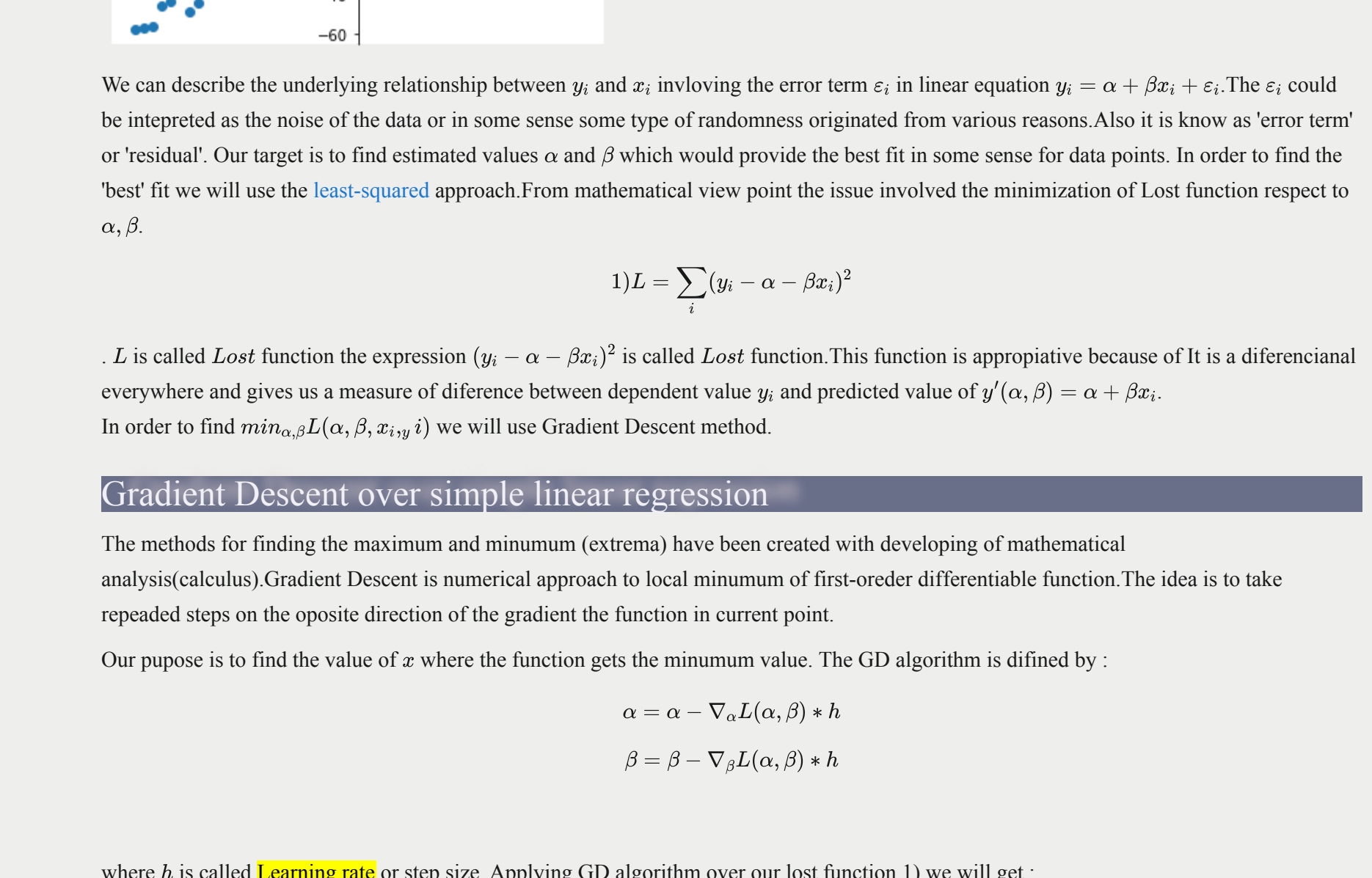
In statistic, **LA** is an linear approach to modeling the relationship between one or more explanatory variables (independent or dependent). The case of one explanatory variable is called Simple Linear Regression for more then one is called Multiple Linear Regression. In Linear Regression the relationships are modeled using **linear predictors function** whose unknown parameters are estimated from data.

## Simple linear regression

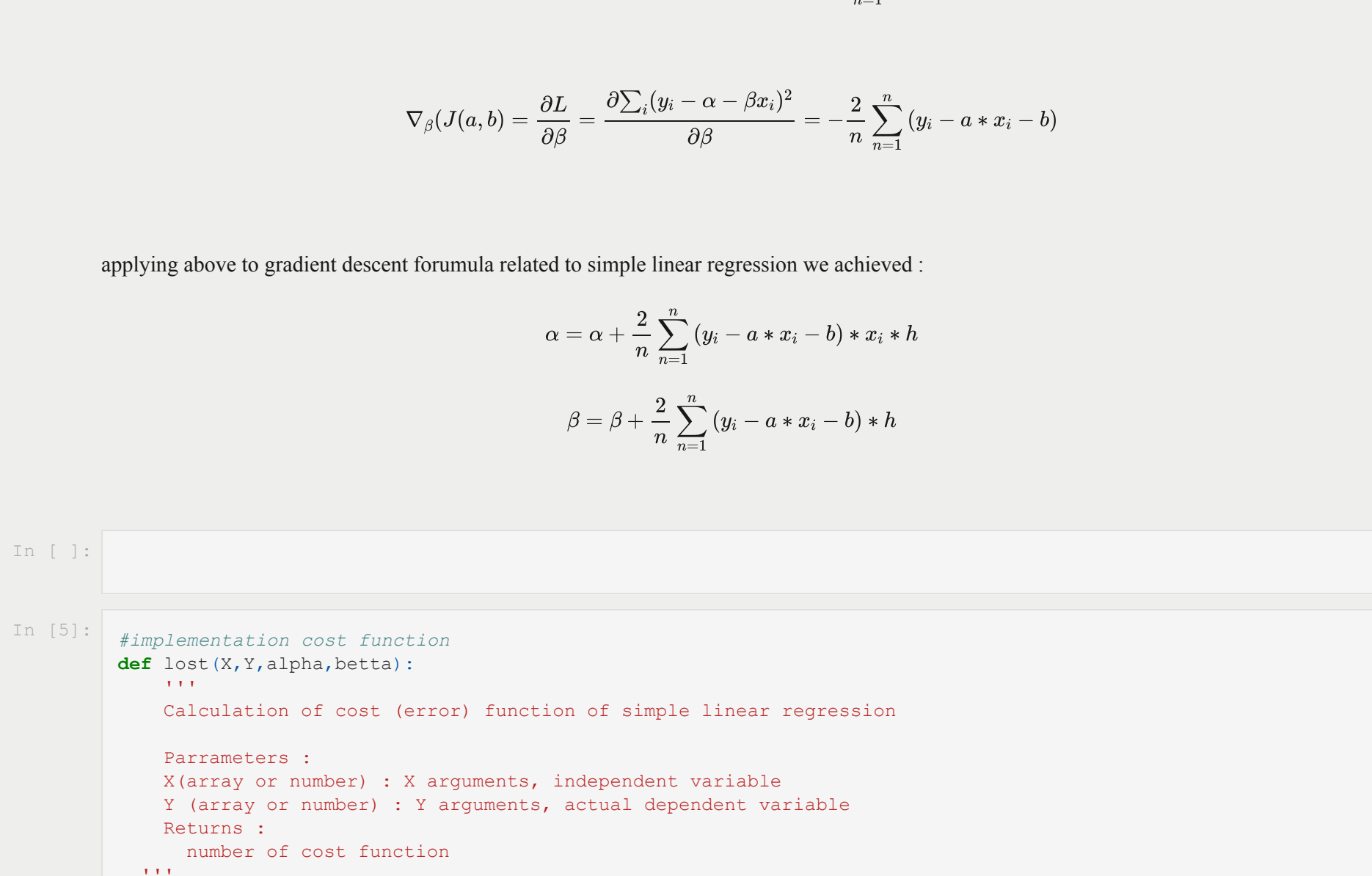
Simple linear regression has a single explanatory variable. It concerns two-dimensional sample points with one dependent and one independent variable. (Conventionally  $x$  and  $y$  data point in Cartesian coordinate system). The adjective simple refers to the fact that the outcome variable is related to a single predictor.

## Fitting the regression line with least-squared

Consider the model function  $y = \alpha + \beta x$  which describes a line with slope  $\beta$  and intercept  $\alpha$ .  
example: let  $\alpha = 7$  and  $\beta = 3$ , geometrically  $y = 7 + 3x$  appears to be as below graphic.



Suppose we observe  $n$  data pairs denoted with  $(x_i, y_i), i = 1, 2, \dots, n$ .



We can describe the underlying relationship between  $y_i$  and  $x_i$  involving the error term  $\epsilon_i$  in linear equation  $y_i = \alpha + \beta x_i + \epsilon_i$ . The  $\epsilon_i$  could be interpreted as the noise of the data or in some sense some type of randomness originated from various reasons. Also it is known as 'error term' or 'residual'. Our target is to find estimated values  $\alpha$  and  $\beta$  which would provide the best fit in some sense for data points. In order to find the 'best' fit we will use the **least-squared** approach. From mathematical view point the issue involved the minimization of *Lost* function respect to  $\alpha, \beta$ .

$$1) L = \sum_i (y_i - \alpha - \beta x_i)^2$$

$L$  is called *Lost* function the expression  $(y_i - \alpha - \beta x_i)^2$  is called *Lost* function. This function is appropriate because of It is a differential everywhere and gives us a measure of difference between dependent value  $y_i$  and predicted value of  $y(\alpha, \beta) = \alpha + \beta x_i$ . In order to find  $\min_{\alpha, \beta} L(\alpha, \beta, x_{1:n})$  we will use Gradient Descent method.

## Gradient Descent over simple linear regression

The methods for finding the maximum and minimum (extrema) have been created with developing of mathematical analysis (calculus). Gradient Descent is numerical approach to local minimum of first-order differentiable function. The idea is to take repeated steps on the opposite direction of the gradient the function in current point.

Our purpose is to find the value of  $x$  where the function gets the minimum value. The GD algorithm is defined by:

$$\alpha = \alpha - \nabla_{\alpha} L(\alpha, \beta) * h$$

$$\beta = \beta - \nabla_{\beta} L(\alpha, \beta) * h$$

where  $h$  is called **Learning rate** or step size. Applying GD algorithm over our lost function 1) we will get:

$$\nabla_{\alpha} (J(\alpha, b)) = \frac{\partial L}{\partial \alpha} = \frac{\partial \sum_i (y_i - \alpha - \beta x_i)^2}{\partial \alpha} = -2 \sum_{i=1}^n (y_i - \alpha - \beta x_i) * x_i$$

$$\nabla_{\beta} (J(\alpha, b)) = \frac{\partial L}{\partial \beta} = \frac{\partial \sum_i (y_i - \alpha - \beta x_i)^2}{\partial \beta} = -2 \sum_{i=1}^n (y_i - \alpha - \beta x_i) * h$$

applying above to gradient descent formula related to simple linear regression we achieved:

$$\alpha = \alpha + \frac{2}{n} \sum_{i=1}^n (y_i - \alpha - \beta x_i) * x_i * h$$

$$\beta = \beta + \frac{2}{n} \sum_{i=1}^n (y_i - \alpha - \beta x_i) * h$$

```
In [5]: #implementation cost function
def lost(X,Y,alpha,beta):
    """
    Calculation of cost (error) function of simple linear regression

    Parameters :
    X(array or number) : X arguments, independent variable
    Y (array or number) : Y arguments, actual dependent variable
    Returns :
    number of cost function
    """
    return np.sum((Y - (alpha + beta*X))**2 )

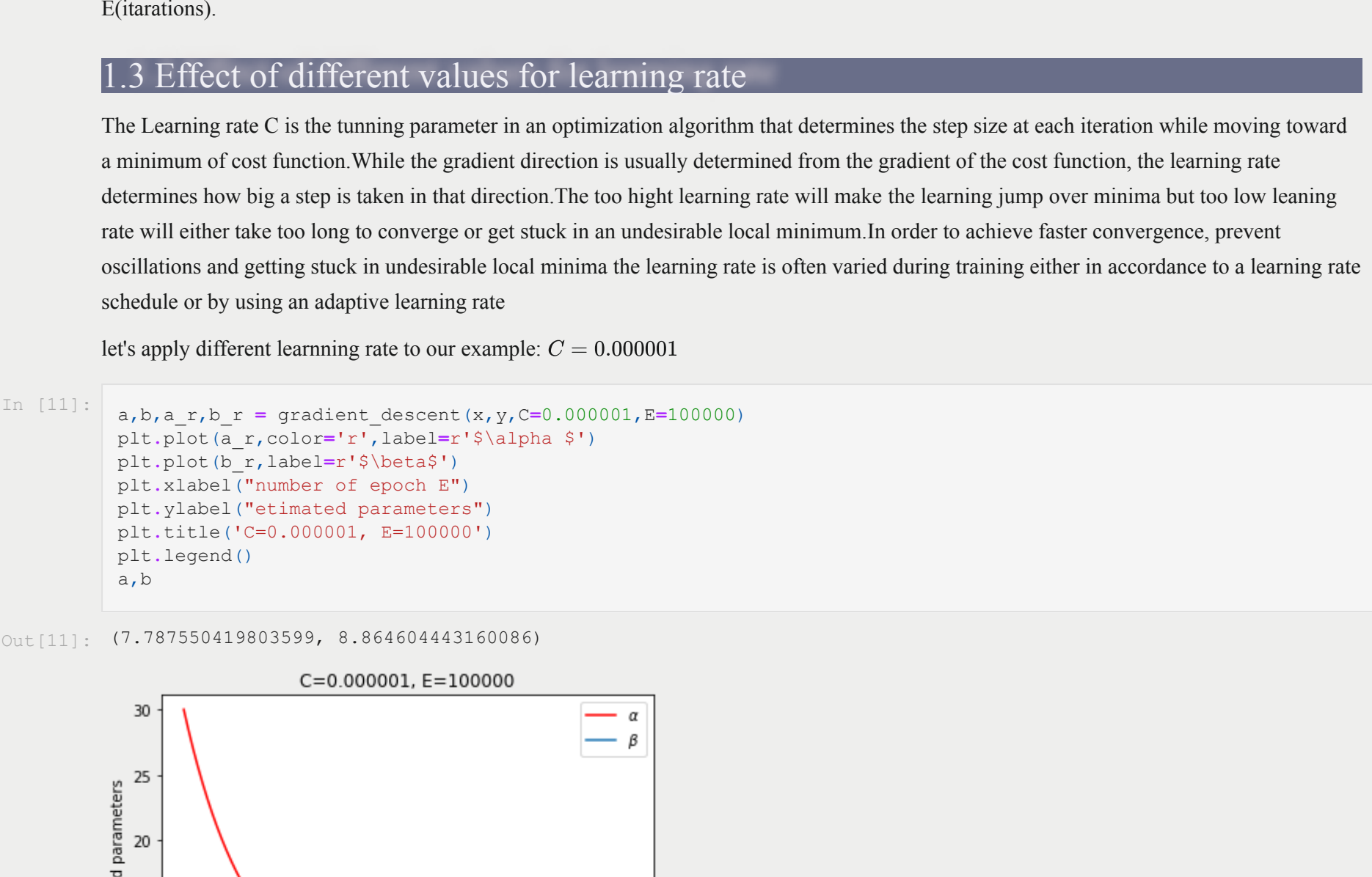
In [6]: def gradient_descent(X,Y,C=0.1,E=100):
    """
    Calculation of cost (error) function of simple linear regression

    Parameters :
    X(array or number) : X arguments, independent variable
    Y (array or number) : Y arguments, actual dependent variable
    C (number) : learning rate
    E (float) : number of iteration (Epoch)
    Returns :
    number of cost function
    """
    a = 30
    b = 10
    a_args = []
    b_args = []
    n = X.size

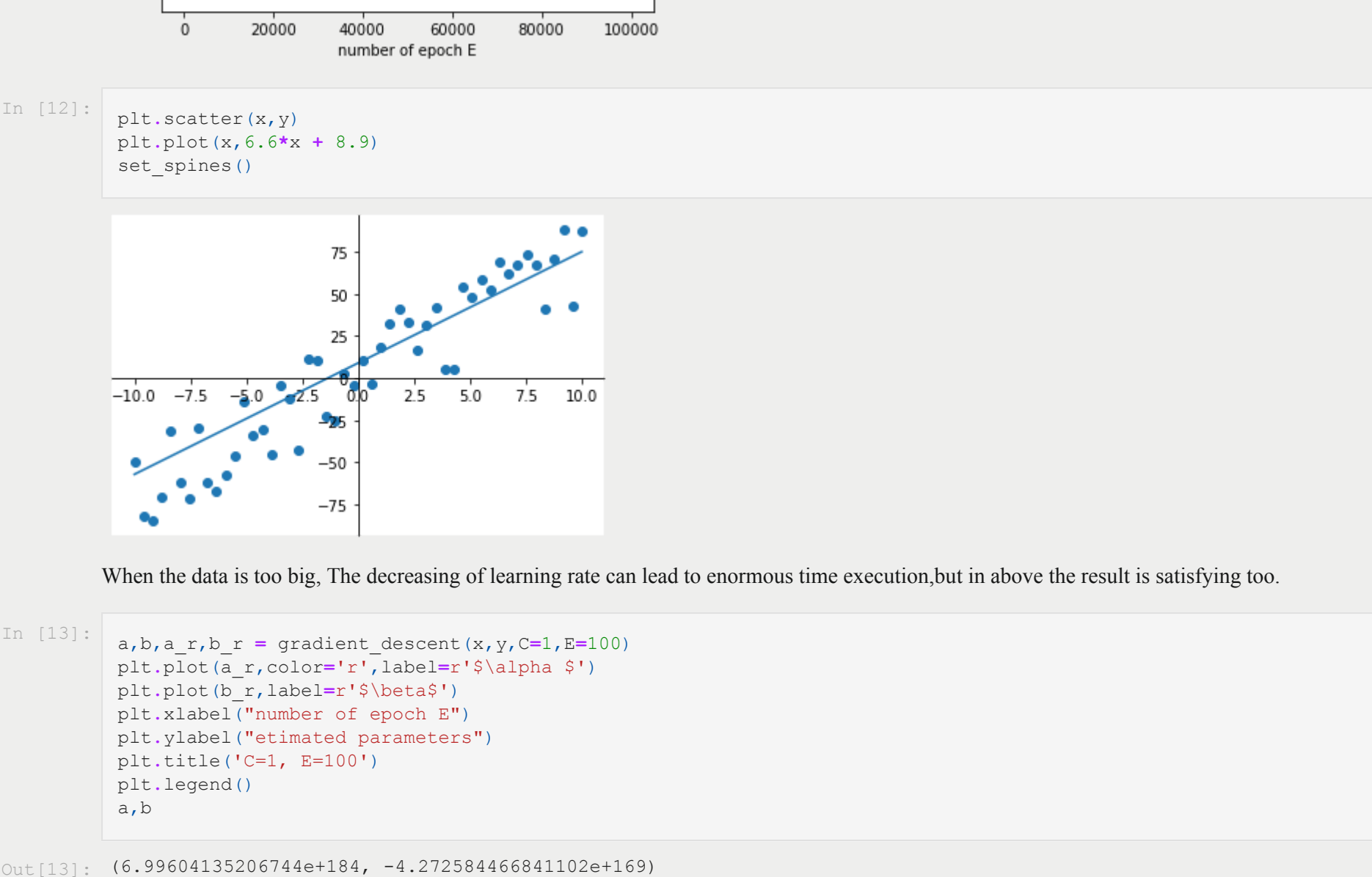
    for i in range(E):
        a = a + 2/n*(np.sum((Y - a*X - b)*X))*C
        b = b + 2/n*(np.sum((Y - a*X - b)))*C
        a_args.append(a)
        b_args.append(b)
    return a,b,a_args,b_args

In [7]: x = np.linspace(-10,10)
f = lambda x : 7*x + 3 + np.random.uniform(-30,30)
y = [f(i) for i in x]
plt.scatter(x,y)

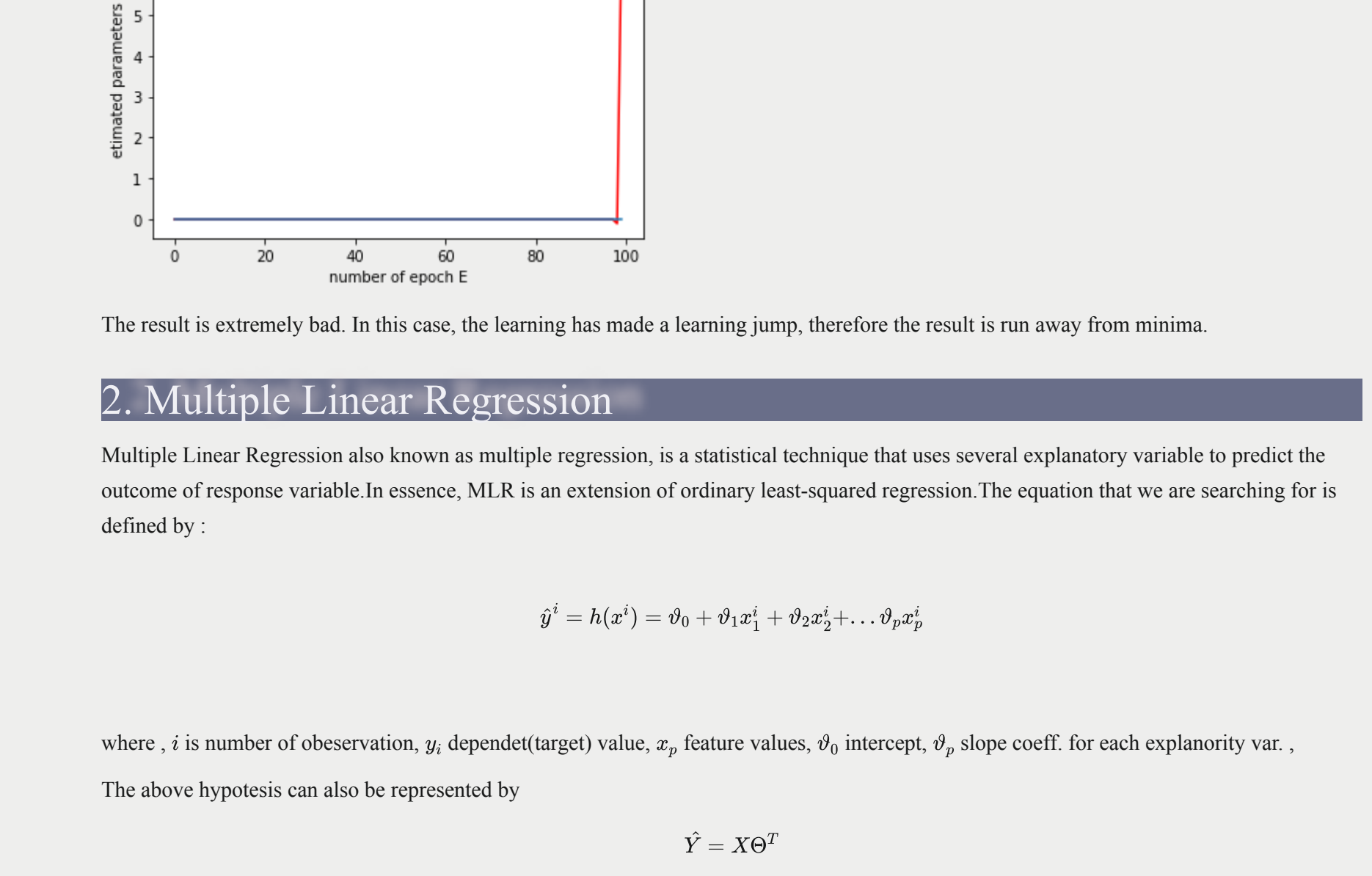
Out [7]: <matplotlib.collections.PathCollection at 0x5af59aa560>
```



Let to apply gradient descent of above points.



From graphics we can see that the curves tend to result ( $\alpha = 6.6603, \beta = 4.431$ ).  
The fit line will be:

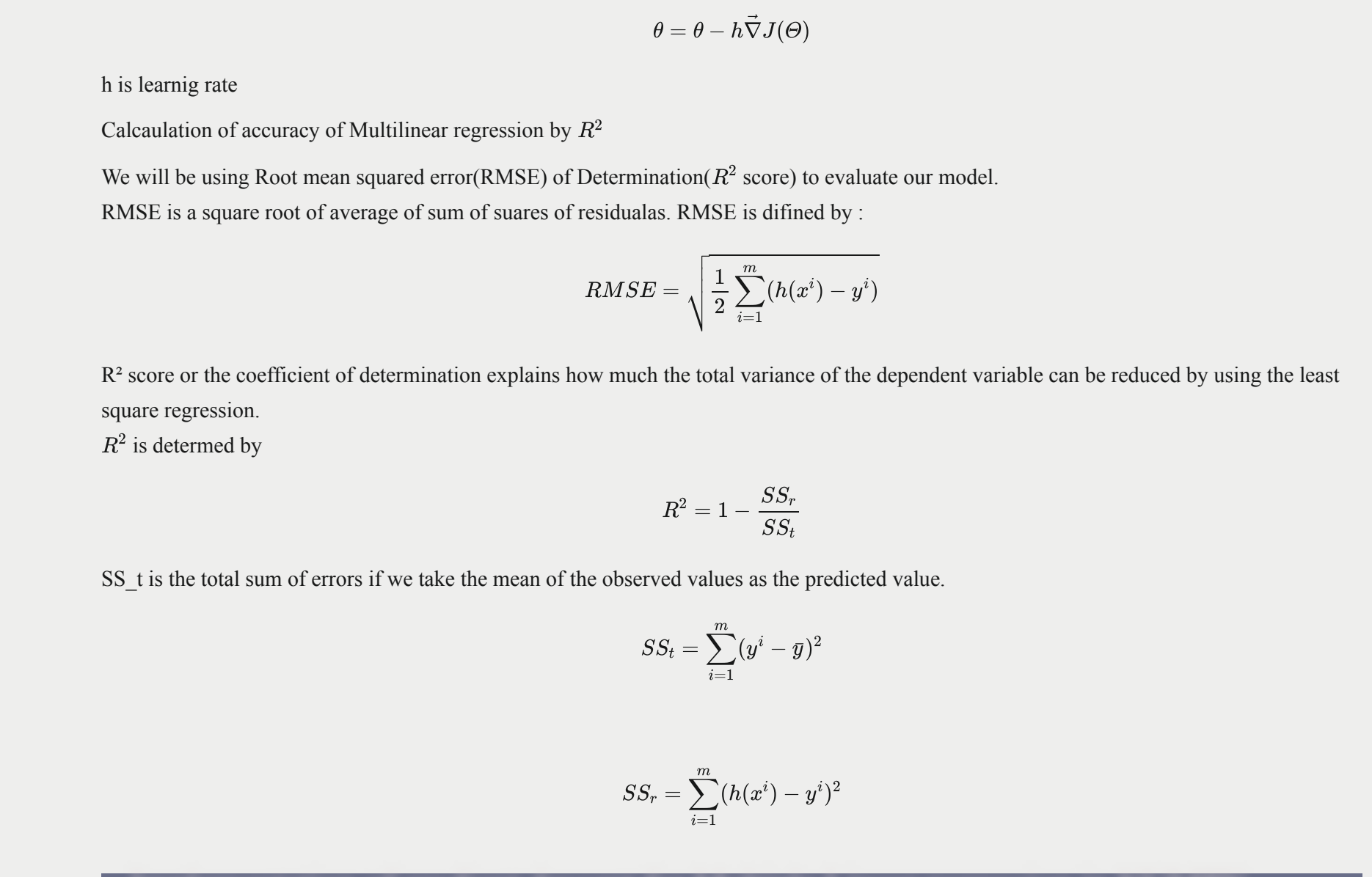
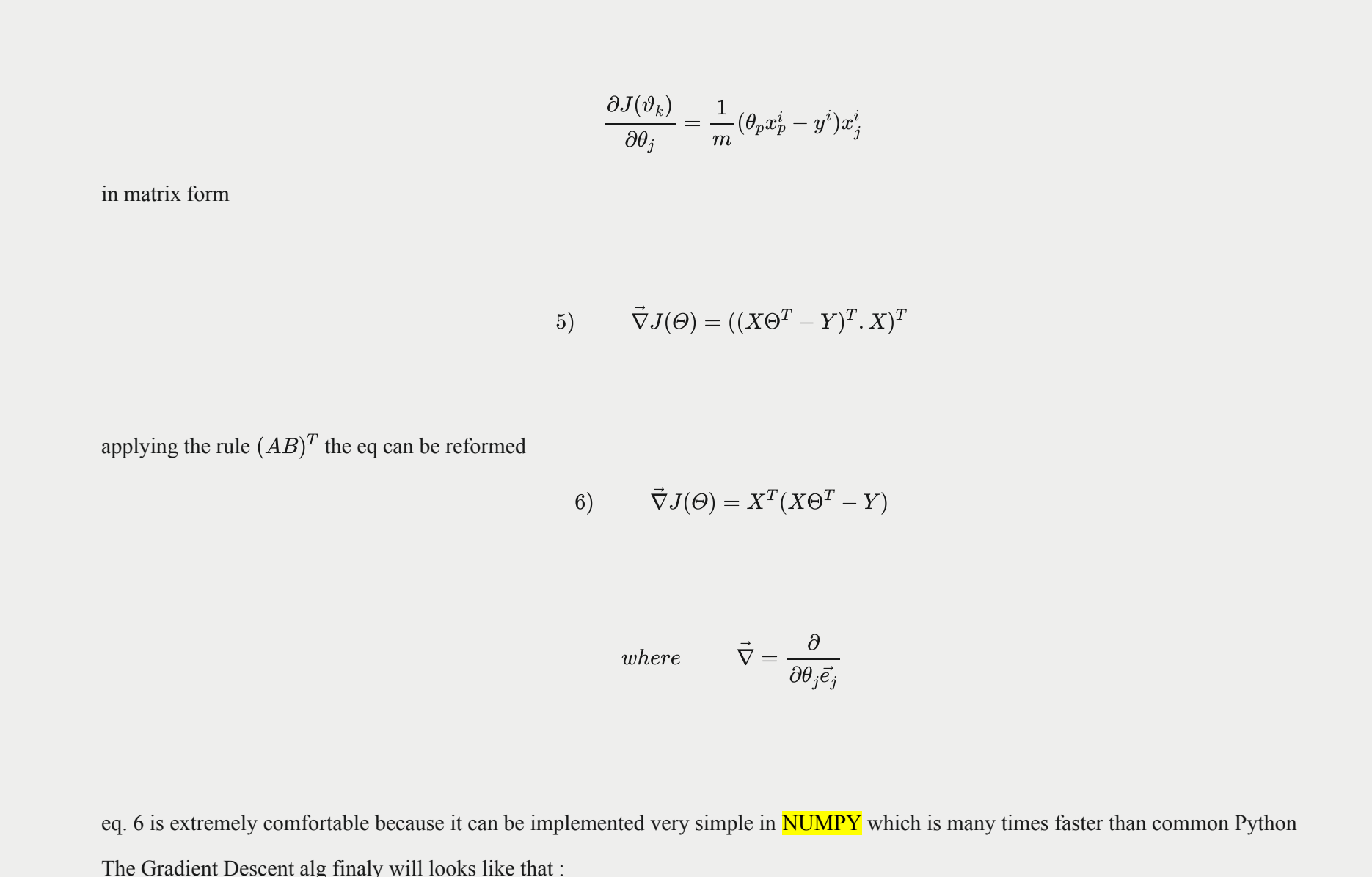


It seems to be somehow appropriate line which describe the data well, but our origin values of \$5 are ( $\alpha = 7, \beta = 1$ ) not ( $\alpha = 6.6603, \beta = 4.431$ ). In below we will examine how will chaged the ( $\alpha, \beta$ ) respect to different Learning rate  $C$  and nubor of epoch  $E$  (iterations).

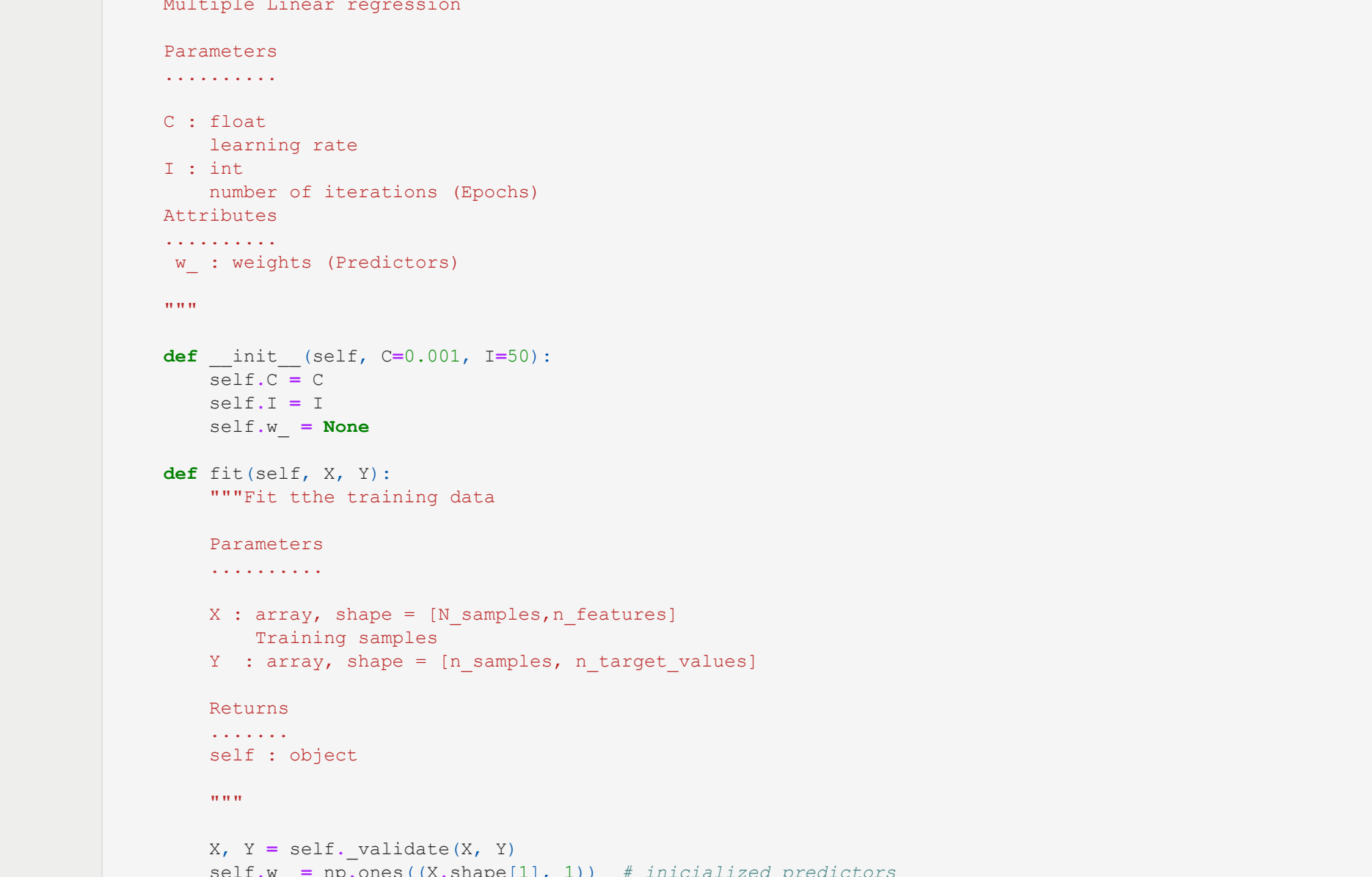
## 1.3 Effect of different values for learning rate

The Learning rate  $C$  is the tuning parameter in an optimization algorithm that determines the step size at each iteration while moving toward a minimum of cost function. While the gradient direction is usually determined from the gradient of the cost function, the learning rate determines how big a step is taken in that direction. The too high learning rate will make the learning jump over minima but too low learning rate will either take too long to converge or get stuck in an undesirable local minimum. In order to achieve faster convergence, prevent oscillations and getting stuck in undesirable local minima the learning rate is often varied during training either in accordance to a learning rate schedule or by using an adaptive learning rate.

let's apply different learning rate to our example:  $C = 0.000001$



When the data is too big, The decreasing of learning rate can lead to enormous time execution, but in above the result is satisfying too.



The result is extremely bad. In this case, the learning has made a learning jump, therefore the result is run away from minima.

## 2. Multiple Linear Regression

Multiple Linear Regression also known as multiple regression, is an ordinary technique that uses several explanatory variable to predict the outcome of response variable. In essence, MLR is an extension of statistical least-squared regression. The equation that we are searching for is defined by:

$$\hat{y}^i = h(x^i) = \theta_0 + \theta_1 x_1^i + \theta_2 x_2^i + \dots + \theta_p x_p^i$$

where,  $i$  is number of observation,  $y_i$  dengetal (target) value,  $x_p$  feature values,  $\theta_0$  intercept,  $\theta_p$  slope coeff. for each explanatory var.

The above hypothesis can also be represented by

$$\hat{Y} = X\Theta^T$$

where  $\Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}$ ;  $X = \begin{bmatrix} 1 & x_1^1 & x_1^2 & \dots & x_1^p \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_p^1 & x_p^2 & \dots & x_p^p \end{bmatrix}$  and  $\hat{Y} = \begin{bmatrix} \hat{y}^1 \\ \hat{y}^2 \\ \vdots \\ \hat{y}^p \end{bmatrix}$

we've append column  $\hat{Y} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$  to  $X$  in order to be used in matrix multiplication directly.

To define and measure the error of our model we define the cost function as the sum of the squares of the residuals. The cost function is denoted by:

$$3) \quad J(\Theta) = \frac{1}{2m} \sum_{i=1}^m (h(x^i) - y^i)^2$$

We have to initialize the model parameter with some random values (random initialization). To use Gradient Descent we need to measure how the cost function changes with change in its parameters. Therefore we compute the partial derivatives of cost function

$$4) \quad J(\theta_0, \theta_1, \dots, \theta_n)$$
$$\frac{\partial J(\Theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h(x^i) - y^i) x_j^i$$

In more computable form using matrix in order to be implemented using **NUMPY**

Using the **Eisten notation** we can rewrite eq. 4)

$$\frac{\partial J(\theta_k)}{\partial \theta_j} = \frac{1}{m} (\theta_p x_p^j - y^j) x_j^i$$

in matrix form

$$5) \quad \vec{\nabla} J(\Theta) = ((X\Theta^T - Y)^T \cdot X)^T$$

applying the rule  $(AB)^T$  the eq can be reformed

$$6) \quad \vec{\nabla} J(\Theta) = X^T (X\Theta^T - Y)$$

where  $\vec{\nabla} = \frac{\partial}{\partial \theta_j}$

eq. 6 is extremely comfortable because it can be implemented very simple in **NUMPY** which is many times faster than common Python

The Gradient Descent algorithm finally will look like that:

$$\theta = \theta - h \vec{\nabla} J(\Theta)$$

**h** is learning rate

Calculation of Accuracy from Squared error by Determination of  $R^2$

We will be using Root mean Squared error (RMSE) of Determination ( $R^2$  score) to evaluate our model.

RMSE is a square root of average of sum of squares of residuals. RMSE is defined by:

$$RMSE = \sqrt{\frac{1}{2} \sum_{i=1}^n (h(x^i) - y^i)^2}$$

$R^2$  score or the coefficient of determination explains how much the total variance of the dependent variable can be reduced by using the least square regression.

$R^2$  is determined by

$$R^2 = 1 - \frac{SS_r}{SS_t}$$

$SS_t$  is the total sum of errors if we take the mean of the observed values as the predicted value.

$$SS_t = \sum_{i=1}^n (y^i - \bar{y})^2$$

$$SS_r = \sum_{i=1}^m (h(x^i) - y^i)^2$$

## Implementation of gradient descent for Multiple Linear regression in NUMPY

```
In [19]: class MultipleLinearRegression:
    """
    Multiple Linear regression

    Parameters
    .....
    C : float
        learning rate
    I : int
        number of iterations (Epochs)
    Attributes
    .....
    w_ : weights (Predictors)

    """
    def __init__(self, C=0.001, I=50):
        self.C = C
        self.I = I
        self.w_ = None

    def fit(self, X, Y):
        """fit the training data

        Parameters
        .....
        X : array, shape = [N_samples, n_features]
            Training samples
        Y : array, shape = [n_samples, n_target_values]

        Returns
        .....
        self : object
        """
        X, Y = self._validate(X, Y)
        self.w_ = np.ones((X.shape[0], 1)) # initialized predictors
        for i in range(self.I):
            # number of samples
            M = X.shape[0]
            self.w_ = self.w_ + self.C * (1 / M) * X.T.dot((X.dot(self.w_) - Y))
            #self.w_ = self.C * (1 / M) * X.T.dot((X.dot(self.w_) - Y))

        return self

    def predict(self, X):
        """Predicts the value after the model has been trained.

        Parameters
        .....
        X : array-like, shape = [n_samples, n_features]
            Test samples

        Returns
        Predicted value
        """
        n = np.append(np.ones((X.shape[0], 1)), X, axis=1)
        return np.dot(n, self.w_)

    def score(self, X, Y):
        """Calculation of accuracy using (R^2 score)

        param y: array-like, shape = [n_samples, n_features]
        return: validated X, y
        """
        n = np.append(np.ones((X.shape[0], 1)), X, axis=1)
        y_prime = np.dot(n, self.w_)
        ssr = np.sum((y_prime - Y)**2)
        sst = np.sum((Y - np.mean(Y))**2)
        r2_score = 1 - (ssr / sst)
        return r2_score

    def _validate(self, X, Y):
        """Added tow with one number to X data
        reshape Y data

        param X: array-like, shape = [n_samples, n_features]
        param y: array-like, shape = [n_samples, n_target_values]
        return: validated X, y
        """
        n = np.append(np.ones((X.shape[0], 1)), X, axis=1)
        n = np.array(y.reshape(y.shape[0], 1))

We will test our implementation over insurance.csv data set
```



Converting Categories to Numbers. The linear regression can be performed only on numbers, so we should convert these categorical features into numbers. To do that, we can make use of a function called **get\_dummies**. So let's convert the "sex," "smoker," and "region" columns into numerically represented features.

```
In [22]: cols = ['sex', 'smoker', 'region']
new_df = pd.get_dummies(df, cols, drop_first=True)
new_df.head()
```

age	bmi	children	charges	sex_male	smoker_yes	region_northwest	region_southeast	region_southwest
0	19	27900	0	16884.92400	0	1	0	0
1	18	33770	1	1725.55230	1	0	0	1
2	28	33000	3	4449.46200	1	0	0	1
3	33	22705	0	21984.47061	1	0	1	0
4	32	28880	0	3866.85520	1	0	1	0

Now, let's only select the features that are the most relevant. Feature selection is one of the important tasks in any machine learning project. You must know which features are most correlated with the targets (the "charges" column in our case) and must only use those features that have a high correlation with your target. This can be done through experimentation. For example, in this problem, I tried using the "sex" and "region" features to predict "charges" but didn't find much of an improvement in the prediction performance of the model. So I decided to omit these features from the model. Through small experimentation like that, I found the "age," "bmi," and "smoker" columns to be most relevant when predicting insurance costs (the "charges" column in our data frame).

```
In [24]: X = new_df[['age', 'bmi', 'smoker_yes']]
y = new_df['charges']

We will perform standardization over X

In [28]: scaler = StandardScaler()

In [29]: X = scaler.fit_transform(X)

Splitting data to Train and Test

In [30]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2)

End let to train the data using our implementation

In [32]: model = MultipleLinearRegression(C=0.01, I=1000)
z = model.fit(X_train, y_train)

In [ ]:

In [37]: model.score(X_test, y_test)

Out [37]: 0.7323173758003395

We've achieved according to R2 score 73% accuracy

Let see how will be in standard LinearRegression in sklearn in Python

In [41]: lr = LinearRegression(fit_intercept=True)
lr.fit(X_train, y_train)

Out [41]: LinearRegression()

In [42]: lr.score(X_test, y_test)

Out [42]: 0.7322470193784116

The score using implementation in sklearn in Python is the same as our implementation. Let to compare the predicted values

In [92]: our_predicted_data = model.predict(X_test)[:,0]
python_predicted_data = lr.predict(X_test)
pd.DataFrame(our_predicted_data, python_predicted_data )

Out [92]:
```

	0
26963.233582	26962.539696
32466.037854	32465.293221
16386.515466	16386.117349
10313.491299	10312.927540
4554.159257	4553.882338
...	...
6066.659315	6066.487660
1301.918644	1301.573708
9178.571257	9178.181162
6629.357965	6629.146516
28596.837404	28596.167084

268 rows x 1 columns

The differences is too small, but we can perform  $R^2$  score over the both Python and Our predicted data

```
In [93]: ssr = np.sum((our_predicted_data - python_predicted_data)**2)
sst = np.sum((our_predicted_data - np.mean(python_predicted_data))**2)
r2_score = 1 - (ssr / sst)

Out [93]: 0.9999999976969305

0.9999999976969305 shows that the difference is negligible
```