

OPTIMIZATION IN BIG DATA RESEARCH
GRADIENT DESCENT METHOD REPORT

Zixuan Li

student number 3118103163

professor: Jianyong Sun

December 12, 2018

Contents

| | | |
|----------|--|-----------|
| 1 | research background | 2 |
| 1.1 | research purpose | 2 |
| 1.2 | Compile environment | 4 |
| 1.3 | Procedure | 4 |
| 2 | Gradient Descent Method | 5 |
| 2.1 | algorithm introduction | 5 |
| 2.2 | result | 5 |
| 2.3 | source code | 6 |
| 3 | Momentum Method | 7 |
| 3.1 | algorithm introduction | 7 |
| 3.2 | result | 7 |
| 3.3 | source code | 7 |
| 4 | Nesterov accerlerated gradient Method | 9 |
| 4.1 | algorithm introduction | 9 |
| 4.2 | result | 9 |
| 4.3 | source code | 9 |
| 5 | Adaptive Moment Estimation Method | 11 |
| 5.1 | algorithm introduction | 11 |
| 5.2 | result | 11 |
| 5.3 | source code | 12 |

1 research background

1.1 research purpose

The purpose of this experiment is to optimize Beale function by using gradient descent, momentum, Nesterov accelerated gradient(NAG), and Adaptive Moment Estimation (Adam) algorithms.

According to the definition by Khon nen, neural network is an extensive parallel interconnected network composed of simple units with adaptability, and its organization can simulate the interaction of biological nervous system with real world objects [Koho nen, 1988]. In machine learning, the neural network means neural network leaning which is an interdisciplinary field of study.

The most basic component of the neural network is the neuron model. In the neural network, each neuron is connected to other neurons, and when it is stimulated, it will send signal to the other connected neurons. And the connection have different weight, thus has different impact on the other neurons. The total input value received by the neuron will be compared with the threshold value of the neuron and then processed by the activation function to produce the output of the neuron. Training neural network is to give a training data set, and continuously adjust the transmission weight and activation function in the neural network to reduce the gap between the predicted value and the data set. If the error of the neural network in the training set is expressed by ϵ and it is obviously a function of the connection weight ω and the threshold value θ . At this time, the training process of the neural network can be regarded as a parameter optimization process. That is, in the parameter space, finding a set of optimal parameters makes ϵ minimum.

Gradient descent is the most widely used parameter optimization method, in which we start from some initial solutions. In each iteration, we first calculate the gradient of the error function $\epsilon(x)$ at the current point and then determine the search direction according to the gradient. The learning rate η determines the step of each iteration. Many deep learning project use various algorithms to optimize gradient descent, but in these large models, the optimization algorithm is often a black box model. Therefore, the purpose of this experiment is to optimize a test function by using several commonly used gradient descent algorithms, and to gain a more intuitive understanding of several gradient descent algorithms by writing code and visualizing the descent process.

The test function used in this research is Beale function:

$$f(x, y) = (1.5 - x + xy)^2 + (2.25 - x + xy^2)^2$$

Global minimum is $f(3, 0.5) = 0$ in search domain $-4.5 \leq x, y \leq 4.5$.

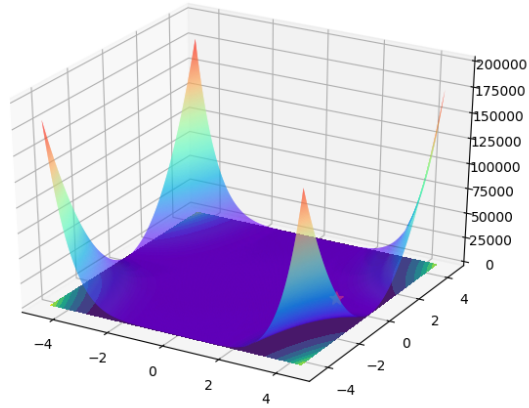


Figure 1: 3d plot of Beale function

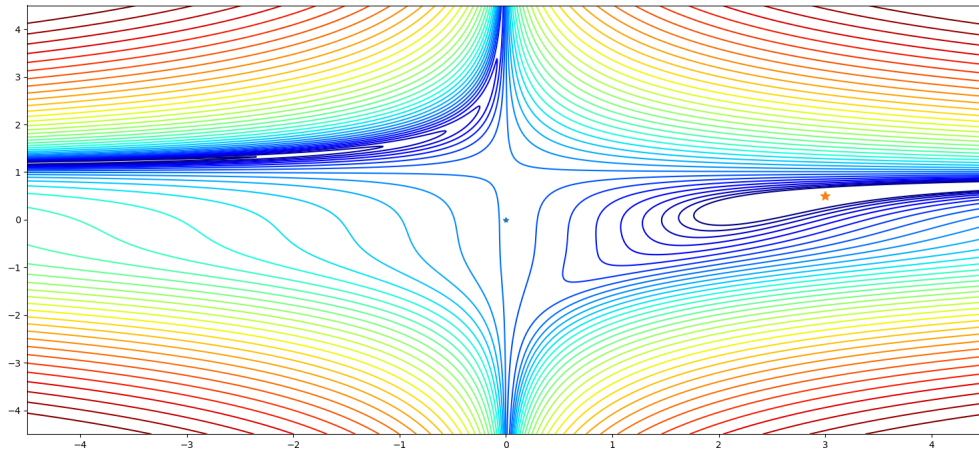


Figure 2: contour plot of Beale function and global minimum

In this paper, Gradient descent, Momentum, Nesterov accelerated gradient, and Adaptive Moment Estimation will be used to solve the global lowest point of Beale Function. The characteristics of various gradient descent methods are known through experiments. The report will consist of the following five main parts:

- Gradient descent method
- Momentum methods
- Nesterov accelerated gradient(NAG) method
- Adaptive Moment Estimation(ADAM) methon
- conclusion and Discussion

1.2 Compile environment

The Python language is used for experiments. The relevant compilation environment is as follows:

- OS Windows 10
- Compiler python3.6
- numpy scientific computing package
- matplotlib python visualize package

1.3 Procedure

1. programming to visualize Beale function.
2. programming the four algorithm.
3. Adjust parameters in each algorithms, such as learning rate, to achieve fastest convergence rate.
4. Run the program, examine the route of convergence to gain insight into each algorithm.
5. Output related graphs and data
6. Finish the experiment report.

2 Gradient Descent Method

2.1 algorithm introduction

The gradient descent method is very easy to understand. The gradient direction indicates the direction where the function grows fastest, so the opposite direction is the direction where the function decreases fastest. For the problem of machine learning model optimization, when we need to solve the minimum value, we can find the optimal value by moving in the direction of gradient descent.

The main process of GD algorithm can be summarized as follows:

1. input: Function to be solved $f(x, y)$, derivative of function to be solved $f(x, y)'$, start point $P(x_0, y_0)$, the step size along the gradient descent direction $step$
2. for $(abs(grad) > 1e - 6)$ do
3. compute the derivative of current point $P(x, y)$
4. $gradient = -f'(P_t)$
5. updated the next point as $P_{t+1} = P_t - step * gradient$
6. end
7. output: the convergence point $P_{convergence}$

2.2 result

The start point as the exercise requirement is set at $P(3, 4)$. As for the step size, If it is less than $1/L$. Then it can guarantee convergence (L is the Lipschitz constant of the gradient of the objective function).

$$|f(x_1) - f(x_2)| = |f'(\xi)(x_1 - x_2)| \leq L|x_1 - x_2|$$

I computed the numerical solution of Beale function's Lipschitz constant. And the result is $L = 248533.18$. so as long as $dp \leq 4 \times 10^{-6}$ the algorithm will convergence. After weighing convergence and convergence speed, I finally choose $step = 5 \times 10^{-6}$.

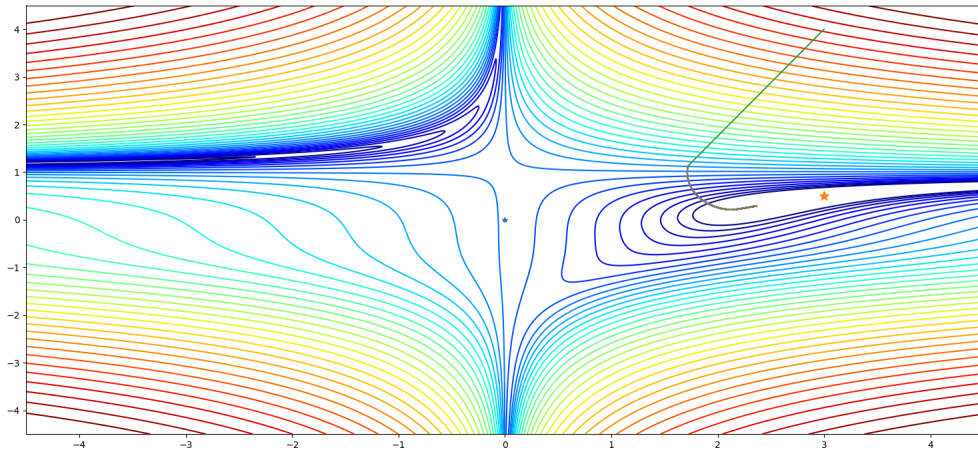


Figure 3: gradient descent process

2.3 source code

The part of source code of this algorithm is as follows.

```
1 def gd(x_start, step, g): # Gradient Descent
2     x = np.array(x_start, dtype='float64')
3     passing_dot = [x.copy()]
4     for i in range(n_iter):
5         grad = g(x)
6         x -= grad * step
7         passing_dot.append(x.copy())
8         if abs(sum(grad)) < 1e-6:
9             break;
10    return x, passing_dot
```

GD.py

3 Momentum Method

3.1 algorithm introduction

As its name suggests, momentum acts as a constant catalyst for the previous optimization in the process of optimization. The influence of an already declining direction will not disappear immediately, but will decay in a certain form.

The main process of Momentum algorithm can be summarized as follows:

1. input: Function to be solved $f(x, y)$, derivative of function to be solved $f(x, y)'$, start point $P(x_0, y_0)$, the step size along the gradient descent direction $step$, momentum term γ
2. for $(abs(grad) > 1e-6)$ do
3. compute the derivative of current point x
4. compute update vector $v_t = \gamma v_{t-1} + step f(x, y)'$
5. updated the next point as $P_{t+1} = P_t - v_t$
6. end
7. output: the convergence point $P_{convergence}$

The momentum term γ is usually set to 0.9.

3.2 result

The start point as the exercise requirement is set at $P(3, 4)$. As for the step size, After weighing convergence and convergence speed, I finally choose $step = 5 \times 10^{-7}$. To restrict the iteration time, the program is set to stop when it iterates 10000 times.

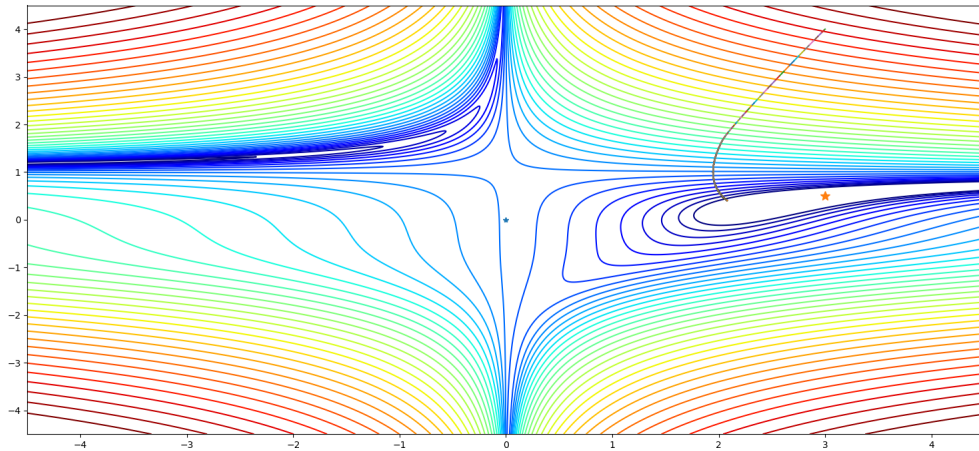


Figure 4: momentum descent process

3.3 source code

The part of source code of this algorithm is as follows.


```

1 def momentum(x_start, step, g, discount = 0.9): # momentum
2     x = np.array(x_start, dtype='float64')
3     passing_dot = [x.copy()]
4     pre_grad = np.zeros_like(x)
5     for i in range(n_iter):
6         grad = g(x)
7         pre_grad = pre_grad * discount + grad
8         x -= pre_grad * step
9
10        passing_dot.append(x.copy())
11        if abs(sum(grad)) < 1e-6:
12            break;
13    return x, passing_dot

```

momentum.py

4 Nesterov accerlerated gradient Method

4.1 algorithm introduction

The momentum method can be thought as a ball rolling down a hill with momentum. Then the NAG method is smarter ball which can slowing down before climbing up another hill.

The main process of Momentum algorithm can be summarized as follows:

1. input: Function to be solved $f(x, y)$, derivative of function to be solved $f(x, y)'$, start point $P(x_0, x_0)$, the step size along the gradient descent direction $step$, momentum term γ
2. for $(abs(grad) > 1e-6)$ do
3. compute the derivative of current point $f(P_t)'$
4. compute update vector $v_t = \gamma v_{t-1} + step P_t - \gamma v_{t-1}'$
5. updated the next point as $P_{t+1} = P_t - v_t$
6. end
7. output: the convergence point $P_{convergence}$

The difference between Nag and Momentum is that, NAG use the momentum of present point to get a proximation of the next gradient rather than using the present gradient directly as Momentum method.

4.2 result

The start point as the exercise requirement is set at $P(3, 4)$. The momentum term γ is also usually set to 0.9. The step size is same as momentum method $step = 5 \times 10^{-7}$. To restrict the iteraion time, the program is set to stop when it ierates 10000 times.

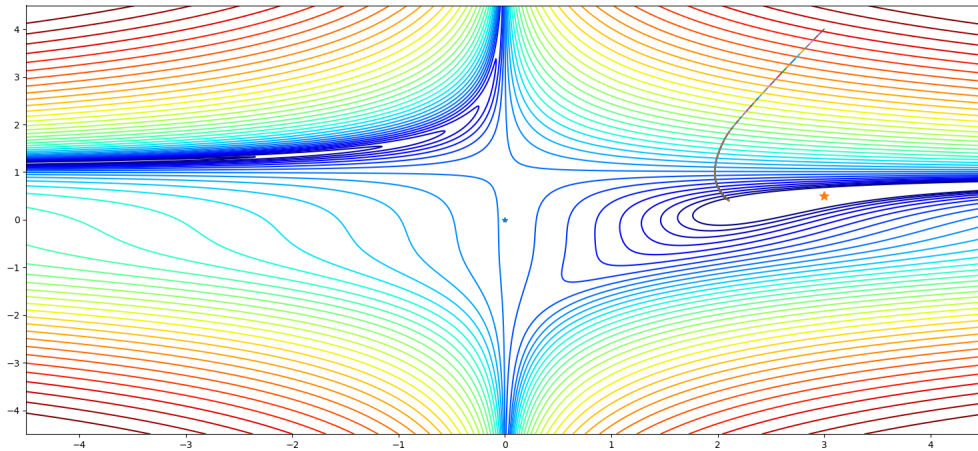


Figure 5: NAG descent process

4.3 source code

The part of source code of this algorithm is as follows.

```

1 def nesterov(x_start, step, g, discount=0.9):    #Nesterov accelerated gradient
2     x = np.array(x_start, dtype='float64')
3     passing_dot = [x.copy()]
4     pre_grad = np.zeros_like(x)
5     for i in range(n_iter):
6         x_future = x - step * discount * pre_grad
7         grad = g(x_future)
8         pre_grad = pre_grad * discount + grad
9         x -= pre_grad * step
10
11         passing_dot.append(x.copy())
12         if abs(sum(grad)) < 1e-6:
13             break;
14     return x, passing_dot
15
16 i

```

nag.py

5 Adaptive Moment Estimation Method

5.1 algorithm introduction

Adaptive Moment Estimation (Adam) is a method that computes adaptive learning rates for each parameter. Adaptive combines the advantages of RMSprop and Momentum method. It stores the an exponentially decaying average of past squared gradients v_t and exponentially decaying average of past gradients m_t .

The main process of Momentum algorithm can be summarized as follows:

1. input: Function to be solved $f(x, y)$, derivative of function to be solved $f(x, y)'$, start point $P(3, 4)$, the step size along the gradient descent direction $step$, super-parameter β_1, β_2 .
2. for $(abs(grad) > 1e - 6)$ do
3. compute the gradient of current point $g_t = -f(P_t)'$
4. compute the decaying averages of past gradient $m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$
5. compute the decaying averages of past squared gradient $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$
6. computing bias-corrected first moment estimates $\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$
7. computing bias-corrected second moment estimates $\hat{v}_t = \frac{m_t}{1 - \beta_2^t}$
8. compute update vector $v_t = \frac{step}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t$
9. updated the next point as $P_{t+1} = P_t - v_t$
10. end
11. output: the convergence point $P_{convergence}$

5.2 result

The default values proposed by the authors of this methd is $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$ And The start point as the exercise requirement is set at $P(3, 4)$. The step size is set to 1 by testing to have a relatively fast convergence rate. To restrict the iteraion time, the program is set to stop when it ierates 10000 times.

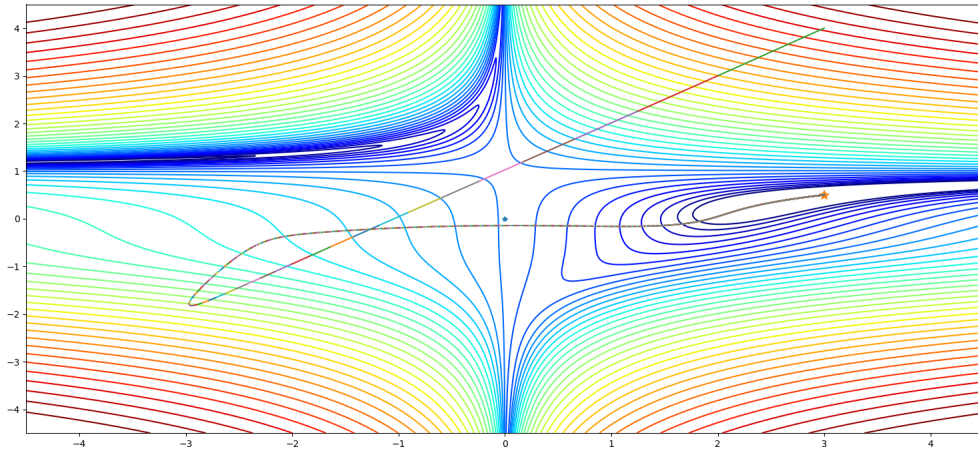


Figure 6: ADAM descent process

We can see that the Adam algorithm behaves like a heavy ball with friction because it stores an exponentially decaying average of past gradients.

5.3 source code

The part of source code of this algorithm is as follows.

```
1 def adam(x_start, step, g):                                #Adaptive Moment Estimation
2     x = np.array(x_start, dtype='float64')
3     beta1=0.9
4     beta2=0.999
5     mt = np.zeros_like(x)
6     vt = np.zeros_like(x)
7     dp = np.zeros_like(x)
8     passing_dot = [x.copy()]
9     for i in range(n_iter):
10         beta1t=beta1**(i+1)
11         beta2t=beta2**(i+1)
12         temp=g(x)
13         mt=beta1*mt + (1-beta1) * g(x)
14         vt=beta2*vt + (1-beta2) * g(x)**2
15         mt1=mt/(1-beta1t)
16         vt1=vt/(1-beta2t)
17         dp[0]= step * mt1[0]/(math.sqrt(vt1[0])+0.00000001)
18         dp[1] = step * mt1[1] / (math.sqrt(vt1[1]) + 0.00000001)
19         x -= dp
20         passing_dot.append(x.copy())
21         if abs(sum(mt1)) < 1e-6:
22             break;
23     return x, passing_dot
```

adam.py

List of Figures

| | | |
|---|---|----|
| 1 | 3d plot of Beale function | 3 |
| 2 | contour plot of Beale function and global minimum | 3 |
| 3 | gradient descent process | 5 |
| 4 | momentum descent process | 7 |
| 5 | NAG descent process | 9 |
| 6 | ADAM descent process | 11 |