# MATLAB Deep Learning Notes VI

**Key Words : Momentum**

## 1. Momentum

1. In fact, there are various weight adjustment forms available. The simplest adjustment form such as equation (9) and (19) doesn't have higher stability or faster speed. So, we should find some more effective weight adjustment forms. In this notes we introduce **Momentum**. The basic concept of Momentum is (20) as follows

$$
\begin{aligned}
\Delta w &= \alpha \delta x \\
m &= \Delta w + \beta m^- \\
w &= w + m \\
m^- &= m
\end{aligned}
$$

where $m^-$ is the previous momentum and $\beta$ is a positive constant that is less than 1.

2. The following steps (21) show how the momentum changes over time.

$$
\begin{aligned}
m(0) &= 0 \\
m(1) &= \Delta w(1) + \beta m(0) = \Delta w(1) \\
m(2) &= \Delta w(2) + \beta m(1) = \Delta w(2) + \beta \Delta w(1) \\
m(3) &= \Delta w(3) + \beta m(2) = \Delta w(3) + \beta \Delta w(2) + \beta^2 \Delta w(1) \\
m(4) &= \cdots
\end{aligned}
$$

3. It is noticeable from these steps that the previous weight update, i.e., $\Delta w(1)$, $\Delta w(2)$, $\Delta w(3)$, etc., is added to each momentum over the process. Since $\beta$ is less than 1, the older weight updates exerts a lesser influence on the momentum. Although the influence diminishes over time, the old weight updates remain in the momentum. Therefore, the weight is not solely affected by a particular weight update value, and the learning stability improves accordingly. In addition, since the momentum grows, the weight update becomes greater and greater as well.

**Program 12**: BackPropMmt

Listing 1: BackPropMmt.m

```matlab
function [weight1, weight2] = BackPropMmt(weight1, weight2, data_input,...
correct_output)

alpha = 0.9;
beta = 0.9;

mmt1 = zeros(size(weight1));
mmt2 = zeros(size(weight2));

N = 4;
for k = 1 : N
    x = data_input(k, :)';
    d = correct_output(k);

    v1 = weight1 * x;
    y1 = sigmoid(v1);
    v = weight2 * y1;
    y = sigmoid(v);

    e = d - y; % calculate the error of the output to the correct output
    delta = y .* (1 - y) .* e;

    e1 = weight2' * delta; % calculate the error of the hidden layer
    delta1 = y1 .* (1 - y1) .* e1;

    % Adjust the weights use momentum
    dw1 = alpha * delta1 * x';
    mmt1 = dw1 + beta * mmt1;
    weight1 = weight1 + mmt1;

    dw2 = alpha * delta * y1';
    mmt2 = dw2 + beta * mmt2;
    weight2 = weight2 + mmt2;
end
```

```
35 end
```

**Program 13**: testBackPropMmt

Listing 2: testBackPropMmt.m

```
1  clear
2  data_input = [ 0, 0, 1; 0, 1, 1; 1, 0, 1; 1, 1, 1]; % training data
3  correct_output = [0; 1; 1; 0]; % correct outputs(i.e., labels)
4  weight1 = 2 * rand(4, 3) - 1;
5  weight2 = 2 * rand(1, 4) - 1;
6
7
8  for epoch = 1:10000 % train
9  [weight1, weight2] = BackPropMmt(weight1, weight2, data_input,...
10 correct_output);
11 end
12
13 N = 4; % inference
14 for k = 1 : N
15 x = data_input(k, :)';
16 v1 = weight1 * x;
17 y1 = sigmoid(v1);
18 v = weight2 * y1;
19 y = sigmoid(v)
20 end
```

**Output 13**:

```
>> testBackPropMmt

y =

    0.0046

y =

    0.9940

y =

    0.9924

y =
```

```
    0.0119
```

## 2. Comparison between Simple Adjustment and Momentum

4. In this section we compare the average training error performance between Simple Adjustment and Momentum. We use SGD and aim to solve the linear separable problem we talked before.

**Program 14**: SimpleVsMmt

Listing 3: SimpleVsMmt.m

```matlab
1  clear
2
3  data_input = [ 0, 0, 1; 0, 1, 1; 1, 0, 1; 1, 1, 1 ]; % training data
4  correct_output = [ 0; 0; 1; 1 ]; % correct outputs
5  weight11 = 2 * rand(4, 3) − 1;
6  weight12 = 2 * rand(1, 4) − 1;
7  weight21 = weight11;
8  weight22 = weight12;
9
10 E1 = zeros(1000, 1);
11 E2 = E1;
12
13 for epoch = 1 : 1000 % the number of epoch = 1000
14     [weight11, weight12] =  BackPropXOR(weight11, weight12, data_input,
           correct_output);
15     [weight21, weight22] =  BackPropMmt(weight21, weight22, data_input,
           correct_output);
16
17     es1 = 0; es2 = 0;
18     N = 4;
19     for k = 1 : N
20         x = data_input(k, :)';
21         d = correct_output(k);
22         % Simple(i.e., Mmt)
23         v1_simple = weight11 * x;
24         y1_simple = sigmoid(v1_simple);
25         v_simple = weight12 * y1_simple;
26         y_simple = sigmoid(v_simple)
27         es1 = es1 + (d − y_simple)^2;
28         % Mmt
```

```matlab
29          v1_mmt = weight21 * x;
30          y1_mmt = sigmoid(v1_mmt);
31          v_mmt = weight22 * y1_mmt;
32          y_mmt = sigmoid(v_mmt)
33          es2 = es2 + (d - y_mmt)^2;
34      end
35      E1(epoch) = es1 / N;
36      E2(epoch) = es2 / N;
37
38 end
39
40 plot(E1, '--b', 'Linewidth',1);
41 hold on
42 plot(E2, '-.r', 'Linewidth',1);
43 xlabel('the number of Epoch')
44 ylabel('Average of Training error')
45 legend('Simple', 'Mmt')
```
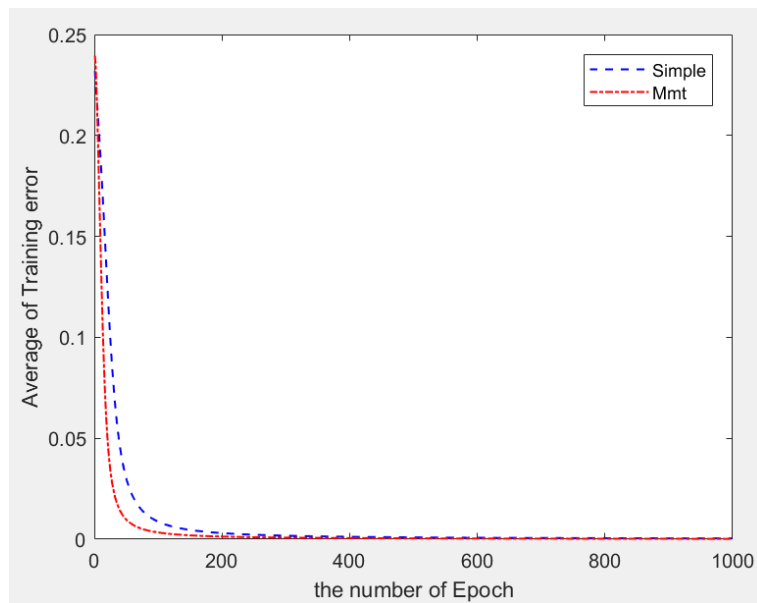
**Output 14**:



Figure 1: Mmt has better performance!.