

Momentum - GD with Momentum

Saturday, October 9, 2021 12:11 PM

What - Momentum is a force, comes from Physics. We use to accelerate our GD

$$m = \beta m + \eta \nabla_0 J(0) \Rightarrow m = \beta m + \eta \frac{\partial L}{\partial W} \Big|_{W_0}$$

$$0 = 0 \cdot m$$

$$W_N = W_0 - m \Rightarrow \textcircled{1}$$

$$W_N = W_0 - \eta \frac{\partial L}{\partial W} \rightarrow \text{GD}$$

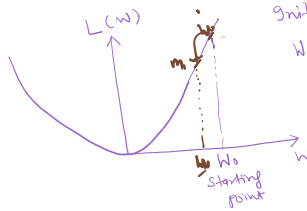
$$W_N = W_0 - \underbrace{\beta m}_{\text{Momentum term}} - \eta \frac{\partial L}{\partial W}$$

β - Coeff. of momentum

m - Affixed variable

Let's take $\beta = 0$

$$W_N = W_0 - \eta \frac{\partial L}{\partial W} \rightarrow \text{Gradient Descent normal}$$



Initialization \Rightarrow

$$W = W_0, m_0 = 0, \beta = 0.9 \text{ default}$$

Step 1 \Rightarrow

$$m_1 = \beta m_0 + \eta \frac{\partial L}{\partial W} \Big|_{W=W_0}$$

$$m_1 = \eta \frac{\partial L}{\partial W} \Big|_{W=W_0} \text{ if } m_0 = 0 \rightarrow \textcircled{1}$$

$$W_1 = W_0 - m_1$$

$$\Rightarrow W_1 = W_0 - \eta \frac{\partial L}{\partial W} \Big|_{W=W_0} \rightarrow \textcircled{1} \text{ GD only}$$

$$\text{Step 2} \Rightarrow m_1 = \eta \frac{\partial L}{\partial W} \Big|_{W=W_0}, W_1 = W_0 - \eta \frac{\partial L}{\partial W} \Big|_{W=W_0}, \beta = 0.9$$

$$m_2 = \beta m_1 + \eta \frac{\partial L}{\partial W} \Big|_{W=W_1}$$

$$m_2 = \beta \cdot \eta \frac{\partial L}{\partial W} \Big|_{W=W_0} + \eta \frac{\partial L}{\partial W} \Big|_{W=W_1}$$

$$m_2 = \eta \left[\underbrace{\beta \cdot \frac{\partial L}{\partial W} \Big|_{W=W_0}}_{\substack{0.9 \\ \text{(fixed value)}}} + \underbrace{\frac{\partial L}{\partial W} \Big|_{W=W_1}}_{\substack{\text{Current} \\ \text{gradient (100\%)}}} \right] \rightarrow \textcircled{m_2}$$

90% past gradient

$$W_2 = W_1 - m_2 = W_1 - \eta \left[\beta \frac{\partial L}{\partial W} \Big|_{W=W_0} + \frac{\partial L}{\partial W} \Big|_{W=W_1} \right]$$

Let's take a special case $\beta = 0$

$$W_2 = W_1 - \eta \frac{\partial L}{\partial W} \Big|_{W=W_1} \text{ — Simple GD weight updates } \rightarrow \textcircled{A}$$

$$\beta = 0.5$$

$$W_2 = W_1 - \eta \left[0.5 \frac{\partial L}{\partial W} \Big|_{W=W_0} + \frac{\partial L}{\partial W} \Big|_{W=W_1} \right] \rightarrow \textcircled{B}$$

Q - Which scenario will have more weight updates?? A OR B

Step 3 \Rightarrow

$$m_3 = \beta m_2 + \eta \frac{\partial L}{\partial W} \Big|_{W=W_2} \quad W = W_2, \beta = 0.9, m_2 = \textcircled{m_2}$$

$$m_3 = \eta \left[\beta \cdot \frac{\partial L}{\partial W} \Big|_{W=W_0} + \frac{\partial L}{\partial W} \Big|_{W=W_1} \right] \rightarrow \textcircled{m_3}$$

$$m_3 = \beta \left[\eta \left[\beta \frac{\partial L}{\partial W} \Big|_{W=W_0} + \frac{\partial L}{\partial W} \Big|_{W=W_1} \right] \right] + \eta \frac{\partial L}{\partial W} \Big|_{W=W_2}$$

$$m_3 = \beta \cdot \eta \left[\beta \frac{\partial L}{\partial W} \Big|_{W=W_0} + \frac{\partial L}{\partial W} \Big|_{W=W_1} \right] + \eta \frac{\partial L}{\partial W} \Big|_{W=W_2}$$

$$m_3 = \eta \left[\underbrace{\beta^2 \frac{\partial L}{\partial W} \Big|_{W=W_0}}_{81\% \text{ of } W_0} + \underbrace{\beta \frac{\partial L}{\partial W} \Big|_{W=W_1}}_{90\% \text{ of } W_1} + \underbrace{\frac{\partial L}{\partial W} \Big|_{W=W_2}}_{100\% \text{ of } W_2} \right] \rightarrow \textcircled{IV}$$

$$m_3 = \eta \left[\underbrace{\beta^2 \frac{\partial L}{\partial \omega}}_{81\% \text{ of } \omega_0} \Big|_{\omega=\omega_0} + \underbrace{\beta \cdot \frac{\partial L}{\partial \omega}}_{90\% \text{ of } \omega_1} \Big|_{\omega=\omega_1} + \underbrace{\frac{\partial L}{\partial \omega}}_{100\% \text{ of } \omega_2} \Big|_{\omega=\omega_2} \right] \rightarrow \textcircled{IV}$$

→ even if this becomes 0

Challenges above?

- 1) β - Additional parameter to tune - 0-9
- 2) Oscillations closer to local or global minima
 \hookrightarrow take $\frac{\partial L}{\partial \omega} = 0$