

# GRADIENT DESCENT

Formula for G.D =  $\theta = \theta - \alpha \nabla J(\theta)$

Formula to calculate  $\nabla J(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - \theta x_i)^2 x_i$

$\alpha$  is the learning rate.

Dataset:

$$(x, y) = (1, 2), (2, 3), (3, 4)$$

Initialize the weight parameter  $\theta$  to 0

$$\theta = 0.$$

$$\nabla J(\theta) = 1/3 [(2-\theta(1))^2 + (3-\theta(2))^2 + (4-\theta(3))^2] \times 1$$

$$\nabla J(\theta) = 2 \quad \therefore \frac{1}{N} \sum_{i=1}^N (y_i - \theta x_i)^2 x_i$$

$$\text{Take } \alpha = 0.1$$

$$\theta = \theta - 0.1(2) = \boxed{-\theta - 2}$$

$$\therefore \theta = \theta - \alpha \nabla J(\theta)$$

# Iteration 1

$$\theta = -0.2$$

$$\nabla J(\theta) = 1/3 [(2-\overset{(-0.2)}{\theta}(1))^2 + (3-\overset{(-0.2)}{\theta}(2))^2 + (4-\overset{(-0.2)}{\theta}(3))^2] \times 1$$

$$\nabla J(\theta) = 0.8$$

$$\theta = -0.2 - 0.1(0.8)$$

$$\boxed{\theta = -0.28}$$

## # Iteration 2

$$\theta = -0.28$$

$$\nabla J(\theta) = 1/3 [(2 - (-0.28)(1))^2 + (3 - (-0.28)(2))^2 + (4 - (-0.28)(3))^2] \times 1$$

$$\nabla J(\theta) = 0.48.$$

$$\theta = -0.28 - 0.1(0.48)$$

$$\boxed{\theta = -0.328}$$

## # Iteration 3

$$\theta = -0.328$$

$$\nabla J(\theta) = 1/3 [(2 - (-0.328)(1))^2 + (3 - (-0.328)(2))^2 + 4(-(-0.328)(3))^2] \times 1 = 0.248.$$

$$\boxed{\theta = -0.36528}$$

## # Iteration 4

$$\theta = -0.36528$$

$$\nabla J(\theta) = 0.1248$$

$$\theta = -0.36528 - (0.1248)(0.1) *$$

$$\boxed{\theta = -0.36528}$$

We can see that the weight parameter converges to the value  $\theta = -0.36528$  after 4 iterations.

# Stochastic Gradient Descent

Formula for SGD:-  $(y_i - \theta x_i) x_i = J(\theta)$

Weight Updation :-  $\theta_{\text{old}} - \alpha J(\theta) = \theta_{\text{new}}$

Dataset:

$$(x, y) = (1, 2), (2, 3), (3, 4)$$

Iteration 1:-  $x = 1, y = 2$

$$J(\theta) = (2 - \theta(1)) \quad \therefore J(\theta) = (y_i - \theta x_i) \quad \alpha = 0.1$$

$$= 2$$

$$\theta_{\text{new}} = 0 - 0.1(2)$$

$$\therefore \theta_{\text{new}} = \theta_{\text{old}} - \alpha J(\theta)$$

$$= -0.2$$

Iteration 2:-  $x = 2, y = 3$

$$J(\theta) = (3 + 0.2(2)) \quad 2$$

$$= 6.4$$

$$\theta_{\text{new}} = -0.2 - (6.4)0.1$$

$$= -0.84$$

Iteration 3:-  $x = 3, y = 4$

$$J(\theta) = (4 + 0.84(3)) \quad 3$$

$$= 19.5$$

$$\theta_{\text{new}} = -0.84 - 0.1(19.5)$$

$$= -2.79$$

After 3 iterations, the weight parameter converges to the value -2.79

# ADAPTIVE GRADIENT DESCENT

Formula for Adagrad =

$$\theta_t = \theta_{t-1} - \frac{\alpha}{\sqrt{G_{t-1} + \epsilon}} \cdot g_t$$

Dataset :

$$(x, y) = (2, 4), (5, 7), (6, 8)$$

We initialize

$$\theta = 0$$

$$\alpha = 0.1$$

$$G_0 = 0$$

$\epsilon$  : Its value is so small that it has negligible impact on the result.

$$x = 2$$

$$y = 4$$

$$g_t = 2(\theta - y \cdot x) \quad \therefore g_t = 2(\theta + x_i - y_i)$$

$$= 2(0 - 4(2))$$

$$\boxed{|g_t = -16|}$$

$$\Rightarrow G_t = G_{t-1} + (g_t)^2$$

$$= 0 + (-16)^2$$

$$\boxed{|G_t = 256|}$$

$$\Rightarrow \theta_t = \theta_{t-1} - \frac{\alpha}{\sqrt{G_{t-1} + \epsilon}} \cdot g_t$$

$$= 0 - \frac{0.1}{\sqrt{256}} \cdot (-16) \Rightarrow \boxed{|\theta_t = 0.1|}$$

## Iteration 2

$$x = 5$$

$$y = 7$$

$$g_t = 2(0.1 - 7(5))$$

$$\boxed{g_t = -69.8}$$

$$G_t = 16 \cdot 256 + (-69.8)^2$$

$$= 256 + 4872.04$$

$$\boxed{G_t = 5128.04}$$

$$\theta = 0.1 - \frac{0.1}{\sqrt{5128.04}} \cdot (-69.8)$$

$$\theta = 0.1 - \frac{0.1}{71.61} (-69.8)$$

$$\theta = 0.1 - 0.0013 (-69.8)$$

$$\theta = 0.1 (0.0974)$$

$$\boxed{\theta = 0.00974, 0.1974}$$

## Iteration 3

$$x = 6$$

$$y = 8$$

$$g_t = 2(0.1974 - 8(6))$$

$$\boxed{g_t = -95.60}$$

$$G_t = 5128.04 + (-95.60)^2$$

$$G_t = 14267.4$$

$$\begin{aligned}
 \theta &= 0.1974 - \frac{0.1}{\sqrt{14267.4}} (-95.60) \\
 &= 0.1974 - \frac{0.1}{119.44} (-95.60) \\
 &= 0.1974 - (0.00083)(-95.60) \\
 &= 0.1974 - \cancel{0.00083} (-0.0800) \\
 \boxed{\theta = 0.2774}
 \end{aligned}$$

After third iteration the parameter vector  $\theta$  has converged to the value 0.2774.

## Adadelta:-

$$\text{Formula :- } \Delta \theta_t = -\alpha \frac{\text{RMS}[g_t]}{\text{RMS}[\Delta \theta_{t-1}]}$$

$\therefore g_t^i = 2(\theta_t x_i - y_i)$   
 $\therefore x_i = \text{feature}$   
 $\theta_t = \text{Predicted output}$

$$\text{Weight Updation :- } \theta_{t+1} = \theta_t + \Delta \theta_t$$

$$\text{Dataset :- } (x, y) = (1, 2), (2, 3), (3, 4)$$

$$\text{Iteration 1 :- } x = 1, y = 2$$

$$g_t^1 = 2(0 \times 1 - 2) = -4 \quad \therefore \theta = 0$$

$$x = 2, y = 3$$

$$g_t^2 = 2(0 \times 2 - 3) = -6$$

$$x = 3, y = 4$$

$$g_t^3 = 2(0 \times 3 - 4) = -8$$

Calculate RMS :-

$$\Rightarrow \text{RMS}[g_t] = \sqrt{\frac{1}{N} \sum_{i=1}^N [g_t^i]^2}$$

$$= \sqrt{\frac{1}{3} ((-4)^2 + (-6)^2 + (-8)^2)}$$

$$= \sqrt{\frac{1}{3} (16 + 36 + 64)}$$

$$= \sqrt{\frac{116}{3}}$$

$$= \sqrt{38.67}$$

$$\text{RMS}[g_t] \approx 6.22$$

Now, we will calculate  $\Delta\theta_t$

$$\Rightarrow \Delta\theta_t = -\alpha \cdot \frac{\text{RMS}[g_t]}{\text{RMS}[\Delta\theta_{t-1}]} \quad \therefore \alpha = 0.1$$

$\therefore \text{RMS}[\Delta\theta_{t-1}]$  represents the RMS of previous iteration  
but for 1st iteration, we will assume it's 1  
because there's no previous update

$$\Delta\theta_t = -0.1 \cdot \frac{6.22}{1}$$

$$\therefore \text{RMS}[\Delta\theta_{t-1}] = 1$$

$$\therefore \Delta\theta_t = -\alpha \cdot \frac{\text{RMS}[g_t]}{\text{RMS}[\Delta\theta_{t-1}]}$$

$$\Delta\theta_t \approx -0.622$$

Now for weight updation:-

$$\theta_{t+1} = \theta_t + \Delta\theta_t$$

$$\theta_{0+1} = \theta_0 + (-0.622)$$

$$\theta_1 = 0 + (-0.622)$$

$$\theta_1 = -0.622$$

Iteration 2:-

Calculate  $g_t^i$

$$g_t^i = 2(-0.622 \times 1 - 2) = -4.244$$

$$g_t^{i+1} = 2(-0.622 \times 2 - 3) = -8.488$$

$$g_t^3 = 2(-0.622 \times 3 - 4) = -11.732$$

$$\begin{aligned}
 \text{RMS}[g_2] &= \sqrt{\frac{1}{3}((-4.244)^2 + (-8.488)^2 + (-11.732)^2)} \\
 &= \sqrt{\frac{1}{3}(18.01 + 72.04 + 137.63)} \\
 &= \sqrt{\frac{227.68}{3}} \\
 &= \sqrt{75.89}
 \end{aligned}$$

$$\text{RMS}[g_2] \approx 8.71$$

$$\Delta\theta_2 = -0.1 \frac{8.71}{1}$$

$$\Delta\theta_2 = -0.871$$

update weight

$$\theta_{t+1} = \theta_t + \Delta\theta_t$$

$$\theta_2 = -0.622 - 0.871$$

$$\theta_2 = -1.493$$

Iteration 3:-

calculate  $g_t^i$

$$g_t^1 = 2(-1.493 \times 1 - 2) = -6.986$$

$$g_t^2 = 2(-1.493 \times 2 - 3) = -11.972$$

$$g_t^3 = 2(-1.493 \times 3 - 4) = -16.958$$

$$\text{RMS}[g_3] = \sqrt{\frac{1}{3} ((-6.986)^2 + (-11.972)^2 + (-16.958)^2)}$$

$$= \sqrt{\frac{1}{3} (48.80 + 143.32 + 287.57)}$$

$$= \sqrt{\frac{479.69}{3}}$$

$$= \sqrt{159.89}$$

$$\text{RMS}[g_3] \approx 12.64$$

$$\Delta\theta_3 = -0.1 \frac{12.64}{1}$$

$$\Delta\theta_3 = -1.264$$

Update weight

$$\theta_3 = -1.493 + (-1.264)$$

$$\theta_3 = -2.757$$

# ROOT MEAN SQUARE PROPAGATION

## Step 1: Initialization

$$\theta = 1$$

$$\alpha = 0.01$$

$$\gamma = 0.9$$

$$E[g^2] = 0$$

$$\epsilon = 1e-8$$

Dataset :

$$(x, y) = (1, 2), (2, 3), (3, 4)$$

## Step 2: Compute MSE and Gradient for $\theta$

For each data point  $(x, y)$ .

For  $\theta$ :

- Compute the predicted value using the current  $\theta$ :  $y_{\text{pred}} = \theta \cdot x$
- calculate the mean squared error (MSE) for that data point :  $MSE = \frac{1}{2} \cdot (y_{\text{pred}} - y)^2$

Compute the gradient with respect to  $\theta$ :

$$\frac{\partial MSE}{\partial \theta} = (y_{\text{pred}} - y) \cdot x$$

## Step 3: Update $E[g^2]$ for $\theta$

Update the moving average of squared gradient for  $\theta$

$$E[g^2] = \gamma \cdot E[g^2] + (1 - \gamma) \cdot (\frac{\partial MSE}{\partial \theta})^2$$

## Step 4: Update $\theta$

Update  $\theta$  using the RMSprop formula,

$$\theta = \theta - (\alpha / \sqrt{E[g^2]} + \epsilon) \cdot \frac{\partial MSE}{\partial \theta}$$

Iteration 1: For the datapoint (1, 2)

$$y_{pred} = 1 \cdot 1 = 1$$

$$MSE = \frac{1}{2} (1-2)^2 = 0.5$$

$$\delta MSE / \delta \theta = (1-2) \cdot 1 = -1$$

$$\text{Update } E[g^2] = 0.9(0) + 0.1(-1)^2 = 0.1.$$

update  $\theta$ :

$$\theta = 1 - (0.01 / \sqrt{0.1 + 1e-6})(-1) \approx 1.005.$$

Iteration 2:

For the datapoint (2, 3)

$$y_{pred} = 1.005(2) \approx 2.010.$$

$$MSE = \frac{1}{2} (2.010 - 3)^2 \approx 0.505$$

$$\delta MSE / \delta \theta \approx (2.010 - 3)2 \approx -1.979$$

Update  $E[g^2]$ :

$$E[g^2] \approx 0.9(0.1) + 0.1(-1.979)^2 \approx 0.2979$$

update  $\theta$ :

$$\theta \approx 1.005 - (0.01 / \sqrt{0.2979 + 1e-8})(-1.979) \approx 1.008$$

Iteration 3:

For the datapoint (3, 4)

$$y_{pred} \approx 1.005 \quad (3) \approx 3.025$$

$$MSE \approx \frac{1}{2} (3.025 - 4)^2 \approx 0.460$$

$$\delta MSE / \delta \theta \approx (3.025 - 4)3 \approx -2.924$$

update  $E[g^2]$ :

$$E[g^2] \approx 0.9(0.2979) + 0.1(-2.924)^2 \approx 0.517.$$

update  $\theta$ :

$$\theta \approx 1.008 - (0.01 / \sqrt{0.517 + 1e-8})(-2.924) \approx 1.017$$

## Adam:-

Formulas:-

$$1^{\text{st}} \text{ Moment} = m_t = \beta_1 \times m_{t-1} + (1 - \beta_1) \times g_t \Rightarrow \text{Mean}$$

$$2^{\text{nd}} \text{ Moment} = v_t = \beta_2 \times v_{t-1} + (1 - \beta_2) (g_t)^2 \Rightarrow \text{Variance}$$

$$\text{Bias} = \hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$\theta_{t+1} = \theta_t - \left( \alpha / (\sqrt{\hat{v}_t} + \epsilon) \right) \times \hat{m}_t$$

$$\therefore w, b \Rightarrow \theta_{t+1} = \theta_t - \frac{\alpha \hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \Rightarrow \text{Parameter Update}$$

Initialization:

$$\alpha = 0.01$$

$$\beta_1 = 0.9$$

$$\beta_2 = 0.999$$

$$\epsilon = 1e-8$$

$$w = 0 \text{ (slope)}$$

$$b = 0 \text{ (intercept)}$$

$$m_w = 0 \text{ (1^{st} moment for slope)}$$

$$m_b = 0 \text{ (1^{st} moment for intercept)}$$

$$v_w = 0 \text{ (2^{nd} moment for slope)}$$

$$v_b = 0 \text{ (2^{nd} moment for intercept)}$$

$$t = 0$$

Iteration 1: -  $(x, y) = (1, 2)$   $\therefore$  mabA

$$y_{\text{pred}} = 0 \cdot 1 + 0 = 0$$

$$\frac{\partial \text{MSE}}{\partial w} = (y_{\text{pred}} - 2) \cdot 1 = -2$$

$$\frac{\partial \text{MSE}}{\partial b} = (y_{\text{pred}} - 2) = -2$$

Update 1<sup>st</sup> Moments for  $w$  &  $b$

$$m_w = 0.9 \times 0 + 0.1 \times (-2) = -0.2$$

$$m_b = 0.9 \times 0 + 0.1 \times (-2) = -0.2$$

Update 2<sup>nd</sup> Moments for  $w$  &  $b$

$$v_w = 0.999 \times 0 + 0.001 \times (-2)^2 = 0.002$$

$$v_b = 0.999 \times 0 + 0.001 \times (-2)^2 = 0.002$$

Calculate bias in 1<sup>st</sup> Moment

$$\hat{m}_w = \frac{m_w}{1 - \beta_1^+} = -0.2$$

$$\hat{m}_b = \frac{m_b}{1 - \beta_2^+} = -0.2$$

Calculate bias in 2<sup>nd</sup> moments

$$\hat{v}_w = \frac{v_w}{1 - \beta_2^+} = 0.002$$

$$\hat{v}_b = \frac{v_b}{1 - \beta_2^+} = 0.002$$

Update parameters  $w$  &  $b$ :

$$w = w - \frac{\alpha}{\sqrt{\hat{v}_w + \epsilon}} \cdot \hat{m}_w$$
$$= 0 - \frac{0.01}{\sqrt{0.002 + 1e-8}} (-0.2)$$

$$w \approx 0.141$$

$$b = b - \frac{\alpha}{\sqrt{\hat{v}_b + \epsilon}} \cdot \hat{m}_b$$
$$= 0 - \frac{0.01}{\sqrt{0.002 + 1e-8}} (-0.2)$$

$$b \approx 0.141$$

Iteration 2:-  $(x, y) = (2, 3)$

Initialization:-

$$w = 0.141 , b = 0.141$$

$$m_w = -0.2 , m_b = -0.2$$

$$v_w = 0.002 , v_b = 0.002$$

$$t = 1 \text{ (incremented from iteration 1)}$$

$$y_{pred} = 0.141 \times 2 + 0.141$$

$$y_{pred} \approx 0.424$$

$$\frac{\partial MSB}{\partial w} \approx (0.424 - 3) 2$$

$$\frac{\partial \text{MSE}}{\partial w} \approx -4.351$$

$$\frac{\partial \text{MSE}}{\partial b} \approx (0.424 - 3)$$

$$\frac{\partial \text{MSE}}{\partial b} \approx -2.575$$

Update 1<sup>st</sup> Moment for w & b.

$$m_w = 0.9 \times (-0.2) + 0.1 \times (-4.351)$$

$$m_w \approx -0.385$$

$$m_b = 0.9 \times (-0.2) + 0.1 \times (-2.575)$$

$$m_b \approx -0.298$$

Update 2<sup>nd</sup> Moment for w & b.

$$v_w = 0.999 \times 0.002 + 0.001 \times (-4.351)^2$$

$$v_w \approx 0.004$$

$$v_b = 0.999 \times 0.002 + 0.001 \times (-2.575)^2$$

$$v_b \approx -0.298$$

Calculate bias in 1<sup>st</sup> Moment.

$$\hat{m}_w \approx \frac{-0.385}{1-0.9^2} \approx -0.053$$

$$\hat{m}_b \approx \frac{-0.298}{1-0.9^2} \approx -0.041$$

Calculate bias for 2nd Moment

$$v_w \approx \frac{0.004}{1 - 0.999^2} \approx 0.111$$

$$v_b \approx \frac{0.0023}{1 - 0.999^2} \approx 0.058$$

update  $w \quad \epsilon_p \quad b$

$$w \approx 0.141 - \frac{0.01}{\sqrt{0.111 + 1e-8}} \cdot (-0.053)$$

$$w \approx 0.283$$

$$b \approx 0.141 - \frac{0.01}{\sqrt{0.058 + 1e-8}} \cdot (-0.041)$$

$$b \approx 0.184$$

Iteration 3:  $(x, y) = (3, 4)$

$$y_{pred} \approx 0.283 \times 3 + 0.184$$

$$y_{pred} \approx 0.917$$

$$\frac{\partial \text{MSE}}{\partial w} \approx (0.917 - 4) \times 3$$

$$\frac{\partial \text{MSE}}{\partial w} \approx -9.748$$

$$\frac{\partial \text{MSE}}{\partial b} \approx (0.917 - 4) \approx -3.082$$

Update 1<sup>st</sup> Moment for  $\omega$  &  $b$

$$m_{\omega} \approx 0.9 \times (-0.053) + 0.1(-9.748)$$
$$\approx -0.964$$

$$m_b \approx 0.9 \times (-0.041) + 0.1(-3.082)$$
$$\approx -0.434$$

Update 2<sup>nd</sup> Moment for  $\omega$  &  $b$

$$v_{\omega} \approx 0.999(0.111) + 0.001(-9.748)^2$$
$$\approx 0.166$$

$$v_b \approx 0.999(0.058) + 0.001(-3.082)^2$$
$$\approx 0.084$$

Calculate bias for 1<sup>st</sup> Moment

$$\hat{m}_{\omega} \approx \frac{-0.964}{1 - 0.9^3} \approx -0.268$$

$$\hat{m}_b \approx \frac{-0.434}{1 - 0.9^3} \approx -0.119$$

Calculate bias for 2<sup>nd</sup> Moment

$$\hat{v}_{\omega} \approx \frac{0.166}{1 - 0.999^3} \approx 0.333$$

$$\hat{v}_b \approx \frac{0.084}{1 - 0.999^3} \approx 0.167$$

Update  $\mu$  &  $b$

$$w \approx 0.283 - \frac{0.01}{\sqrt{0.333 + 1e-8}} \cdot (-0.268)$$

$$w \approx 0.043$$

$$b \approx 0.184 - \frac{0.01}{\sqrt{0.167 + 1e-8}} \cdot (-0.119)$$

$$b \approx 0.259$$