

Gradient Descent vs Stochastic Gradient Descent vs Mini Batch Gradient Descent

Let us have 5 training examples :

$$\langle x_1^{(1)}, y^{(1)} \rangle, \langle x_1^{(2)}, y^{(2)} \rangle, \langle x_1^{(3)}, y^{(3)} \rangle, \langle x_1^{(4)}, y^{(4)} \rangle, \langle x_1^{(5)}, y^{(5)} \rangle$$

predicted output be denoted by $o^{(i)}$

$$\text{Let, } o^{(i)} = w_0 + w_1 x_1^{(i)}$$

$$o^{(1)} = w_0 + w_1 x_1^{(1)}$$

$$o^{(2)} = w_0 + w_1 x_1^{(2)}$$

$$o^{(3)} = w_0 + w_1 x_1^{(3)}$$

$$o^{(4)} = w_0 + w_1 x_1^{(4)}$$

$$o^{(5)} = w_0 + w_1 x_1^{(5)}$$

Here, i goes from 1 to 5

Weight Update,

$$w_j = w_j + \Delta w_j$$

Here, j goes from 0 to 1

where,

$$\Delta w_j = -\eta \left(\partial E / \partial w_j \right)$$

In case of Gradient Descent,

$$\begin{aligned} E &= (1/2) * [[y^{(1)} - o^{(1)}]^2 + [y^{(2)} - o^{(2)}]^2 + [y^{(3)} - o^{(3)}]^2 + [y^{(4)} - o^{(4)}]^2 + [y^{(5)} - o^{(5)}]^2] \\ &= (1/2) * [\\ &\quad [y^{(1)} - (w_0 + w_1 x_1^{(1)})]^2 + \\ &\quad [y^{(2)} - (w_0 + w_1 x_1^{(2)})]^2 + \\ &\quad [y^{(3)} - (w_0 + w_1 x_1^{(3)})]^2 + \\ &\quad [y^{(4)} - (w_0 + w_1 x_1^{(4)})]^2 + \\ &\quad [y^{(5)} - (w_0 + w_1 x_1^{(5)})]^2 \\ &\quad] \\ &= (1/2) * \sum [y^{(i)} - o^{(i)}]^2 \end{aligned}$$

Again, i goes from 1 to 5. Same value of w_0 and w_1 is being used for all the 5 examples.

We update weight only after we have seen all the examples once.

Where,

$$\begin{aligned} \partial E / \partial w_0 &= (1/2) * [\\ &\quad 2 * [y^{(1)} - (w_0 + w_1 x_1^{(1)})] * (-1) + \\ &\quad 2 * [y^{(2)} - (w_0 + w_1 x_1^{(2)})] * (-1) + \\ &\quad 2 * [y^{(3)} - (w_0 + w_1 x_1^{(3)})] * (-1) + \\ &\quad 2 * [y^{(4)} - (w_0 + w_1 x_1^{(4)})] * (-1) + \\ &\quad 2 * [y^{(5)} - (w_0 + w_1 x_1^{(5)})] * (-1) \\ &\quad] \\ &= (-1/2) * 2 * \sum [y^{(i)} - (w_0 + w_1 x_1^{(i)})] \\ &= (-1/2) * 2 * \sum [y^{(i)} - o^{(i)}] \end{aligned}$$

Again, i goes from 1 to 5.

$$\begin{aligned} \partial E / \partial w_1 &= (1/2) * [\\ &\quad 2 * [y^{(1)} - (w_0 + w_1 x_1^{(1)})] * (-x_1^{(1)}) + \end{aligned}$$

$$\begin{aligned}
& 2*[y^{(2)} - (w_0 + w_1 x_1^{(2)})]*(-x_1^{(2)}) + \\
& 2*[y^{(3)} - (w_0 + w_1 x_1^{(3)})]*(-x_1^{(3)}) + \\
& 2*[y^{(4)} - (w_0 + w_1 x_1^{(4)})]*(-x_1^{(4)}) + \\
& 2*[y^{(5)} - (w_0 + w_1 x_1^{(5)})]*(-x_1^{(5)}) \\
&] \\
& = (-1/2)*2*\sum [y^{(i)} - (w_0 + w_1 x_1^{(i)})]*(x_1^{(i)}) \\
& = (-1/2)*2*\sum [y^{(i)} - o^{(i)}]*(x_1^{(i)})
\end{aligned}$$

which is equivalent to

$$\partial E / \partial w_j = (-1/2)*2*\sum [y^{(i)} - (w_0 + w_j x_j^{(i)})]*(x_j^{(i)})$$

Again, i goes from 1 to 5.

Therefore,

In Gradient Descent we update the weights only after seeing all the examples.

$$w_j = w_j + \Delta w_j$$

where

$$\begin{aligned}
\Delta w_j \text{ when } j = 0 & \text{ is } \eta \sum [y^{(i)} - o^{(i)}] \\
j = 1 & \text{ is } \eta \sum [y^{(i)} - (w_0 + w_1 x_1^{(i)})]*(x_1^{(i)})
\end{aligned}$$

Again, i goes from 1 to 5 (i.e., over all the training examples).

In Stochastic Gradient Descent, we update the weight after each training example:

$$w_j = w_j + \Delta w_j$$

where

$$\begin{aligned}
\Delta w_j \text{ when } j = 0 & \text{ is } \eta [y^{(i)} - o^{(i)}] \\
j = 1 & \text{ is } \eta [y^{(i)} - (w_0 + w_1 x_1^{(i)})]*(x_1^{(i)})
\end{aligned}$$

Considering Batch Size of 2 examples,

According to Mini-batch Gradient Descent,

We would take

$$\text{Batch_one} = \langle x_1^{(1)}, y^{(1)} \rangle, \langle x_1^{(2)}, y^{(2)} \rangle$$

and apply gradient descent to it

Then take

$$\text{Batch_two} = \langle x_1^{(3)}, y^{(3)} \rangle, \langle x_1^{(4)}, y^{(4)} \rangle$$

and apply gradient descent to it

Then take

$$\text{Batch_three} = \langle x_1^{(5)}, y^{(5)} \rangle$$

and apply gradient descent to it.