

Backpropagation

- A technique used to compute gradients efficiently using chain rule of partial derivatives
- Calculated by moving backwards in the network.

Gradient Descent "Apply the changes to descent"

- Updates weights & bias using gradients computed during backpropagation

Learning?

- Finding the optimal value for weights & biases such that cost is minimum.
- Finds the derivative of cost fn (gradient) & move in that dir.

Weights \rightarrow cost $f_2 \rightarrow$ cost : Since input to cost f_2 is weights & bias,

Bias Gradient descent $\rightarrow \nabla C$ will also have size of weights & bias.

$\nabla C = \begin{bmatrix} 0.1 \\ 3.2 \\ 4.5 \end{bmatrix}$: magnitude of gradient indicates which weight is more important in network

performed for each weights & bias per node, then

3.2 is 32 times more than 0.1

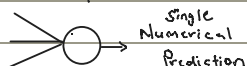
\rightarrow Effect more on cost f_2 .

Cost Function

- Average loss for all samples
- Gradient descent performed on weights & biases of each node

Understanding Learning intuitively.

Numerical prediction



Suppose predict = 12000 & Activation from last layer.
Actual = 10000

Goal = -2000

$$\text{Activation} = \sigma(W \cdot A + B)$$

$$-2000 = \sigma(W_1 a_1 + \dots W_n a_n + B)$$

Things Adjustable :

1) Activation f_2 (Not very effective)

most imp update \rightarrow 2) Alter weights w_i (where gradient magnitudes are higher)

3) Adjust a_i - +ve values reduces \uparrow (No direct influence, dependent on previous layer weights & bias)
-ve values increases

4) decrease B_i

Gradient vector

* Repeat same steps for all nodes in previous layer till input layer.

Repeat same step for all examples.

Take the average change of all examples.

	Sample 1	Sample 2	Sample 3 ...	
w_0	-0.8	+0.2	-0.4	$\left. \begin{array}{l} \\ \\ \\ \end{array} \right\} -\nabla C = \begin{bmatrix} -\frac{0.1}{3} \\ \frac{0.11}{3} \\ 0 \\ \vdots \end{bmatrix}$ changes required for all weight & bias of each node for current iteration
w_1	+0.7	+0.3	+0.1	
w_2	-0.1	-0.4	+0.5	
\vdots				

Stochastic Gradient Descent

- 1) Randomly shuffle training data
- 2) Split data into mini batches
- 3) compute gradient descent for each batch

\rightarrow Not the Actual gradient, thus not the most efficient but a good approximation

\rightarrow faster computation

Stochastic gradient descent
 \leftarrow fast movement
 \leftarrow min

Gradient descent
 \leftarrow slow movement
 \leftarrow min