

1 Gradient Descent

If we keep decreasing the ϵ in our Finite Difference approach we effectively get the Derivative of the Cost Function.

$$C'(w) = \lim_{\epsilon \rightarrow 0} \frac{C(w + \epsilon) - C(w)}{\epsilon} \quad (1)$$

Let's compute the derivatives of all our models. Throughout the entire paper n means the amount of samples in the training set.

1.1 Linear Model



$$y = x \cdot w \quad (2)$$

1.1.1 Cost

$$C(w) = \frac{1}{n} \sum_{i=1}^n (x_i w - y_i)^2 \quad (3)$$

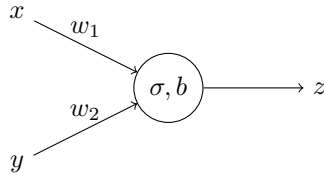
$$C'(w) = \left(\frac{1}{n} \sum_{i=1}^n (x_i w - y_i)^2 \right)' = \quad (4)$$

$$= \frac{1}{n} \left(\sum_{i=1}^n (x_i w - y_i)^2 \right)' \quad (5)$$

$$= \frac{1}{n} \sum_{i=1}^n ((x_i w - y_i)^2)' \quad (6)$$

$$= \frac{1}{n} \sum_{i=1}^n 2(x_i w - y_i)x_i \quad (7)$$

1.2 One Neuron Model with 2 inputs



$$z = \sigma(xw_1 + yw_2 + b) \quad (8)$$

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (9)$$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x)) \quad (10)$$

1.2.1 Cost

$$a_i = \sigma(x_iw_1 + y_iw_2 + b) \quad (11)$$

$$\partial_{w_1} a_i = \partial_{w_1} (\sigma(x_iw_1 + y_iw_2 + b)) = \quad (12)$$

$$= a_i(1 - a_i)\partial_{w_1} (x_iw_1 + y_iw_2 + b) = \quad (13)$$

$$= a_i(1 - a_i)x_i \quad (14)$$

$$\partial_{w_2} a_i = a_i(1 - a_i)y_i \quad (15)$$

$$\partial_b a_i = a_i(1 - a_i) \quad (16)$$

$$C = \frac{1}{n} \sum_{i=1}^n (a_i - z_i)^2 \quad (17)$$

$$\partial_{w_1} C = \frac{1}{n} \sum_{i=1}^n \partial_{w_1} ((a_i - z_i)^2) = \quad (18)$$

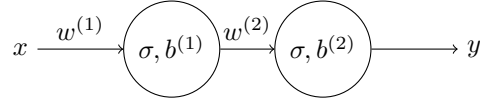
$$= \frac{1}{n} \sum_{i=1}^n 2(a_i - z_i)\partial_{w_1} a_i = \quad (19)$$

$$= \frac{1}{n} \sum_{i=1}^n 2(a_i - z_i)a_i(1 - a_i)x_i \quad (20)$$

$$\partial_{w_2} C = \frac{1}{n} \sum_{i=1}^n 2(a_i - z_i)a_i(1 - a_i)y_i \quad (21)$$

$$\partial_b C = \frac{1}{n} \sum_{i=1}^n 2(a_i - z_i)a_i(1 - a_i) \quad (22)$$

1.3 Two Neurons Model with 1 input



$$x = a^{(0)} \quad (23)$$

$$a^{(1)} = \sigma(xw^{(1)} + b^{(1)}) \quad (24)$$

$$y = a^{(2)} = \sigma(a^{(1)}w^{(2)} + b^{(2)}) \quad (25)$$

The superscript in parenthesis denotes the current layer. For example $a_i^{(l)}$ denotes the activation from the l -th layer on i -th sample.

Then it follows that for a model with L layers, we have:

$$x = a^{(0)} \quad (26)$$

$$a^{(1)} = \sigma(xw^{(1)} + b^{(1)}) \quad (27)$$

$$a^{(2)} = \sigma(a^{(1)}w^{(2)} + b^{(2)}) \quad (28)$$

$$a^{(3)} = \sigma(a^{(2)}w^{(3)} + b^{(3)}) \quad (29)$$

$$\dots \quad (30)$$

$$a^{(l)} = \sigma(a^{(l-1)}w^{(l)} + b^{(l)}) \quad (31)$$

$$y = a^{(l+1)} = \sigma(a^{(l)}w^{(l+1)} + b^{(l+1)}) \quad (32)$$

1.3.1 Feed-Forward

$$a_i^{(0)} = x_i \quad (33)$$

$$(34)$$

$$a_i^{(1)} = \sigma(x_iw^{(1)} + b^{(1)}) \quad (35)$$

$$\partial_{w^{(1)}} a_i^{(1)} = a_i^{(1)}(1 - a_i^{(1)})x_i \quad (36)$$

$$\partial_{b^{(1)}} a_i^{(1)} = a_i^{(1)}(1 - a_i^{(1)}) \quad (37)$$

$$\partial_{a_i^{(0)}} a_i^{(1)} = \partial_{x_i} a_i^{(1)} = a_i^{(1)}(1 - a_i^{(1)})w^{(1)} \quad (38)$$

$$(39)$$

$$a_i^{(2)} = \sigma(a_i^{(1)}w^{(2)} + b^{(2)}) \quad (40)$$

$$\partial_{w^{(2)}} a_i^{(2)} = a_i^{(2)}(1 - a_i^{(2)})a_i^{(1)} \quad (41)$$

$$\partial_{b^{(2)}} a_i^{(2)} = a_i^{(2)}(1 - a_i^{(2)}) \quad (42)$$

$$(43)$$

1.3.2 Back-Propagation

$$C^{(2)} = \frac{1}{n} \sum_{i=1}^n (a_i^{(2)} - y_i)^2 \quad (44)$$

$$\partial_{w^{(2)}} C^{(2)} = \frac{1}{n} \sum_{i=1}^n \partial_{w^{(2)}} ((a_i^{(2)} - y_i)^2) \quad (45)$$

$$= \frac{1}{n} \sum_{i=1}^n 2(a_i^{(2)} - y_i) \partial_{w^{(2)}} a_i^{(2)} = \quad (46)$$

$$= \frac{1}{n} \sum_{i=1}^n 2(a_i^{(2)} - y_i) a_i^{(2)} (1 - a_i^{(2)}) a_i^{(1)} \quad (47)$$

$$\partial_{b^{(2)}} C^{(2)} = \frac{1}{n} \sum_{i=1}^n 2(a_i^{(2)} - y_i) a_i^{(2)} (1 - a_i^{(2)}) \quad (48)$$

$$\partial_{a_i^{(1)}} C^{(2)} = \frac{1}{n} \sum_{i=1}^n 2(a_i^{(2)} - y_i) a_i^{(2)} (1 - a_i^{(2)}) w^{(2)} \quad (49)$$

$$e_i = a_i^{(1)} - \partial_{a_i^{(1)}} C^{(2)} \quad (50)$$

$$C^{(1)} = \frac{1}{n} \sum_{i=1}^n (a_i^{(1)} - e_i)^2 \quad (51)$$

$$\partial_{w^{(1)}} C^{(1)} = \partial_{w^{(1)}} \left(\frac{1}{n} \sum_{i=1}^n (a_i^{(1)} - e_i)^2 \right) = \quad (52)$$

$$= \frac{1}{n} \sum_{i=1}^n \partial_{w^{(1)}} ((a_i^{(1)} - e_i)^2) = \quad (53)$$

$$= \frac{1}{n} \sum_{i=1}^n 2(a_i^{(1)} - e_i) \partial_{w^{(1)}} a_i^{(1)} = \quad (54)$$

$$= \frac{1}{n} \sum_{i=1}^n 2(\partial_{a_i^{(1)}} C^{(2)}) a_i^{(1)} (1 - a_i^{(1)}) x_i \quad (55)$$

$$\partial_{b^{(1)}} C^{(1)} = \frac{1}{n} \sum_{i=1}^n 2(\partial_{a_i^{(1)}} C^{(2)}) a_i^{(1)} (1 - a_i^{(1)}) \quad (56)$$

1.4 Arbitrary Neurons Model with m input

Let's assume that we have m layers.

1.4.1 Feed-Forward

Let's assume that $a_i^{(0)}$ is x_i .

$$a_i^{(l)} = \sigma(a_i^{(l-1)} w^{(l)} + b^{(l)}) \quad (57)$$

$$\partial_{w^{(l)}} a_i^{(l)} = a_i^{(l)} (1 - a_i^{(l)}) a_i^{(l-1)} \quad (58)$$

$$\partial_{b^{(l)}} a_i^{(l)} = a_i^{(l)} (1 - a_i^{(l)}) \quad (59)$$

$$\partial_{a_i^{(l-1)}} a_i^{(l)} = a_i^{(l)} (1 - a_i^{(l)}) w^{(l)} \quad (60)$$

1.4.2 Back-Propagation

Let's denote $a_i^{(m)} - y_i$ as $\partial_{a_i^{(m)}} C^{(m+1)}$.

$$C^{(l)} = \frac{1}{n} \sum_{i=1}^n (\partial_{a_i^{(l)}} C^{(l+1)})^2 \quad (61)$$

$$\partial_{w^{(l)}} C^{(l)} = \frac{1}{n} \sum_{i=1}^n 2(\partial_{a_i^{(l)}} C^{(l+1)}) a_i^{(l)} (1 - a_i^{(l)}) a_i^{(l-1)} = \quad (62)$$

$$\partial_{b^{(l)}} C^{(l)} = \frac{1}{n} \sum_{i=1}^n 2(\partial_{a_i^{(l)}} C^{(l+1)}) a_i^{(l)} (1 - a_i^{(l)}) \quad (63)$$

$$\partial_{a_i^{(l-1)}} C^{(l)} = \frac{1}{n} \sum_{i=1}^n 2(\partial_{a_i^{(l)}} C^{(l+1)}) a_i^{(l)} (1 - a_i^{(l)}) w^{(l)} \quad (64)$$