# Excersise 1:

i) $\sigma(x) = \dfrac{1}{1+e^{-x}}$

$\dfrac{\partial \sigma(x)}{\partial x} = \dfrac{0 - 1 * e^{-x}(-1)}{(1+e^{-x})^2}$

$= \dfrac{e^{-x}}{(1+e^{-x})^2} \Rightarrow \dfrac{(1+e^{-x})}{(1+e^{-x})^2} - \dfrac{1}{(1+e^{-x})^2}$

$= \sigma(x) - \sigma^2(x) \qquad \because \dfrac{1}{1+e^{-x}} = \sigma(x)$

ii) Taylor Series $\Rightarrow$

$P(x) = f(a) + \dfrac{1}{1!}\dfrac{\partial f}{\partial x}(a)(x-a)^1 + \dfrac{1}{2!}\dfrac{\partial^2 f}{\partial x^2}(a)(x-a)^2 + \ldots$

Here $a = 0$

$P(x) = f(0) + \dfrac{1}{1!}\dfrac{\partial f}{\partial x}(0)\cdot x + \dfrac{1}{2!}\dfrac{\partial^2 f}{\partial x^2}(0)\cdot x^2$

$= \dfrac{1}{2} + \left(\dfrac{1}{2} - \dfrac{1}{4}\right)\cdot x + \dfrac{1}{2!}\cdot 0 \cdot x^2 = \dfrac{1}{2} + \dfrac{1}{4}x.$

$\dfrac{\partial^2 f}{\partial x^2} = \dfrac{\partial}{\partial x}\left(\sigma(x)(1-\sigma(x))\right) = \sigma(x)\cdot(-1)\left(\sigma(x) - \sigma^2(x)\right)$

$\qquad\qquad\qquad\qquad\qquad + \sigma(x)(1-\sigma(x))(1-\sigma(x))$

$\qquad\qquad = \sigma(x)(1-\sigma(x))^2 - \sigma^2(x)(1-\sigma(x))$

$\qquad\qquad = \sigma(x)(1-\sigma(x))\left((1-\sigma(x)) - \sigma(x)\right)$

At $x=0$: $\Rightarrow \dfrac{1}{2}*\dfrac{1}{2}\left(\dfrac{1}{2} - \dfrac{1}{2}\right) = 0$

(c)

Lets approximate the $f$ function at $w_n$:

$$f(x) = f(w_n) + \frac{1}{1!} \frac{\partial f}{\partial w}(w_n)(x - w_n) + \ldots$$

$$f(x) = f(w_n) + g(n) \cdot (x - w_n)$$

Lets approximate it at $x = w_{n+1}$

$$f(w_{n+1}) = f(w_n) + g(n)(w_{n+1} - w_n)$$

$$\therefore w_{n+1} = w_n - \epsilon \, g(n)$$

$$f(w_{n+1}) = f(w_n) + g(n)(w_n - \epsilon g(n) - w_n)$$

$$\boxed{f(w_{n+1}) = f(w_n) - \epsilon (\nabla f(w_n))^2}$$

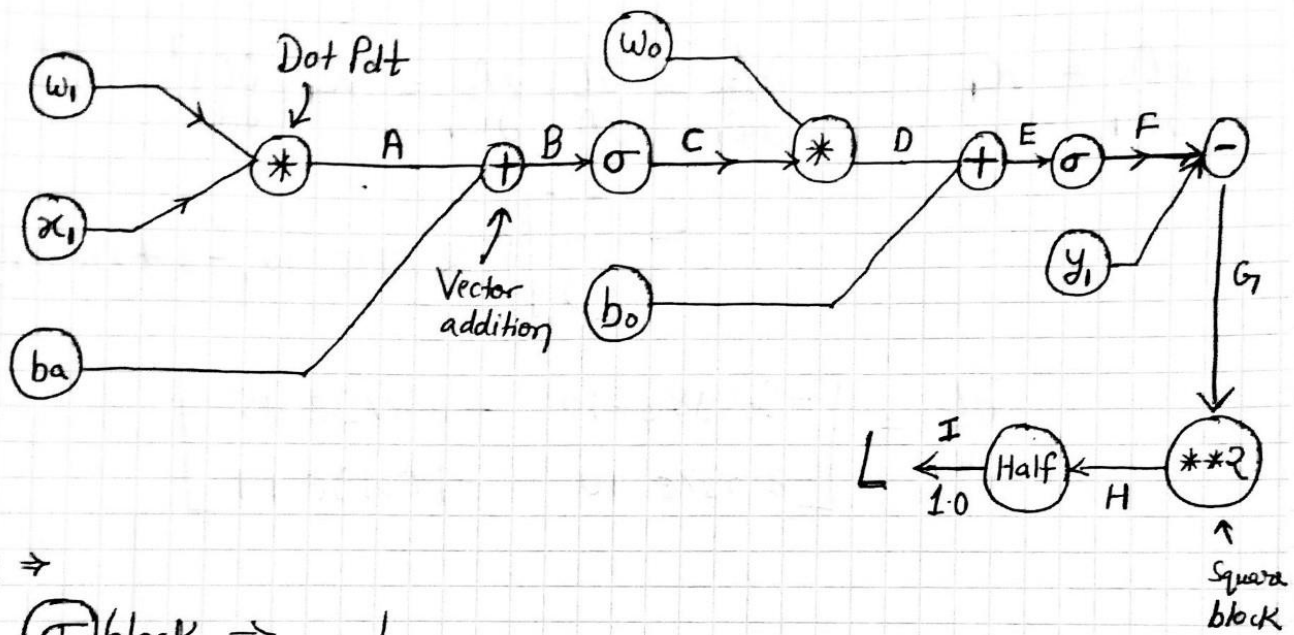This is similar to Gradient descent step, but instead of gradient we are using the Square of gradient.

Above Equation shows that the value of function at next step will decrease, and suddenly hit local/Global minima.

Excersise 2 ⇒

(a) $L = \frac{1}{2} \sum_{i=1}^{N} \left( \sigma\left(w_0^T\left(\sigma\left(w_i^T x_n + b_a\right)\right) + b_0\right) - y_n \right)^2$

for a single instance

$$= \frac{1}{2}\left( \sigma\left(w_0^T\left(\sigma\left(w_i^T x_i + b_a\right)\right) + b_0\right) - y_i \right)^2$$



Note ⇒

$\sigma$ block ⇒ $\dfrac{1}{1+e^{-x}}$

2.b : Lets assume some values $\rightarrow$ We could also assume Constants

$$W_1 = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}_{2\times2} \qquad W_0 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}_{2\times1} \qquad b_a = \begin{bmatrix} 2 \\ 4 \end{bmatrix} \qquad b_o = \begin{bmatrix} 3 \\ 0 \end{bmatrix}_{1\times1}$$

$$x_1 = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

$$\therefore y_1 = 1 \qquad \text{Assuming}$$

**Forward Pass** $\rightarrow$

$$A = W_1 x_1 = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 8 \\ 18 \end{bmatrix}$$

$$B = W_1 x_1 + b_a = \begin{bmatrix} 8 \\ 18 \end{bmatrix} + \begin{bmatrix} 2 \\ 4 \end{bmatrix} = \begin{bmatrix} 10 \\ 22 \end{bmatrix}$$

$$C = \sigma(B) = \begin{bmatrix} 0.999 \\ 1 \end{bmatrix}$$

- $D = W_0 \cdot C = [2.999]$

- $E = W_0 \cdot C + b_o = [5.999]$

- $F = \sigma(E) = [0.9975]$

- $G = [0.9975 - 1] = -0.00249$

- $H = G^2 = 6.24e-06$

$I = 3.124e-06 = L$

**Backward Pass** $\rightarrow$

- $\dfrac{\partial L}{\partial L} = \dfrac{\partial t}{\partial L} = 1 \qquad \therefore I = L$

- $I = \dfrac{H}{2} \Rightarrow \dfrac{\partial I}{\partial H} = \dfrac{1}{2} \qquad \therefore \dfrac{\partial t}{\partial I} * \dfrac{\partial I}{\partial H} \Rightarrow \boxed{\dfrac{\partial L}{\partial H} = \dfrac{1}{2}} \checkmark$

- $H = G^2 \qquad \therefore \dfrac{\partial H}{\partial G} = 2*G \qquad \therefore \dfrac{\partial L}{\partial H} * \dfrac{\partial H}{\partial G} \Rightarrow \dfrac{\partial L}{\partial G} = \dfrac{1}{2} * 2 G$

$$\boxed{\dfrac{\partial L}{\partial G} = G}$$

$\cdot\ G = F - y_1$

$$\frac{\partial G}{\partial F} = 1 \qquad \therefore\ \frac{\partial L}{\partial G}\frac{\partial G}{\partial F} = G * 1 = G \qquad \left[\ \therefore\ \frac{\partial L}{\partial F} = G\right]$$

$\cdot\ F = \sigma(E)$

$$\frac{\partial F}{\partial E} = \sigma(E)(1 - \sigma(E)) \Rightarrow F(1-F).$$

$$\therefore\ \frac{\partial L}{\partial F} * \frac{\partial F}{\partial E} = G\,F(1-F) = [-0.00249]\,[0.9975]$$

$$(1 - 0.9975)$$

$$\therefore\ \left[\frac{\partial L}{\partial E} = -6.1665e-06\right]$$

$E = D + b_0$ $\qquad \frac{\partial E}{\partial D} = 1 \qquad \therefore\ \left[\frac{\partial L}{\partial D} = \frac{\partial L}{\partial E} * \frac{\partial E}{\partial D} = -6.166e-06\right]$

$$\frac{\partial E}{\partial b_0} = 1 \qquad \therefore\ \left[\frac{\partial L}{\partial b_0} = \frac{\partial L}{\partial E} \cdot \frac{\partial E}{\partial b_0} = -6.166e-06\right]$$

$D = W_0 \cdot C$ $\qquad \frac{\partial D}{\partial W_0} = C \qquad \frac{\partial L}{\partial W_0} = \frac{\partial L}{\partial D} * \frac{\partial D}{\partial W_0} = -6.166e-06 * C$

$$\left[\frac{\partial L}{\partial W_0} = \begin{bmatrix} -6.165e-06 \\ -6.166e-06 \end{bmatrix}\right] \checkmark$$

$\frac{\partial D}{\partial C} = W_0$

$$\therefore\ \frac{\partial L}{\partial C} = -6.166e-06 * W_0 = \begin{bmatrix} -6.166* \\ e-06 \\ -1.23e-05 \end{bmatrix}_{2\times 1} \checkmark$$

$C = \sigma(B)$ $\qquad \therefore\ \frac{\partial C}{\partial B} = C(1-C)$

$$\frac{\partial L}{\partial B} = \frac{\partial L}{\partial C} * \frac{\partial C}{\partial B} = C(1-C)\,\overline{\frac{\partial L}{\partial C}} = \begin{bmatrix} -2.799e-10 \\ -3.439e-15 \end{bmatrix}$$

$\Rightarrow$ $B = A + ba$

$$\frac{\partial B}{\partial ba} = 1 \qquad \frac{\partial L}{\partial ba} = \frac{\partial L}{\partial B} \cdot \frac{\partial B}{\partial ba} = \begin{bmatrix} -2.799e-10 \\ -3.439e-15 \end{bmatrix}$$

$$\frac{\partial B}{\partial A} = 1 \qquad \therefore \frac{\partial L}{\partial A} = \begin{bmatrix} -2.799e-10 \\ -3.439e-15 \end{bmatrix}$$

Because of Consistency.

$A = \omega_1^T x_1$

$$\frac{\partial A}{\partial \omega_1} = x_1 \qquad \frac{\partial A}{\partial \omega_1} = \frac{\partial L}{\partial A} \cdot \frac{\partial A}{\partial \omega_1} = x_1 \left(\frac{\partial L}{\partial A}\right)^T$$

$$= \begin{bmatrix} 2 \\ 3 \end{bmatrix} \begin{bmatrix} -2.799e-10, & -3.439e-15 \end{bmatrix}$$

$$\frac{\partial L}{\partial \omega_1} = \begin{bmatrix} -5.598e-10 & -6.879e-15 \\ -8.397e-10 & -1.0319e-14 \end{bmatrix}$$

**Problem 6.3**

**A)** The eigenvalues of the Hessian characterize the local curvature of the loss which, for example, determine how fast models can be optimized via first-order methods (atleast for convex problems), and is also conjectured to influence the generalization properties.

So if the training error decreases so does the fraction of negative eigen value decreases. It happens because we know that with the training error decreases, we close toward the minimum where we get all positive eigenvalues.

## B) In Gradient Descent:

Step of this method points in the right direction close to a saddle point and if an eigenvalue is positive then step will move away in the direction of eigenvector and achieve minimum value. So eigenvalues are very small hence the steps taken are very small. Gradient descent drawback is the small size of the step along with each eigenvector.

## In Newton method:

Newton methods solve the slowness problem by rescaling the gradients in each direction with the inverse of the corresponding eigenvalue. However, this approach can result in moving in the wrong direction. So for the negative eigenvalues the newton step moves in the direction opposite to the gradient descent step, and thus moves in the direction of increasing error.

## In Trust region:

In Trust Region if the minimum eigenvalues is very low then the damping factor has to be large and this results in potentially small step size.

**C)** To achieve the optimal rescaling as in the case of Newton's Method the gradient descent was scaled along each eigen direction by $1/\lambda i$ where $\lambda i$ represents eigenvalues. This also preserves the direction, which was not possible in the case of Newton's Method .

# Team Members:

Priyanka Upadhyay, Email Id : s8prupad@stud.uni-saarland.de
Gopal Bhattrai, Email Id : gobh00001@stud.uni-saarland.de

 Tutor: Redion Xhepa

 Matriculation Number :

 Priyanka : 2581714
 Gopal : 7013547