

CS 7015 Assignment - 2

Total sheet = 6

Name: Anubhav Agarwal
Roll no.: BE16B002

PAGE:
DATE:

(1)

Part - B

b. Sigmoid ; $f(x) = 1/(1+e^{-x})$

By quotient rule,

$$f'(x) = \frac{(1+e^{-x})(0) - 1 \times (-e^{-x})}{(1+e^{-x})^2} = \frac{e^{-x}}{(1+e^{-x})^2}$$

$$f'(x) = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1+e^{-x}-1}{(1+e^{-x})^2} = \frac{1}{1+e^{-x}} - \frac{1}{(1+e^{-x})^2}$$

$$f'(x) = \frac{1}{1+e^{-x}} \left(1 - \frac{1}{1+e^{-x}} \right) = f(x)(1-f(x))$$

$$\therefore \text{Gradient of Sigmoid } \sigma'(x) = \sigma(x)(1-\sigma(x))$$

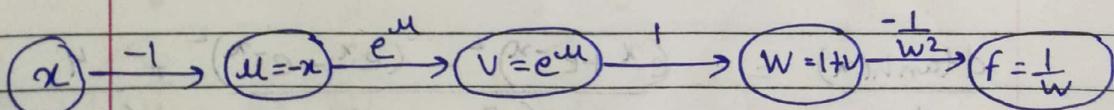
We can use a computation graph to represent a composite function as a network of nodes, where each node is a function.

Also, each edge represents the change in output node w.r.t input node i.e., the partial derivatives.

$$\sigma(x) = \frac{1}{1+e^{-x}}. \text{ Let } u = -x, v = e^u, w = 1+v, f = \frac{1}{w}. \text{ Here, } u, v, w \in$$

f are newly-defined intermediate variables.

The computational graph is given by,



$$\frac{\partial u}{\partial x} = -1; \quad \frac{\partial v}{\partial u} = e^u; \quad \frac{\partial w}{\partial v} = 1; \quad \frac{\partial f}{\partial w} = -\frac{1}{w^2}$$

To obtain the gradient of the sigmoid, we perf

- form reverse-mode differentiation (backpropagation) on the computational graph.

$$\begin{array}{c}
 \text{Graph: } \\
 \begin{array}{ccccccc}
 \textcircled{1} & \frac{\partial F}{\partial u} = -e^u & \leftarrow e^u & \textcircled{2} & \frac{\partial F}{\partial v} = -1 & \leftarrow 1 & \textcircled{3} \frac{\partial F}{\partial w} = -1 \\
 \downarrow -1 & & & & & & \\
 \textcircled{4} & \frac{\partial F}{\partial x} = \frac{e^{-x}}{w^2} & & & \frac{\partial F}{\partial x} = \frac{e^{-x}}{w^2} & & \textcircled{5} \frac{\partial F}{\partial x} = \frac{1}{1+e^{-x}}
 \end{array}
 \end{array}$$

Substituting u & w in terms of x

$$\frac{\partial F}{\partial x} = \frac{e^{-x}}{(1+e^{-x})^2}$$

This is same as the derivative obtained by using quotient rule.

And the decomposition used $\Rightarrow [g(x) = f(w(v(u(x)))]$

$$\text{Chain rule: } \frac{\partial f}{\partial x} = \frac{\partial f}{\partial F} \times \frac{\partial F}{\partial w} \times \frac{\partial w}{\partial v} \times \frac{\partial v}{\partial u} \times \frac{\partial u}{\partial x}$$

2. Hyperbolic Tangent.

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

By quotient rule,

$$\begin{aligned}
 \frac{d \tanh(x)}{dx} &= \frac{(e^x + e^{-x})(e^x + e^{-x}) - (e^x - e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2} \\
 &\quad \square (1 *)
 \end{aligned}$$

$$\begin{aligned}
 \frac{d \tanh(x)}{dx} &= \frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2} = \frac{4}{(e^x + e^{-x})^2}
 \end{aligned}$$

$$(a+b)^2 - (a-b)^2 = 4ab ; a = e^x \text{ & } b = e^{-x}$$

$$\therefore \frac{d \tanh(x)}{dx} = \frac{4}{(e^x + e^{-x})^2}$$

In (1*), On splitting the fraction we obtain,

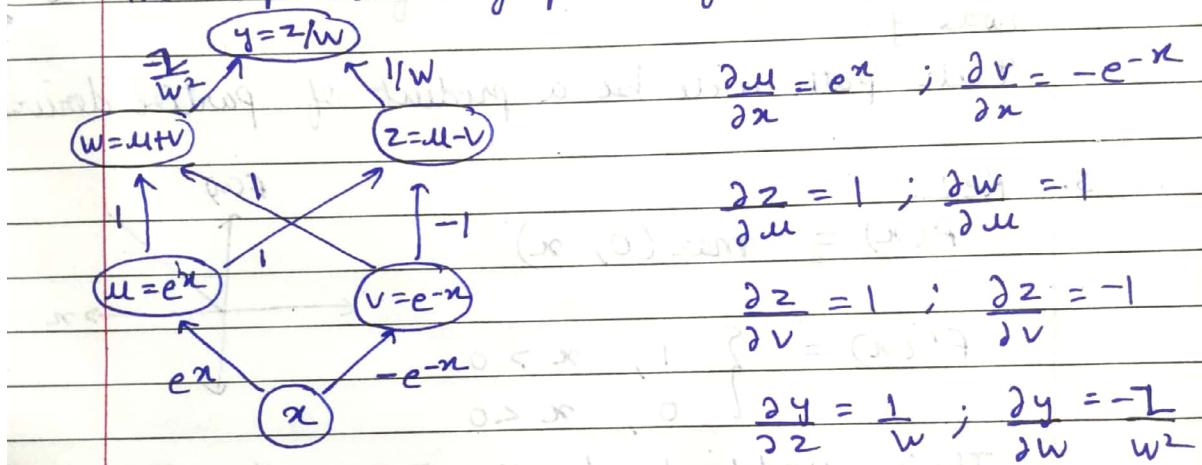
$$\frac{d \tanh(x)}{dx} = 1 - \left(\frac{e^x - e^{-x}}{e^x + e^{-x}} \right)^2 = 1 - \tanh^2(x)$$

Let's define new intermediate variables,

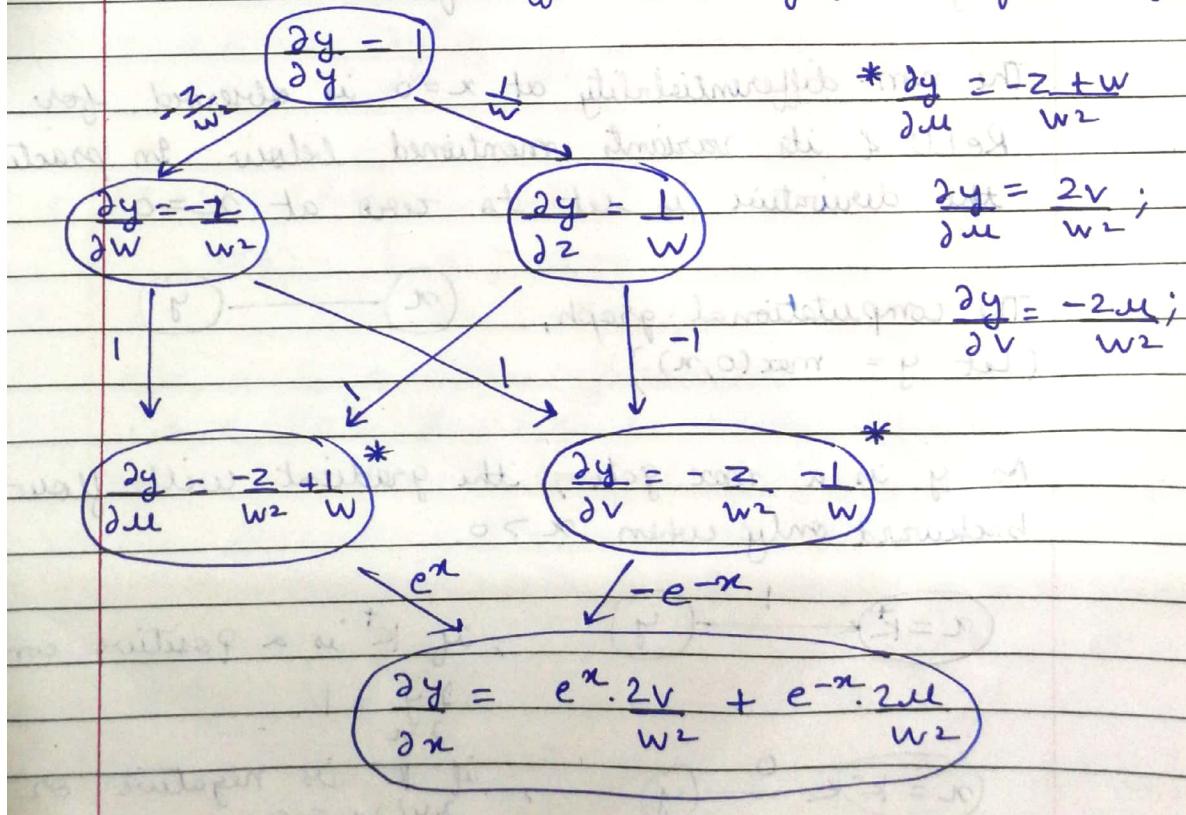
$$u = e^x, v = e^{-x}, w = u+v, z = u-v \text{ and}$$

$$y = \frac{z}{w}.$$

The computational graph is given by,



The reverse-mode differentiation graph is given by,



Substituting u , v & w in $\frac{\partial y}{\partial x}$,

$$\frac{\partial y}{\partial x} = \frac{e^x \cdot 2e^{-x} + e^{-x} \cdot 2e^x}{(e^{-x} + e^x)^2} = \frac{1}{(e^{-x} + e^x)^2}$$

PAGE: / /
DATE: / /

This is same as the derivative obtained by using quotient rule.

And the decomposition used \Rightarrow

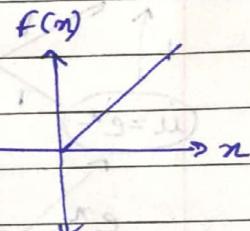
$$\tanh(x) = w_y(z(u(x), v(x)), w(u(x), v(x)))$$

The chain rule will involve the summation of all 'paths' (4 in this case) from node x to node y .

Each path will be a product of partial derivatives

3. ReLU

$$f(x) = \max(0, x)$$

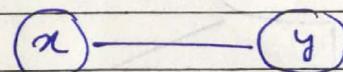


$$f'(x) = \begin{cases} 1, & x > 0 \\ 0, & x \leq 0 \end{cases}$$

It is undefined at $x=0$ as the left and right derivatives are not equal.

The non-differentiability at $x=0$ is observed for ReLU & its variants mentioned below. In practice, the derivative is set to zero at $x=0$.

The computational graph,
(Let $y = \max(0, x)$)



As y is a max gate, the gradient will flow backward only when $x > 0$.

$$(x = k^+) \xrightarrow{1} (y) \quad , \text{ if } k^+ \text{ is a positive constant}$$

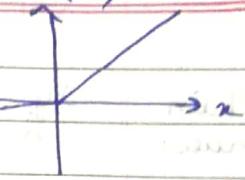
$$\frac{\partial y}{\partial x} = 1$$

$$(x = k^-) \xrightarrow{0} (y) \quad , \text{ if } k^- \text{ is negative or zero}$$

$$\frac{\partial y}{\partial x} = 0$$

4. Leaky ReLU

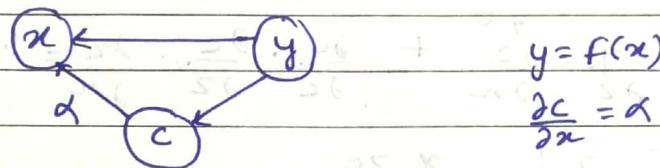
$$f(x) = \begin{cases} x, & x > 0 \\ \alpha x, & \text{otherwise}; \alpha \geq 0 \end{cases}$$



Here, α is a positive constant set by the user.

The computational graph is given by :

(Let $c = \alpha x$ be an intermediate variable)



$$\text{if } x > 0, \quad \begin{array}{c} x \leftarrow 1 \\ \alpha \leftarrow 0 \\ c \leftarrow x \\ y \leftarrow c \end{array} \quad \therefore f'(x) = \begin{cases} 1, & x > 0 \\ \alpha, & \text{otherwise} \end{cases}$$

$$\text{if } x < 0, \quad \begin{array}{c} x \leftarrow 0 \\ \alpha \leftarrow 1 \\ c \leftarrow 0 \\ y \leftarrow c \end{array} \quad \therefore f'(x) = \begin{cases} 1, & x > 0 \\ \alpha, & \text{otherwise} \end{cases}$$

Chain rule :

$$(i) x > 0 \quad \frac{\partial y^*}{\partial x} = \frac{\partial y}{\partial x} + \frac{\partial y}{\partial c} \cdot \frac{\partial c}{\partial x} = 1$$

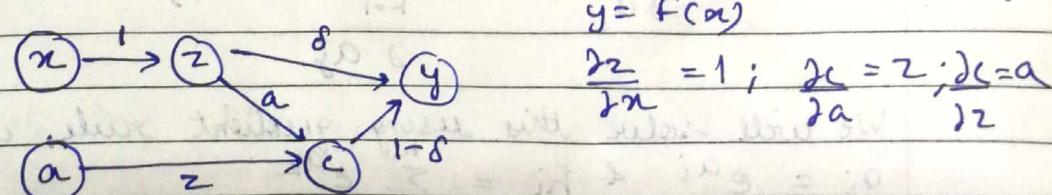
$$(ii) x < 0 \quad \frac{\partial y^*}{\partial x} = \frac{\partial y}{\partial x} + \frac{\partial y}{\partial c} \cdot \frac{\partial c}{\partial x} = 1 \cdot \alpha = \alpha$$

5. PReLU (Parametric ReLUs)

$$f(x) = \begin{cases} x, & x > 0 \\ ax, & \text{otherwise} \end{cases}$$

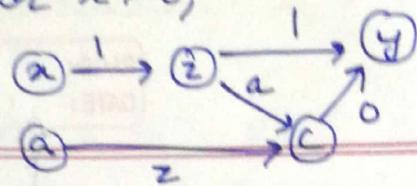
Here, a is a learned parameter.

Let $c = az$ & $z = x$ be newly defined intermediate variables. The computational graph is given by:



$$\text{if } x > 0, \delta = 1 \\ \text{else } \delta = 0$$

For $x > 0$,



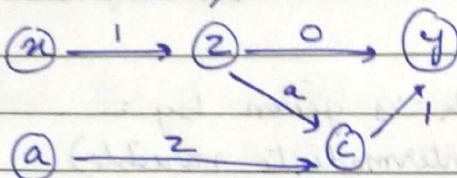
PAGE : / /

DATE : / /

Chain rule:

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial x} + \frac{\partial y}{\partial c} \cdot \frac{\partial c}{\partial z} \cdot \frac{\partial z}{\partial x} = 1 \cdot 1 = 1$$

For $x \leq 0$,



Function decomposition

$$f(x) = y(z(a), c(z(a), a))$$

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial x} + \frac{\partial y}{\partial c} \cdot \frac{\partial c}{\partial z} \cdot \frac{\partial z}{\partial x} = 1 \cdot a \cdot 1 = a //$$

$$f'(x) = \begin{cases} 1, & x > 0 \\ a & \text{otherwise} \end{cases}$$

2. Softmax: The function takes an N -dimensional vector of arbitrary real values and produces another N -dimensional vector with real values in the range $(0, 1)$ that adds up to 1.

Notation:

$$s(a) = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_N \end{bmatrix} \rightarrow \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_N \end{bmatrix} \quad s(a) : \mathbb{R}^N \rightarrow \mathbb{R}^N$$

$\therefore N$ gradients w.r.t N

$$\text{where } s_j = \frac{e^{a_j}}{\sum_{k=1}^N e^{a_k}} \quad \forall j \in 1..N \quad \begin{array}{l} \text{variables need} \\ \text{to be calculated} \\ \text{i.e., } N \times N \text{ total} \\ \text{gradients} \end{array}$$

The partial derivative of the i^{th} output w.r.t the (a_j) j^{th} input is given by $\frac{\partial s_i}{\partial a_j}$.

$$\frac{\partial s_i}{\partial a_j} = \frac{\partial s_i}{\partial a_j} = \frac{\frac{\partial}{\partial a_j} e^{a_i}}{\sum_{k=1}^N e^{a_k}} \quad \left. \begin{array}{l} \text{This is } g(x) \text{ form} \\ h(x) \end{array} \right.$$

$$\frac{\partial a_j}{\partial a_j} = 1$$

We will solve this using quotient rule, where
 $g_i = e^{a_i}$ & $h_i = \sum_{k=1}^N e^{a_k}$

h. Leaky ReLU

$$f(x) = \begin{cases} x, & x > 0 \\ \alpha x, & \text{otherwise}, \alpha \geq 0 \end{cases}$$

$$h_i = \sum_{k=1}^N e^{a_k}$$

$$\frac{dh_i}{da_j} = e^{a_j} + j \in 1 \dots N$$

$$\frac{dg_i}{da_j} = \begin{cases} e^{a_j} & \text{if } i=j, \text{ otherwise } 0 \\ 0 & \text{otherwise} \end{cases}$$

① For $i=j$ case, the quotient rule gives us:

$$\frac{D_j s_i}{(i=j)} = \frac{\frac{d}{da_j} \sum_{k=1}^N e^{a_k}}{\frac{d}{da_j} e^{a_i}} = \frac{e^{a_i} \sum_{k=1}^N e^{a_k} - e^{a_j} e^{a_i}}{\left(\sum_{k=1}^N e^{a_k}\right)^2}$$

$$\frac{D_j s_i}{(i=j)} = \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}} - \frac{e^{a_j}}{\sum_{k=1}^N e^{a_k}} \cdot \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}} = s_i - s_j s_i$$

$$\therefore D_j s_i = s_i (1 - s_j) \quad (1)$$

② For $i \neq j$ case, the quotient rule gives us:

$$\frac{D_j s_i}{(i \neq j)} = \frac{\sum_{k=1}^N e^{a_k} \cdot 0 - e^{a_j} e^{a_i}}{\left(\sum_{k=1}^N e^{a_k}\right)^2} = -\frac{e^{a_j} e^{a_i}}{\sum_{k=1}^N e^{a_k} \sum_{k=1}^N e^{a_k}}$$

$$D_j s_i = -s_i s_j \quad (2)$$

From (1) & (2), $D_j s_i = \begin{cases} s_i (1 - s_j), & i=j \\ -s_i s_j, & i \neq j \end{cases}$

Let's try obtaining the computational graph,

Let $u_i = e^{a_i}$ for $i=1 \dots N$ be intermediate variables

$$t = \sum_{i=1}^N u_i$$

variables

Remember the softmax function is the following

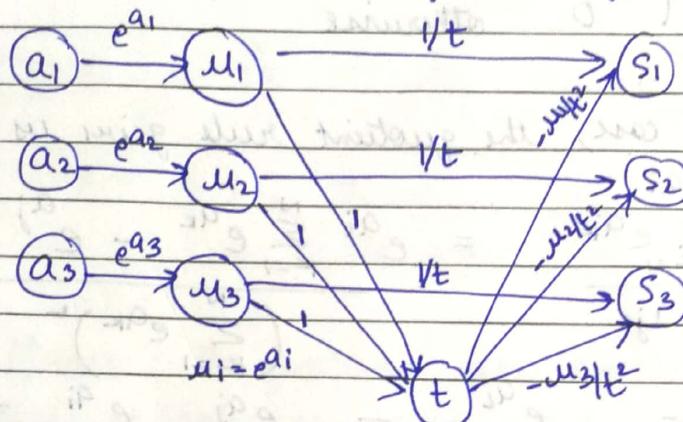
$$S(a) : \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_N \end{bmatrix} \rightarrow \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_N \end{bmatrix}$$

PAGE: / /
DATE: / /

Let $s_i = \frac{e^{a_i}}{t} + i = 1..N$. and a_i are the input variables.

We define $t = \sum_{i=1}^N a_i$ to reduce the number of paths in the graph.

For $N = 3$, (so too can be easily extended for large N), the graph is given by:



$$\frac{\partial t}{\partial a_i} = \frac{\partial \sum a_k}{\partial a_i} = 1$$

$$\frac{\partial s_i}{\partial a_i} = \frac{1}{t}$$

$$\frac{\partial s_i}{\partial t} = \frac{\partial \left(\frac{a_i}{t} \right)}{\partial t} = -\frac{a_i}{t^2}$$

If we densely connected all a_i 's with s_i 's we would have $n \times n$ connections. By introducing $t = \sum a_i$, we have n connections from a_i to t) + n (a_i to s_i) + n (t to s_i) = $3n$.

For $N = 3$, $n^2 = 3n$ will have same no. of connections. But for large N , this won't be the case. Also, the partial derivatives in the graph are simplified after introducing intermediate t .

$$\frac{\partial s_1}{\partial a_1} = \sum_{\text{all paths from } a_1 \text{ to } s_1} \text{Product of the PD's in path}$$

$$\frac{\partial S_1}{\partial a_1} = \frac{\partial S_1}{\partial u_1} \cdot \frac{\partial u_1}{\partial a_1} + \frac{\partial S_1}{\partial t} \cdot \frac{\partial t}{\partial u_1} \cdot \frac{\partial u_1}{\partial a_1}$$

$$\frac{\partial S_1}{\partial a_1} = \frac{1}{t} \cdot e^{a_1} + \left(-\frac{u_1}{t^2}\right) \cdot 1 \cdot e^{a_1}$$

$$\frac{\partial S_1}{\partial a_1} = \frac{e^{a_1}}{t} \left(1 - \frac{u_1}{t}\right) = \frac{u_1}{t} \left(1 - \frac{u_1}{t}\right) = s_1(1-s_1)$$

$$\therefore \boxed{\frac{\partial S_1}{\partial a_1} = s_1(1-s_1)}$$

follows the $\frac{\partial S_i}{\partial a_i} = s_i(1-s_i)$
($i=j$)

$$\frac{\partial S_3}{\partial a_2} = \sum_{\substack{\text{all paths} \\ \text{from } a_2 \text{ to } S_3}} \text{Product of PDs in the path}$$

$$\frac{\partial S_3}{\partial a_2} = \frac{\partial S_3}{\partial t} \cdot \frac{\partial t}{\partial u_2} \cdot \frac{\partial u_2}{\partial a_2} = -\frac{u_3}{t^2} \cdot 1 \cdot e^{a_2}$$

$$\frac{\partial S_3}{\partial a_2} = -\frac{u_3 \cdot u_2}{t^2} = -S_3 S_2$$

$$\therefore \boxed{\frac{\partial S_3}{\partial a_2} = -S_3 S_2}$$

follows the $\frac{\partial S_i}{\partial a_i} = -S_i S_j$
($i \neq j$)

\therefore The function decomposition of
 $y_i = s_i \Rightarrow y_i = s_i(u_i(a_i), t(u_1, \dots, u_N))$

Q.2 solved later

Q3. Hand-calculation of gradients

Notation (From Nielsen's book)

- 1) We'll use w_{jk}^l to denote the weight for the connection from the k^{th} neuron in the $(l-1)^{\text{th}}$ layer to the j^{th} neuron in the l^{th} layer.
 e.g., $w_{12}^3 = 0.1$, $w_{21}^2 = -0.2$

- 2) b_j^l for the bias of the j^{th} neuron in the l^{th} layer. a_j^l for the activation of the j^{th} neuron in the l^{th} layer.

3) we define an intermediate quantity, z_j^l , this is the weighted input to the j^{th} neuron in layer l .

PAGE NO. / /
DATE: / /

$$z_j^l = \sum_k w_{jk}^l \times a_k^{l-1} + b_j^l$$

$$a_j^l = \sigma(z_j^l)$$

4) w^l is the weight matrix of a layer, the entries of which are w_{jk}^l . Similarly, b^l , a^l , z^l , σ^l .

5) we define the error δ_j^l of neuron j in layer l by $\delta_j^l = \frac{\partial C}{\partial z_j^l}$, where C is the cost function

6) The equations being used,

i) $\delta^L = \frac{\partial C}{\partial a^L} \odot \sigma'(z^L)$; Error in last layer.

L : no. of layers, \odot is the hadamard product

ii) $\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$; Error in layer l given the error in layer $l+1$.

iii) $\frac{\partial C}{\partial b_j^l} = \frac{\partial C}{\partial z_j^l} \times \frac{\partial z_j^l}{\partial b_j^l} = \delta_j^l \times \frac{\partial (\sum_k w_{jk}^l a_k^{l-1} + b_j^l)}{\partial b_j^l}$

$\frac{\partial C}{\partial b_j^l} = \delta_j^l \times 1 = \delta_j^l$; Gradient of cost function w.r.t b_j^l is just δ_j^l .

iv) $\frac{\partial C}{\partial w_{jk}^l} = \frac{\partial C}{\partial z_j^l} \times \frac{\partial z_j^l}{\partial w_{jk}^l} = \delta_j^l \times \frac{\partial (\sum_k w_{jk}^l a_k^{l-1} + b_j^l)}{\partial w_{jk}^l}$

$\frac{\partial C}{\partial w_{jk}^l} = \delta_j^l \times a_k^{l-1}$

Basic steps in the algorithm:

- 1) Feedforward
- 2) Calculate output error
- 3) Backpropagate the errors
- 4) Gradient descent
- 5) Repeat

Step 1: Feedforward

$$z_1^2 = 0.2 \times 0.8 + 0.3 \times 0.7 + 0.4 \times 0.1 - 0.2$$

$$\underline{z_1^2 = 0.210}$$

PAGE: / /
DATE: / /

$$a_1^2 = \sigma(z_1^2) = \underline{0.552}$$

$$z_2^2 = \cancel{0.2} - 0.2 \times 0.8 + 0.1 \times 0.7 + 0.5 \times 0.1 - 0.2$$

$$\underline{z_2^2 = -0.240}$$

$$a_2^2 = \sigma(z_2^2) = \underline{0.440}$$

$$z_1^3 = 0.552 \times -0.3 + 0.440 \times 0.1 - 0.2$$

$$\underline{z_1^3 = -0.322}$$

$$a_1^3 = \sigma(z_1^3) = \underline{0.420}$$

Step 2: Calculate output error

$$\text{Loss function: } L(a^L) = (y - a^L)^2$$

$$L = (y - a_1^3)^2 = (0.9 - 0.42)^2 = \underline{0.2304}$$

Step 3: Backpropagate the error

(Eqn(i) in Notation)

$$\frac{\partial L}{\partial z_1^3} = \delta_1^3 = \frac{\partial L}{\partial a_1^3} \times \frac{\partial a_1^3}{\partial z_1^3} = \frac{\partial L}{\partial a_1^3} \times \sigma'(z_1^3)$$

$$\delta_1^3 = \frac{\partial (y - a_1^3)^2}{\partial a_1^3} \times \sigma(z_1^3) (1 - \sigma(z_1^3))$$

$$\delta_1^3 = (a_1^3 - y) \times \sigma(-0.322) (1 - \sigma(-0.322))$$

$$\underline{\delta_1^3 = -0.117} \quad (\text{substitute } a_1^3 = 0.42, y = 0.9 \text{ and } \sigma(-0.322) = 0.42)$$

From eqn (iii) in Notation,

$$\frac{\partial L}{\partial b_1^3} = \delta_1^3 = \underline{-0.117}$$

From eqn (iv) in Notation,

$$\frac{\partial L}{\partial w_{11}^3} = \delta_1^3 \times a_1^2 = -0.117 \times 0.552 = \underline{-0.065}$$

$$\frac{\partial L}{\partial w_{12}^3} = \delta_1^3 \times a_2^2 = -0.117 \times 0.440 = \underline{-0.051}$$

From eqn (ii) in Notation,

$$\frac{\partial L}{\partial z^2} = \delta^2 = (w^3)^T \beta^3 \odot \sigma'(z^2)$$

$$\delta^2 = \begin{bmatrix} \delta_1^2 \\ \delta_2^2 \end{bmatrix} = [w_{11}^3 \ w_{12}^3]^T [\delta_1^3] \odot \begin{bmatrix} \sigma'(z_1^2) \\ \sigma'(z_2^2) \end{bmatrix}$$

$$\delta^2 = \begin{bmatrix} \delta_1^2 \\ \delta_2^2 \end{bmatrix} = \begin{bmatrix} -0.3 \\ 0.1 \end{bmatrix} \begin{bmatrix} -0.117 \\ 0.240 \end{bmatrix} \odot \begin{bmatrix} 0.210 \\ -0.240 \end{bmatrix}$$

PAGE: _____
DATE: _____

$$\delta^2 = \begin{bmatrix} 0.035 \\ -0.012 \end{bmatrix} \odot \begin{bmatrix} 0.247 \\ 0.246 \end{bmatrix} = \begin{bmatrix} 0.009 \\ -0.003 \end{bmatrix}$$

$$\frac{\partial L}{\partial b_1^2} = \delta_1^2 = 0.009 \quad \frac{\partial L}{\partial b_2^2} = \delta_2^2 = -0.003$$

$$\frac{\partial L}{\partial w_{jk}^2} = \delta_j^2 \times a_k^1 \quad ; \quad j=1,2 \quad k=1,2,3$$

$$\frac{\partial L}{\partial w^2} = \begin{bmatrix} \delta_1^2 \cdot a_1^1 & \delta_1^2 \cdot a_2^1 & \delta_1^2 \cdot a_3^1 \\ \delta_2^2 \cdot a_1^1 & \delta_2^2 \cdot a_2^1 & \delta_2^2 \cdot a_3^1 \end{bmatrix}_{2 \times 3}$$

$$\frac{\partial L}{\partial w^2} = \begin{bmatrix} 0.009 \times 0.8 & 0.009 \times 0.7 & 0.009 \times 0.1 \\ -0.003 \times 0.8 & -0.003 \times 0.7 & -0.003 \times 0.1 \end{bmatrix}$$

$$\frac{\partial L}{\partial w^2} = \begin{bmatrix} 7.2 \times 10^{-3} & 6.3 \times 10^{-3} & 9 \times 10^{-4} \\ -2.4 \times 10^{-3} & -2.1 \times 10^{-3} & -3 \times 10^{-4} \end{bmatrix}$$

Step 4: Gradient descent

$$w_{jk}^l \rightarrow w_{jk}^l - \eta \frac{\partial L}{\partial w_{jk}^l}$$

$$b_j^l \rightarrow b_j^l - \eta \frac{\partial L}{\partial b_j^l}$$

For $\eta = 1$,

$$w_{11}^2 = 0.2 - 1 \times 7.2 \times 10^{-3} = 0.1928 \quad b_1^2 = -0.2 + 0.009$$

$$w_{12}^2 = 0.3 - 1 \times 6.3 \times 10^{-3} = 0.2937 \quad b_1^2 = -0.29$$

$$w_{13}^2 = 0.4 - 9 \times 10^{-4} = 0.3991 \quad b_2^2 = -0.2 + 0.003 = -0.197$$

$$b_2^2 = -0.197$$

$$w_{21}^2 = -0.2 + 1 \times 2.4 \times 10^{-2} = -0.1976$$

$$w_{22}^2 = 0.1 + 1 \times 2.1 \times 10^{-3} = 0.1021$$

$$w_{23}^2 = 0.5 + 1 \times 3 \times 10^{-4} = 0.5003$$

$$w_{11}^3 = -0.3 + 1 \times 0.065 = -0.235$$

$$w_{12}^3 = 0.1 + 1 \times 0.051 = 0.151$$

Step 5: Repeat

4

EIGB002

moulshaw

Forward pass with updated weights:

$$z_1^2 = 0.1928 \times 0.8 + 0.2937 \times 0.7 + 0.3991 \times 0.1 - 0.209$$

PAGE: 0.1

DATE: 1/1

$$z_1^2 = 0.1907$$

$$a_1^2 = \sigma(z_1^2) = 0.5475$$

$$z_2^2 = -0.1976 \times 0.8 + 0.1021 \times 0.7 + 0.5003 \times 0.1 - 0.197$$

$$z_2^2 = -0.2336$$

$$a_2^2 = 0.4419 = \sigma(z_2^2)$$

$$z_1^3 = 0.5475 \times -0.235 + 0.4419 \times 0.151 - 0.083$$

$$z_1^3 = -0.1449$$

$$a_1^3 = \sigma(z_1^3) = 0.4638$$

$$\text{Loss} = (0.9 - 0.4638)^2 = 0.1903$$

∴ After backpropagating + G.D after 1 training example update, we get

	Original	After 1 update	
\hat{y}	0.42	0.4638	Target = 0.9
L	0.2304	0.1903	

2. Gradient calculation of common loss functions

y_i : true output of neuron i , a_i^L : final activation of neuron i

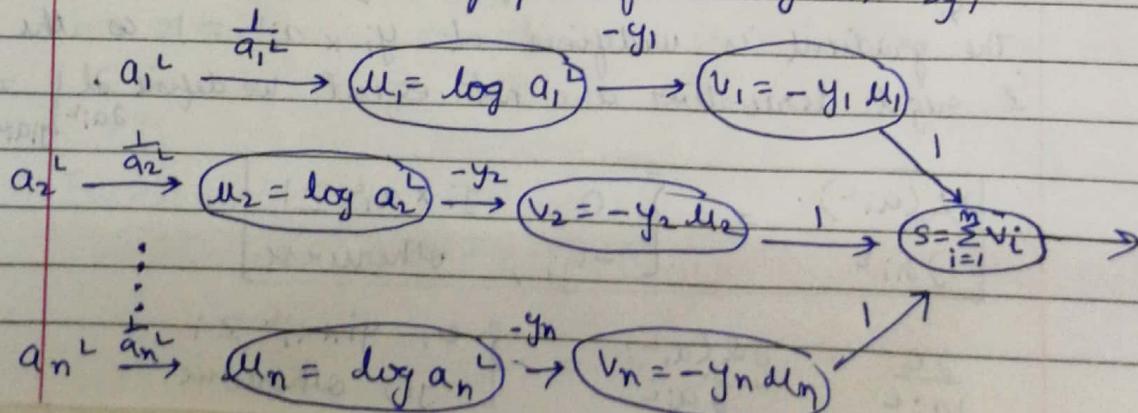
i) Cross entropy loss

$$C = - \sum_{i=1}^m y_i \log(a_i^L)$$

n is the no. of neurons in final activation layer.
 L is the last layer.

$$\frac{\partial C}{\partial a_i^L} = -\frac{y_i}{a_i^L}$$

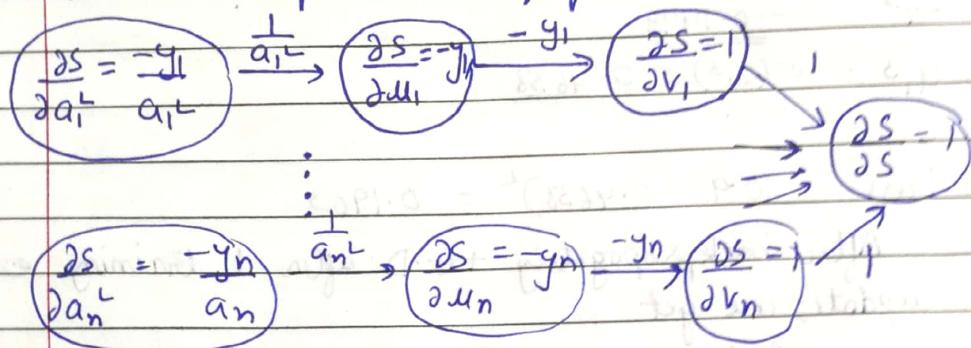
The computational graph of C is given by,



Let $u_i = \log a_i^L$, $v_i = -y_i u_i$ and
 $S = \sum_{i=1}^n v_i$ be intermediate variables.

The edges of the graph represent the change in output node w.r.t input node, i.e., partial derivative. $\frac{\partial u_i}{\partial a_i^L} = \frac{1}{a_i^L}$, $\frac{\partial v_i}{\partial u_i} = -y_i$, $\frac{\partial S}{\partial v_i} = 1$

To obtain the gradient of the cross-entropy loss with respect to output activations, we perform reverse-mode differentiation (backpropagation) on the computational graph.



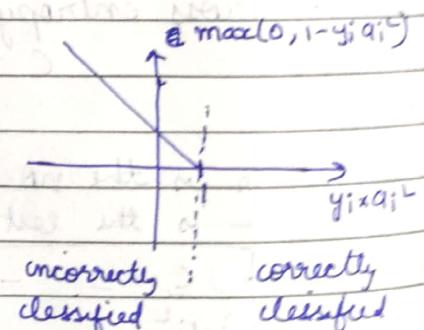
$\frac{\partial C}{\partial a_i^L} = \frac{\partial S}{\partial a_i^L} = -\frac{y_i}{a_i^L}$, this is same as what was obtained by differentiation earlier.

Decomposition of function used \Rightarrow

$$C(a_1^L, \dots, a_n^L) = S(v_1(u_1(a_1^L)), \dots, v_n(u_n(a_n^L)))$$

2) Hinge loss

$$\begin{aligned} C(a_1^L, \dots, a_n^L) &= \sum_{i=1}^n \max(0, 1 - y_i \cdot a_i^L) \\ &= \sum_{i=1}^n l_{y_i}(a_i^L) \end{aligned}$$



For a single neuron,

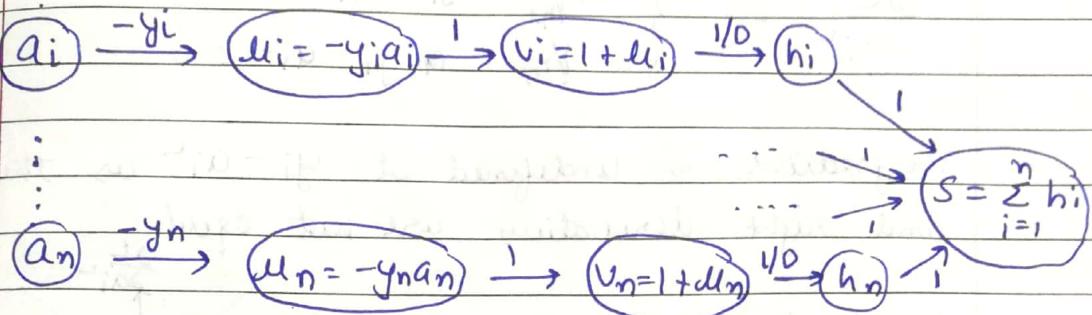
$$l_{y_i}(a_i^L) = \max(0, 1 - y_i \cdot a_i^L)$$

The gradient is undefined at $y_i \cdot a_i^L = 1$ as the left & right derivatives are not equal. We define $\frac{\partial l}{\partial a_i^L} \Big|_{y_i \cdot a_i^L = 1} = 0$.

$$\frac{\partial l_{y_i}(a_i^L)}{\partial a_i^L} = \begin{cases} 0, & y_i \cdot a_i^L \geq 1 \\ -y_i, & \text{otherwise} \end{cases}$$

$$\frac{\partial C}{\partial a_i^L} = \frac{\partial l_{y_i}(a_i^L)}{\partial a_i^L} = \begin{cases} 0, & y_i \cdot a_i^L \geq 1 \\ -y_i, & \text{otherwise} \end{cases}$$

Let $u_i = -y_i a_i$, $v_i = 1 + u_i$, $h_i = \max(0, v_i)$
and $S = \sum h_i$ be intermediate variables.
DATE: / /
The computational graph is given by,

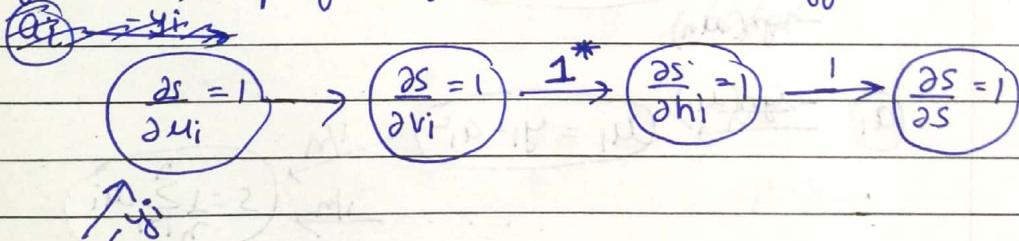


$$\frac{\partial u_i}{\partial a_i} = -y_i ; \frac{\partial v_i}{\partial u_i} = 1, \frac{\partial S}{\partial h_i} = 1$$

h_i is a max gate. The gradient will be 1 if $v_i > 0$ and 0 if $v_i < 0$

∴ if a_i (the predicted label) is same as the target y_i , i.e., it is correctly classified, ~~with high confidence~~ then ~~because~~ the gradient from $(v_i) \xrightarrow{?} (h_i)$ will be zero.

If $v_i > 0$, performing reverse-mode differentiation,



$$\frac{\partial S}{\partial a_i} = -y_i$$

$$\therefore v_i > 0, \text{i.e., } 1 - y_i a_i > 0 \Rightarrow y_i a_i < 1$$

we get $\frac{\partial L}{\partial a_i} = -y_i$

(* becomes 0)

$$\text{When } v_i \leq 0, \text{i.e., } 1 - y_i a_i \leq 0 \Rightarrow y_i a_i \geq 1$$

$$\text{we get } \frac{\partial L}{\partial a_i} = 0$$

We obtain same results as differentiation which is

$$\frac{\partial L}{\partial a_i} = \begin{cases} 0, & y_i a_i \geq 1 \\ -y_i, & \text{otherwise} \end{cases}$$

Functional decomposition:

$$C(a_1^L, \dots, a_n^L) = S(h_1(v_1(u_1(a_1^L))), \dots, h_n(v_n(u_n(a_n^L))))$$

3) L1 loss

$$C = \frac{1}{n} \sum_{i=1}^n |y_i - a_i^L|$$

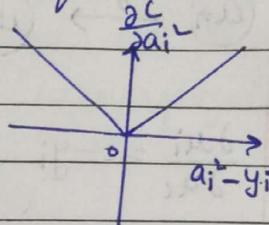
PAGE: / /
DATE: / /

$$\frac{\partial C}{\partial a_i^L} = \begin{cases} -1/n, & y_i - a_i^L > 0 \\ 1/n, & y_i - a_i^L < 0 \end{cases}$$

The gradient is undefined at $y_i = a_i^L$ as the left and right derivatives are not equal.

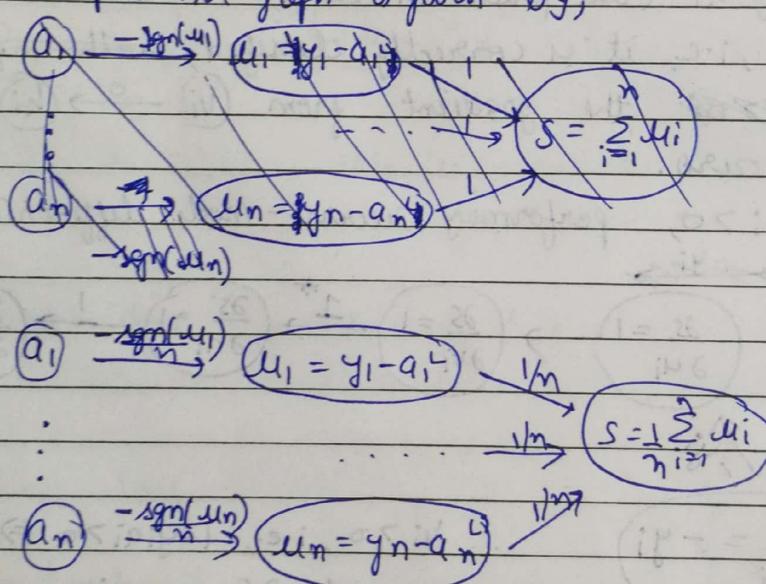
Shorthand: $\frac{\partial C}{\partial a_i^L} = -\frac{1}{n} \operatorname{sgn}(y_i - a_i^L)$

$\operatorname{sgn}(x) = \begin{cases} 0, x = 0 \\ 1, x > 0 \\ -1, x < 0 \end{cases}$



Let $u_i = y_i - a_i^L$ and $S = \frac{1}{n} \sum_{i=1}^n u_i$ be intermediate variables.

The computational graph is given by,



After reverse-mode differentiation,

$$\begin{aligned} \frac{\partial S}{\partial a_1} &= -\frac{\operatorname{sgn}(u_1)}{n} \xrightarrow{-\operatorname{sgn}(u_1)/n} \frac{\partial S}{\partial u_1} = 1 \xrightarrow{1/n} \frac{\partial S}{\partial S} = 1 \\ &\vdots \\ \frac{\partial S}{\partial a_n} &= -\frac{\operatorname{sgn}(u_n)}{n} \xrightarrow{-\operatorname{sgn}(u_n)/n} \frac{\partial S}{\partial u_n} = 1 \end{aligned}$$

$$\therefore \frac{\partial S}{\partial a_i} = -\frac{\operatorname{sgn}(u_i)}{n} = -\frac{\operatorname{sgn}(y_i - a_i^L)}{n}$$

We obtain the same result as differentiation.

5

The functional decomposition is given by,

$$C(a_1^L, \dots, a_n^L) = S(u_1(a_1^L), \dots, u_n(a_n^L))$$

PAGE:

DATE: / /

knowledge

BE16B002

4) Huber loss

$$(Wikipedia) \rightarrow H = \begin{cases} 0.5(y - f(x))^2, & |y - f(x)| \leq \delta \\ \delta|y - f(x)| - \frac{\delta^2}{2}, & \text{otherwise} \end{cases}$$

For our notation, $f(x) = a_i^L$. Let's define R_i as the residual. $R_i = y_i - f(x_i) = y_i - a_i^L$

$$l(R_i) = \begin{cases} 0.5 R_i^2, & |R_i| \leq \delta \\ \delta|R_i| - \frac{\delta^2}{2}, & |R_i| > \delta \end{cases}; C = \sum_{i=1}^n l(R_i)$$

cost function

For small residual R_i , the huber function $l(R_i)$ reduces to the L2 loss, and for large R_i it reduces to the robust L1 loss. The derivatives are continuous at $|R_i| = \delta$.

$$\frac{dl}{dR_i} = \begin{cases} R_i, & |R_i| \leq \delta \\ \delta \operatorname{sgn}(R_i), & |R_i| > \delta \end{cases} = \begin{cases} -\delta, & R_i < -\delta \\ R_i, & |R_i| \leq \delta \\ \delta, & R_i > \delta \end{cases}$$

$(\frac{dc}{da_i^L} = \frac{dc}{dl} \times \frac{dl}{da_i^L} = n \cdot \frac{dl}{da_i^L})$

$$\frac{dl}{da_i^L} = \frac{dc}{da_i^L} = \frac{dc}{dl} \times \frac{dl}{da_i^L} = -\frac{dc}{dl} = \begin{cases} \delta, & y_i - a_i^L < -\delta \\ R_i, & |y_i - a_i^L| \leq \delta \\ -\delta, & y_i - a_i^L > \delta \end{cases}$$

$$\frac{dc}{da_i^L} \cdot n = \begin{cases} \delta, & y_i - a_i^L < -\delta \\ a_i^L - y_i, & |y_i - a_i^L| \leq \delta \\ -\delta, & y_i - a_i^L > \delta \end{cases}$$

Let's try to decompose this function C using a computational graph. $C = \frac{1}{n} \sum_{i=1}^n l_{s,y_i}(R_i)$

For simplicity, let's first find the computational graph of $l_{s,y_i} = \begin{cases} 0.5 R_i^2, & |R_i| \leq \delta \\ \delta|R_i| - \frac{\delta^2}{2}, & |R_i| > \delta \end{cases}$

$$\text{Let } t = y_i - a_i^*; \quad u = \frac{t^2}{2}; \quad v = \delta |t|;$$

$$w = v - \frac{\delta^2}{2}; \quad x = (\star(1t) \leq \delta) \quad \boxed{\begin{matrix} \text{PAGE: } & u:p \\ \text{DATE: } & j \end{matrix}}$$

$y = (|t| > s) ? w : 0$; be intermediate variables.

We see that t is now same as the residual R_i .

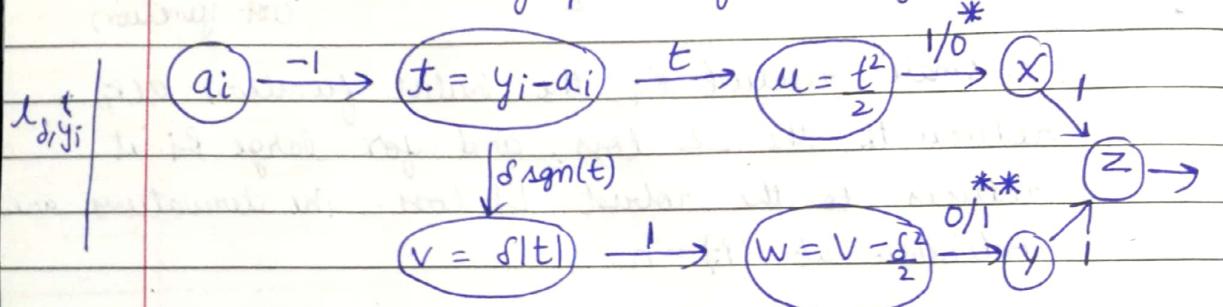
\times & \div are ternary operators

condition ? value_if_true : value_if_false.

This is required as $l_{S,y_i}(R_i)$ has different value when $|R_i| \leq s$ and when $|R_i| > s$

Also, define $Z = X + Y$

The computational graph is given by,



$$\frac{\partial t}{\partial a_i} = -1; \quad \frac{\partial u}{\partial t} = t; \quad \frac{\partial v}{\partial t} = \frac{\partial s|t|}{\partial t} = s \operatorname{sgn}(t)$$

$$\frac{\partial w}{\partial v} = 1 ; \quad \frac{\partial z}{\partial x} = 1 ; \quad \frac{\partial z}{\partial y} = 1 ;$$

$$x = (|t| \leq \delta) ? u : 0 \quad | \quad \text{if } |t| = |R| \leq \delta \text{ then,}$$

$$y = \begin{cases} 1 & |t| > \delta \\ 0 & \text{otherwise} \end{cases}, \quad x = u \text{ and } y = 0$$

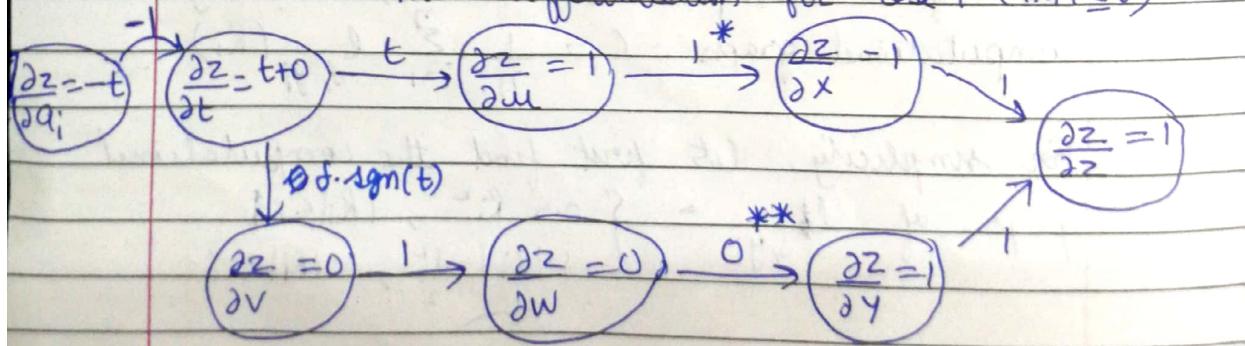
Case 2: If $|t| = |r_i| > \delta$ then,

$$x = 0 \text{ and } y = w$$

In Case 1: If the edge pd at $*$ will be $\frac{dx}{\|x\|} = \frac{1}{d}$

" * will be $\frac{\partial Y}{\partial W} = \underline{0}$

The reverse-mode differentiation for case 1 ($|R_i| \leq s$)



$$\text{We get } \frac{\partial z}{\partial a_i} = -t = -(y_i - a_i) = a_i - y_i$$

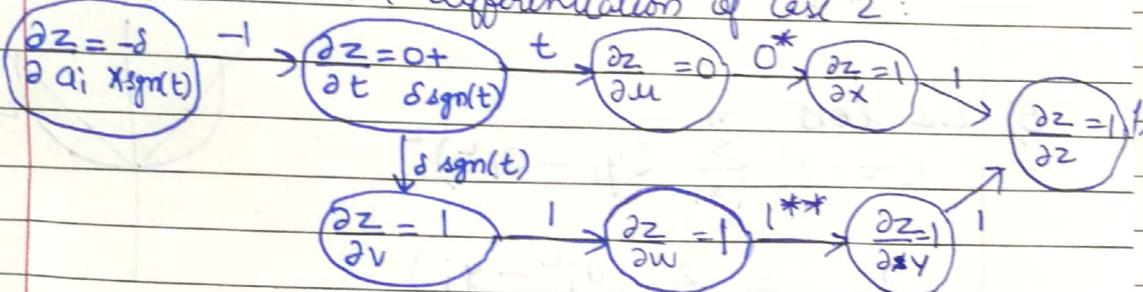
PAGE: / /
DATE: / /

i.e., when $|R_i| \leq \delta \Rightarrow \frac{\partial z}{\partial a_i} = a_i - y_i$

In case 2, $|R_i| > \delta$ i.e., $R_i \geq \delta$ or $R_i \leq -\delta$

$$\frac{\partial x}{\partial u} = 0 \quad \frac{\partial y}{\partial w} = 1$$

The reverse-mode differentiation of case 2:



$$\text{We get } \frac{\partial z}{\partial a_i} = -\delta \operatorname{sgn}(t) = \begin{cases} -\delta, & t > 0 \\ \delta, & t < 0 \end{cases}$$

$$t = R_i \Rightarrow \begin{cases} -\delta, & R_i > 0 \\ \delta, & R_i < 0 \end{cases}$$

But, we know that for case 2, $R_i \geq \delta$ or $R_i \leq -\delta$

$$\therefore \frac{\partial z}{\partial a_i} = \begin{cases} -\delta, & R_i > \delta \\ \delta, & R_i < -\delta \end{cases} \quad \textcircled{2}$$

Combining $\textcircled{1}$ from case i & $\textcircled{2}$ from case ii,

$$\frac{\partial z}{\partial a_i} = \frac{\partial l_{s,y_i}}{\partial a_i} = \begin{cases} \delta, & R_i < -\delta \\ a_i - y_i, & |R_i| \leq \delta \\ -\delta, & R_i > \delta \end{cases}$$

$$\frac{\partial l_{s,y_i}}{\partial a_i} = \begin{cases} \delta, & y_i - a_i < -\delta \\ a_i - y_i, & |y_i - a_i| \leq \delta \\ -\delta, & y_i - a_i > \delta \end{cases}$$

$$\frac{\partial C}{\partial a_i} = \frac{\partial C}{\partial l_{s,y_i}} \cdot \frac{\partial l_{s,y_i}}{\partial a_i} = \frac{\partial}{\partial l_{s,y_i}} \left(\sum_{i=1}^n \frac{1}{n} \frac{\partial l_{s,y_i}}{\partial a_i} \right) \cdot \frac{\partial l_{s,y_i}}{\partial a_i}$$

$$\frac{\partial C}{\partial a_i} = \frac{1}{n} \cdot 1 \cdot \frac{\partial l_{s,y_i}}{\partial a_i} = \frac{1}{n} \frac{\partial l_{s,y_i}}{\partial a_i}$$

$$\frac{\partial C}{\partial a_i} \cdot n = \begin{cases} 8, & y_i - a_i^L < -\delta \\ a_i^L - y_i, & |y_i - a_i^L| \leq \delta \\ -8, & y_i - a_i^L > \delta \end{cases}$$

PAGE:
DATE:

This is same as what we obtain using normal differentiation.

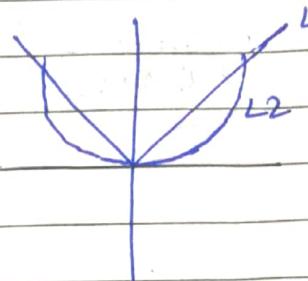
The functional decomposition of the Huber loss:

$$C(a_1^L, \dots, a_n^L) = \sum_{i=1}^n z_i (x_i(\mu_i(t_i(a_i^L))) + y_i(w_i(v_i(t_i(a_i^L)))))$$

5) L2 loss

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

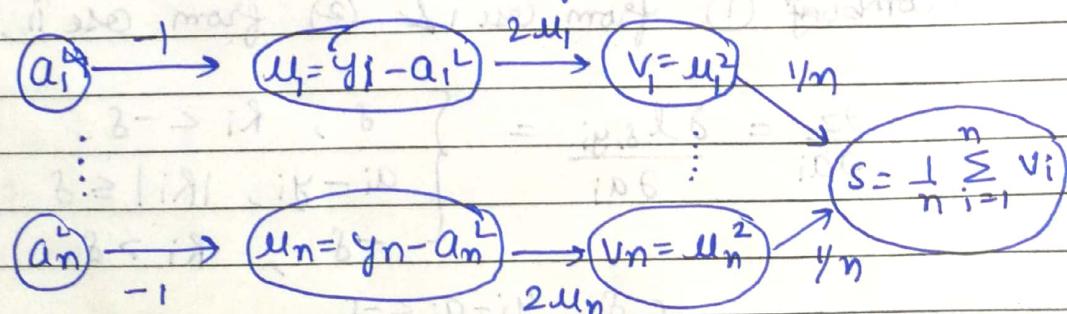
$$\frac{\partial C}{\partial a_i^L} = \frac{2}{n} (a_i^L - y_i)$$



The L1 loss is constant irrespective of whether a_i^L is close to or far from y_i . Whereas, L2 loss decreases more rapidly when a_i^L is far from y_i .

The L2 loss decreases slowly compared to L1 loss when a_i^L is close to y_i .

The computational graph of L2 loss is given by,



$$\frac{\partial \mu_i}{\partial a_i} = -1; \frac{\partial v_i}{\partial a_i} = 2\mu_i; \frac{\partial S}{\partial v_i} = \frac{1}{n}$$

$$\text{Here, } \mu_i = y_i - a_i^L; v_i = \mu_i^2; S = \frac{1}{n} \sum_{i=1}^n v_i$$

are intermediate variables

The reverse-mode differentiation on the computational graph is given by,

$$\begin{array}{c}
 \frac{\partial S}{\partial a_1} = \frac{-2u_1}{n} \quad -1 \quad \frac{\partial S}{\partial u_1} = \frac{2u_1}{n} \quad 2u_1 \\
 \vdots \\
 \frac{\partial S}{\partial a_n} = \frac{-2u_n}{n} \quad -1 \quad \frac{\partial S}{\partial u_n} = \frac{2u_n}{n} \quad 2u_n
 \end{array}
 \rightarrow
 \begin{array}{c}
 \frac{\partial S}{\partial v_1} = \frac{1}{n} \quad 1/n \\
 \vdots \\
 \frac{\partial S}{\partial v_n} = \frac{1}{n} \quad 1/n
 \end{array}
 \rightarrow
 \frac{\partial S}{\partial S} = 1$$

PAGE:

DATE:

$$\text{We get, } \frac{\partial S}{\partial a_i} = -\frac{2u_i}{n} = -\frac{2}{n}(y_i - a_i^L) = \frac{2(a_i^L - y_i)}{n}$$

This is same as normal differentiation.

The functional decomposition of L2 loss is given by,

$$C(a_1^L, \dots, a_n^L) = S(v_1(u_1(a_1^L)), \dots, v_n(u_n(a_n^L)))$$

6) Cosine Similarity

The cosine similarity between two vectors y and \hat{y} is given by

$$\text{similarity} = \frac{y \cdot \hat{y}}{\|y\| \|\hat{y}\|}$$

Assume the final activations a^L have been normalized. Also, $y_{\text{target}} = y$ is normalized. This is done to get clean results.

The loss function can be defined as,

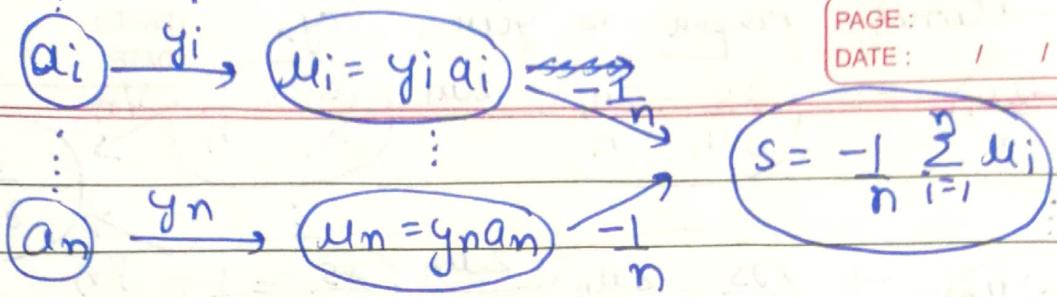
$$C = -\frac{1}{n} \sum_{i=1}^n y_i \cdot a_i^L$$

To measure similarity & thus minimizing the loss, a minus sign is added.

$$\boxed{\frac{\partial C}{\partial a_i^L} = -\frac{y_i}{n}}$$

Let $u_i = y_i a_i^L$; $S = -\frac{1}{n} \sum_{i=1}^n u_i$ be intermediate variables.

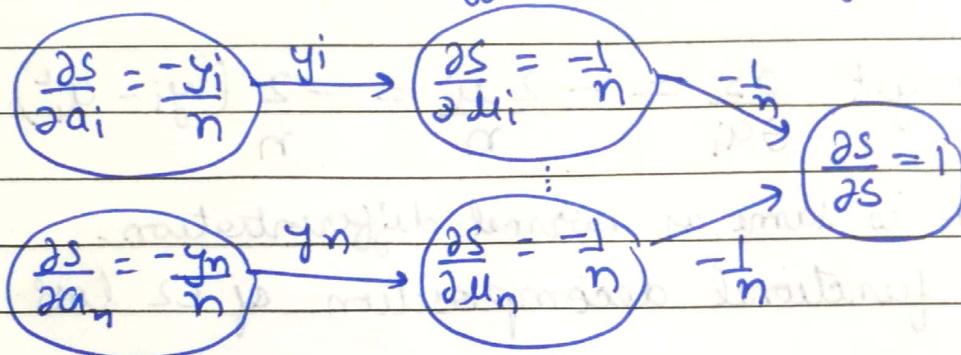
The computational graph is given by,



PAGE : / /

DATE : / /

The reverse-mode differentiation is given by,



We get, $\frac{\partial s}{\partial a_i} = \boxed{\frac{\partial c}{\partial a_i} = -\frac{y_i}{n}}$ i.e., same as differentiation

The functional decomposition is given by,

$$c(a_1^L, \dots, a_n^L) = S(u_1(a_1^L), \dots, u_n(a_n^L))$$

Total 6 sheet bundle

Ansulthan Agarwal

BE16B002

6