

21 Aug '22

HW 1

adarshs@comp.nus.edu.sg

Q.1 Solve gradients for:

1) Sigmoid  $f^h$ :  $f(x) = \frac{1}{1 + e^{-x}}$

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

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

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

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

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

$$= f(x) [1 - f(x)]$$

2) Softmax  $f(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$ ,  $1 \leq i \leq n$

$$f(x_i) = \frac{e^{x_i}}{(e^{x_1} + e^{x_2} + \dots + e^{x_i} + \dots + e^{x_n})}$$

$$\frac{\partial f}{\partial x_i} = \frac{e^{x_i} \cdot \left(\sum_{j=1}^n e^{x_j}\right) + e^{x_i} \cdot (e^{x_i})}{\left(\sum_{j=1}^n e^{x_j}\right)^2}$$

$$= \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} + \left( \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} \right)^2$$

$$= f(x_i) (1 + f(x_i))$$

3) Softplus activation:  $f(x) = \frac{1}{\beta} \cdot \ln(1 + e^{\beta x})$

$$\frac{\partial f}{\partial x} = \frac{1}{\cancel{\beta}} \times \frac{1}{(1 + e^{\beta x})} \times e^{\beta x} \times \beta$$

$$= \frac{e^{\beta x}}{1 + e^{\beta x}}$$

$$= \frac{1}{1 + e^{-\beta x}}$$


---

Q2 Solve gradients & do a shape check.

1)  $f(\vec{x}) = \vec{x}^T (\vec{A} \vec{x} + \vec{z})$

Assuming

$$\begin{aligned} \dim(\vec{x}) &= n \times 1 \\ \dim(\vec{A}) &= n \times n \\ \dim(\vec{z}) &= n \times 1 \end{aligned}$$

Shape check:

$$\begin{aligned} \dim(\vec{A} \vec{x}) &= n \times 1 \\ \dim(\vec{x}^T) &= 1 \times n \\ \dim(\vec{A} \vec{x} + \vec{z}) &= n \times 1 \\ \Rightarrow \dim[\vec{x}^T (\vec{A} \vec{x} + \vec{z})] &= 1 \times 1 = \text{scalar} \end{aligned}$$

Simplifying the expression:

$$f(x) = \bar{x}^T (\bar{A} \bar{x} + \bar{z})$$

$$= \bar{x}^T \cdot \bar{A} \bar{x} + \bar{x}^T \bar{z}$$

Given that  $\bar{x}$  and  $\bar{z}$  are both column vectors of  $\dim(n \times 1)$ , we can write:

$$f(x) = \underbrace{\bar{x}^T \bar{A} \bar{x}}_U + \underbrace{\bar{z}^T \bar{x}}_{\text{(both scalars)}}$$

$$\frac{\partial U}{\partial \bar{x}} = \underbrace{\bar{z}}_{n \times 1} \left( \because \frac{\partial (\bar{a}^T \bar{x})}{\partial \bar{x}} = \bar{a} \right) \left( \begin{array}{l} \text{denominator} \\ \text{layout} \end{array} \right)$$

$$U = \bar{x}^T \bar{A} \bar{x}$$

Now, if

$$\bar{x}^T = [x_1 \ x_2 \ \dots \ x_n]$$

$$\bar{A} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1n} \\ a_{21} & a_{22} & - & - & - \\ \vdots & & & & \\ a_{n1} & a_{n2} & - & - & - a_{nn} \end{bmatrix}$$

Then,

$$\bar{A} \bar{x} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 + \dots + a_{2n}x_n \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n \end{bmatrix}$$

(dim =  $n \times 1$ )

$$\Rightarrow \bar{x}^T \bar{A} \bar{x} = x_1(a_{11}x_1 + a_{12}x_2 + \dots) + x_2(a_{21}x_1 + a_{22}x_2 + \dots) + \dots$$

$$\Rightarrow \vec{x}^T \bar{A} \vec{x} = x_1 \left( \sum_{i=1}^n a_{1i} x_i \right) + x_2 \left( \sum_{i=1}^n a_{2i} x_i \right) + \dots + x_n \left( \sum_{i=1}^n a_{ni} x_i \right)$$

$$(\text{dim} = 1 \times 1)$$

= scalar

$$= \sum_{j=1}^n x_j \left( \sum_{i=1}^n a_{ji} x_i \right)$$

$$\Rightarrow \underbrace{\frac{\partial v}{\partial \vec{x}}}_{\substack{n \times 1 \\ \text{(denominator layout)}}} = \frac{\partial (\vec{x}^T \bar{A} \vec{x})}{\partial \vec{x}} = \begin{bmatrix} \partial v / \partial x_1 \\ \partial v / \partial x_2 \\ \vdots \\ \partial v / \partial x_n \end{bmatrix}_{n \times 1}$$

$$= \begin{bmatrix} (a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n) + (a_{11}x_1 + a_{21}x_2 + a_{31}x_3 + \dots + a_{n1}x_n) \\ (a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n) + (a_{12}x_1 + a_{22}x_2 + a_{32}x_3 + \dots) \\ \vdots \\ \vdots \end{bmatrix}$$

$$(\text{dim} = n \times 1)$$

$$= \bar{A} \vec{x} + A^T \vec{x} = \underbrace{(\bar{A} + A^T)}_{n \times n} \underbrace{\vec{x}}_{n \times 1}_{n \times 1}$$

$$\therefore \underbrace{\frac{\partial f}{\partial \vec{x}}}_{n \times 1} = \frac{\partial v}{\partial \vec{x}} + \frac{\partial v}{\partial \vec{x}} = \underbrace{(\bar{A} + \bar{A}^T)}_{n \times 1} \vec{x} + \underbrace{\vec{z}}_{n \times 1}$$

$$2) \text{ L2 loss: } L(\omega) = \frac{1}{2} (\underbrace{\bar{\omega}^T \bar{x} - y}_v)^2$$

using denominator layout.  $v = \text{scalar}$

$$\dim(y) = 1 \times 1 = \text{scalar}$$

$$\dim(\bar{x}) = n \times 1$$

$$\dim(\bar{\omega}) = n \times 1$$

$$\underbrace{\frac{\partial L}{\partial \bar{\omega}}}_{n \times 1} = \frac{1}{2} \times 2 \times v \times \frac{\partial v}{\partial \bar{\omega}}$$

$$\frac{\partial v}{\partial \bar{\omega}} = \frac{\partial (\bar{\omega}^T \bar{x} - y)}{\partial \bar{\omega}}$$

$$= \frac{\partial (\bar{\omega}^T \bar{x})}{\partial \bar{\omega}}$$

$$= \frac{\partial (\bar{x}^T \bar{\omega})}{\partial \bar{\omega}}$$

(given  $\bar{\omega}$  &  $\bar{x}$  are  $n \times 1$  column vectors)

$$= \bar{x} \quad \left( \because \frac{\partial (\bar{a}^T \bar{x})}{\partial \bar{x}} = \bar{a} \right)$$

$$\Rightarrow \underbrace{\frac{\partial L}{\partial \bar{\omega}}}_{n \times 1} = v \bar{x} = \underbrace{(\bar{\omega}^T \bar{x} - y)}_{\text{scalar}} \underbrace{\bar{x}}_{n \times 1}$$

3) L2 loss (multiple examples)

$$L(\omega) = \frac{1}{2} \|\bar{X} \bar{\omega} - \bar{y}\|^2$$

$$\begin{cases} \dim(y) = n \times 1 \\ \dim(\bar{\omega}) = m \times 1 \\ \dim(\bar{x}) = n \times m \end{cases}$$

$$= \frac{1}{2} (\bar{X} \bar{\omega} - \bar{y})^T (\underbrace{\bar{X} \bar{\omega} - \bar{y}}_v)$$

$$\dim(v) = n \times 1$$

$$= \frac{1}{2} v^T v$$

$$(\dim(L) = 1 \times 1 = \text{scalar})$$

$$\Rightarrow \underbrace{\frac{\partial L}{\partial \bar{w}}}_{m \times 1} = \underbrace{\frac{\partial \bar{u}}{\partial \bar{w}}}_{m \times n} \times \underbrace{\frac{\partial L}{\partial \bar{u}}}_{n \times 1}$$

$$= \frac{\partial (\bar{x} \bar{w} - \bar{y})}{\partial \bar{w}} \times \frac{1}{\cancel{2}} \times \cancel{2} \times \bar{u}$$

$$= \frac{\partial (\bar{x} \bar{w})}{\partial \bar{w}} \times (\bar{x} \bar{w} - \bar{y})$$

$$= \underbrace{\bar{x}^T}_{m \times n} \underbrace{(\bar{x} \bar{w} - \bar{y})}_{n \times 1} \left( \because \frac{\partial (\bar{A} \bar{x})}{\partial \bar{x}} = \bar{A}^T \right)$$

[using denominator layout]

$$\underbrace{\hspace{10em}}_{m \times 1}$$

Q.3 Solve for gradient  $\frac{\partial L}{\partial \bar{w}}$  & do slope check.

$$\bar{z} = \bar{w} \bar{x} + \bar{b}$$

$$L = \|\bar{z} - \bar{y}\|^2$$

$$\dim(\bar{y}) = m \times 1$$

$$\dim(\bar{b}) = m \times 1$$

$$\dim(\bar{x}) = n \times 1$$

$$\dim(\bar{w}) = m \times n$$

$$\bar{z} = \bar{w} \bar{x} + \bar{b}$$

$$\text{let } \bar{w} = \begin{bmatrix} \bar{w}_1 \\ \bar{w}_2 \\ \vdots \\ \bar{w}_m \end{bmatrix}$$

$$\left( \bar{w}_i = \underbrace{[w_{i1} \ w_{i2} \ w_{i3} \ \dots \ w_{in}]}_{1 \times n} \right)$$

$$\bar{z} = \begin{bmatrix} \bar{w}_1 \bar{x} + b_1 \\ \bar{w}_2 \bar{x} + b_2 \\ \vdots \\ \bar{w}_m \bar{x} + b_m \end{bmatrix} \quad m \times 1$$

$$\bar{z} - \bar{y} = \begin{bmatrix} \bar{w}_1 \bar{x} + b_1 - y_1 \\ \bar{w}_2 \bar{x} + b_2 - y_2 \\ \vdots \\ \bar{w}_m \bar{x} + b_m - y_m \end{bmatrix} \quad m \times 1$$

$$L = \|\bar{z} - \bar{y}\|^2 = \sum_{i=1}^m (\bar{w}_i \bar{x} + b_i - y_i)^2$$

$$\text{where } \bar{w}_i \bar{x} = w_{i1}x_1 + w_{i2}x_2 + w_{i3}x_3 + \dots + w_{in}x_n$$

$$\underbrace{\frac{\partial L}{\partial \bar{w}}}_{m \times h} = \begin{bmatrix} \frac{\partial L}{\partial w_{11}} & \frac{\partial L}{\partial w_{12}} & \dots & \frac{\partial L}{\partial w_{1n}} \\ \frac{\partial L}{\partial w_{21}} & \dots & \dots & \frac{\partial L}{\partial w_{2n}} \\ \vdots & & & \\ \frac{\partial L}{\partial w_{m1}} & \dots & \dots & \frac{\partial L}{\partial w_{mn}} \end{bmatrix}_{m \times h}$$

Solving just for  $\frac{\partial L}{\partial w_{11}}$  =

$$\frac{\partial \left( \sum_{i=1}^m (\bar{w}_i \bar{x} + b_i - y_i)^2 \right)}{\partial w_{11}}$$

$$= \frac{\partial \left[ (\bar{w}_1 \bar{x} + b_1 - y_1)^2 \right]}{\partial w_{11}}$$

$$= 2 (\bar{w}_1 \bar{x} + b_1 - y_1) \times \frac{\partial (\bar{w}_1 \bar{x} + b_1 - y_1)}{\partial w_{11}}$$

Similarly,

$$\frac{\partial L}{\partial w_{12}} = 2 (\bar{w}_1 \bar{x} + b_1 - y_1) x_2$$

$\Rightarrow$  in general

$$\frac{\partial L}{\partial w_{ij}} = 2 (\bar{w}_i \bar{x} + b_i - y_i) x_j$$

which can be rewritten as:

$$\frac{\partial L}{\partial \bar{w}} = 2 \left[ \underbrace{\underbrace{\bar{w}}_{m \times h} \underbrace{\bar{x}}_{h \times 1}}_{m \times 1} + \underbrace{\bar{b} - \bar{y}}_{m \times 1} \right] \underbrace{\bar{x}^T}_{1 \times h}$$

$m \times h$



Q.4

$$y = xw$$

$$w, x \in \mathbb{R}$$

$$\text{sample : } x=1, y=100$$

$$w_0 = 0$$

$$\text{loss } f^u = L(w) = \frac{1}{2}(wx - y)^2 \quad \left( \begin{array}{l} \because \text{only} \\ \text{one} \\ \text{example} \end{array} \right)$$

$$\begin{aligned} \frac{\partial L}{\partial w} &= \frac{1}{2} \times 2(wx - y) \times x \\ &= (wx - y)x \end{aligned}$$

1) Doing gradient descent

$$w_{i+1} = w_i - \alpha \left( \frac{\partial L}{\partial w} \right) \Big|_{w_i}$$

$$\underline{\alpha = 0.5}$$

$$\begin{aligned} w_1 &= 0 - (0.5) \times (-100) \\ &= 50 \end{aligned}$$

$$\begin{aligned} w_2 &= 50 - (0.5) \times (-50) \\ &= 75 \end{aligned}$$

$$\begin{aligned} w_3 &= 75 - (0.5) \times (-25) \\ &= 100 \end{aligned}$$

$$\begin{aligned} w_4 &= 100 - (0.5) \times (0) \\ &= 100 \end{aligned}$$

$$w_4 = w_3$$

$$\text{slope } \left| \frac{\partial L}{\partial w} \right|_{w=100} = 0$$

$\Rightarrow$  GD converges

$$\underline{\alpha = 1.5}$$

$$\begin{aligned} w_1 &= 0 - \frac{3}{2} \times (-100) \\ &= 150 \end{aligned}$$

$$\begin{aligned} w_2 &= 150 - \frac{3}{2} \times (50) \\ &= 75 \end{aligned}$$

$$\begin{aligned} w_3 &= 75 - \frac{3}{2} \times (-25) \\ &= 112.5 \end{aligned}$$

$$\begin{aligned} w_4 &= 112.5 - \frac{3}{2} \times (12.5) \\ &= 93.75 \end{aligned}$$

$$\begin{aligned} w_5 &= 93.75 - \frac{3}{2} \times (-6.25) \\ &= 103.125 \end{aligned}$$

slope  $\left| \frac{\partial L}{\partial w} \right|$  keeps getting smaller, suggesting we are moving closer to an optima.

$$\underline{\alpha = 2.5}$$

$$\begin{aligned} w_1 &= 0 - \frac{5}{2} \times (-100) \\ &= 250 \end{aligned}$$

$$\begin{aligned} w_2 &= 250 - \frac{5}{2} \times (150) \\ &= -125 \end{aligned}$$

$$\begin{aligned} w_3 &= -125 - \frac{5}{2} \times (-225) \\ &= 687.5 \end{aligned}$$

$$\begin{aligned} w_4 &= 687.5 - \frac{5}{2} \times (587.5) \\ &= -781.25 \end{aligned}$$

slope  $\frac{\partial L}{\partial w}$  oscillates b/w +ve & -ve with increasing magnitudes

This suggests  $\alpha$  is too large & GD won't converge here.

2) we have, from Q4.1,

$$\frac{\partial L}{\partial w} = (wx - y)x$$

$$\Rightarrow w_{i+1} = w_i - \alpha \frac{\partial L}{\partial w} \Big|_{w_i} \quad (i = \text{iterat}^n w.)$$

$$= w_i - \alpha (w_i x - y)x$$

for the solo example of  $x=1$  &  $y=100$

$$w_{i+1} = w_i - \alpha (w_i - 100)$$

$$\Rightarrow w_{i+1} = w_i (1 - \alpha) + 100\alpha$$

for GD to converge, we must have  $|\partial L / \partial w|$  decreasing through each iterat<sup>n</sup>.

$$\Rightarrow \left| \frac{\partial L}{\partial w} \Big|_{w_{i+1}} \right| < \left| \frac{\partial L}{\partial w} \Big|_{w_i} \right|$$

$$\Rightarrow |(w_{i+1}x - y)x| < |(w_i x - y)x|$$

$$\Rightarrow |w_{i+1} - 100| < |w_i - 100|$$

$$\Rightarrow |w_i(1 - \alpha) + 100\alpha - 100| < |w_i - 100|$$

$$\Rightarrow |(1 - \alpha)(w_i - 100)| < |w_i - 100|$$

for  $w_i \neq 100$ ,

$$|1 - \alpha| < 1$$

$$\Rightarrow -1 < 1 - \alpha < 1$$

$$\Rightarrow -2 < -\alpha < 0$$

$$\Rightarrow \boxed{2 > \alpha > 0}$$

We note that  $w_i = 100$  is the optima for the curve  $\left(\frac{\partial L}{\partial w} = 0\right)$ , at which point GD converges anyway.

Therefore,

GD converges for  $x \in (0, 2)$

GD diverges for  $x \in (-\infty, 0) \cup (2, \infty)$

And for  $x = \{0, 2\}$ ,

$x = 0$

$$w_0 = t$$

$$w_1 = t(1-x) + 100x$$

$(= w_0)$

$$w_2 = t$$

$(= w_0)$

⋮

$x = 2$

$$w_0 = t$$

$$w_1 = t(1-x) + 100x$$

$= -t + 200$

$$w_2 = (200-t)(1-2) + 200$$

$= t \quad (= w_0)$

⋮

Thus, for  $x \in \{0, 2\}$ , GD oscillates & neither converges nor diverges.

3) The above part proves conditions for  $x$  for any stopping threshold  $\varepsilon$   $(|w_{i+1} - w_i|)$ .