

19/8/22

Homework - (1) Part (1)

Date:

Q1. (1.) Sigmoid function.

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

ans. * Gradient of $f(x) \equiv$ derivative of $f(x) \equiv \frac{df(x)}{dx}$ w.r.t x

$$\therefore \frac{df(x)}{dx} = \frac{d\left[\frac{1}{1+e^{-x}}\right]}{dx} \rightarrow \textcircled{1}$$

* By applying quotient rule,

$$\frac{df}{dx} = \frac{d(1) \cdot (1+e^{-x}) - d(1+e^{-x}) \cdot (1)}{(1+e^{-x})^2}$$

$$= \frac{d(1+e^{-x})}{dx} \cdot \frac{(-1)}{(1+e^{-x})^2} \quad \cancel{e^{-x} \cdot (-1)}$$

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

$$= \frac{-1}{(1+e^{-x})^2} \cdot e^{-x} \cdot (-1) = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{(1+e^{-x}) - 1}{(1+e^{-x})^2} = \underline{\underline{f(x)[1-f(x)]}}$$

(2.) Softmax

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

ans. * For simplicity, $\sum_{j=1}^n \equiv \sum$

$$\Rightarrow f(x_i) = \frac{e^{x_i}}{\sum e^{x_j}}, \quad 1 \leq i \leq n$$

$$f'(x_i) = \frac{d[f(x_i)]}{dx_j} = \frac{d\left[\frac{e^{x_i}}{\sum e^{x_j}}\right]}{dx_j} \rightarrow (1)$$

★ By applying quotient rule,

$$f'(x_i) = \frac{\frac{d(e^{x_i})}{dx_j} \cdot \left(\sum e^{x_j}\right) - \frac{d\left[\sum e^{x_j}\right]}{dx_j} \cdot (e^{x_i})}{\left[\sum e^{x_j}\right]^2}$$

$$= \frac{e^{x_i} \cdot \frac{d(e^{x_i})}{dx_j} - e^{x_i} [0 + 0 + \dots + e^{x_j} + \dots + 0]}{\left[\sum e^{x_j}\right]^2} \rightarrow (2)$$

★ if $i=j$,

eq. (2) becomes,

$$= \frac{e^{x_i} \cdot \sum e^{x_i} - e^{x_i} \cdot e^{x_i}}{\left(\sum e^{x_i}\right)^2} = \frac{e^{x_i}}{\sum e^{x_i}} \left[\frac{1 - e^{x_i}}{\sum e^{x_i}} \right] = f(x_i) \cdot [1 - f(x_i)] \rightarrow (3)$$

NOTE:-

★★★ If we did a partial derivative w.r.t. to (x_i) , $(i \neq j)$ will not exist. [this is a generalized case].

★ if $i \neq j$, eq. (2) becomes,

$$= \frac{0 - e^{x_i} \cdot e^{x_j}}{\left(\sum e^{x_j}\right)^2} = -f(x_i) \cdot f(x_j) \rightarrow (4)$$

Combining (3) & (4), [eq. (4) is only for generalized case].

$$f'(x_i) = \begin{cases} f(x_i) \cdot (1 - f(x_i)) & , \text{ if } i = j \\ -f(x_i) \cdot f(x_j) & , \text{ if } i \neq j \end{cases}$$

(3) Softplus activation

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

ans. $f'(x) = \frac{d(f(x))}{dx} = \frac{d\left[\frac{1}{\beta} \cdot \ln(1 + e^{\beta x})\right]}{dx}$

$$= \frac{1}{\beta} \left[\frac{d(\ln(1 + e^{\beta x}))}{dx} \right] = \frac{1}{\cancel{\beta}} \left[\frac{1}{1 + e^{\beta x}} \right] \cdot [0 + e^{\beta x} \cdot \cancel{\beta}]$$

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

Q2. (1) $f(x) = x^T(Ax + z)$, $x, z \in \mathbb{R}^{n \times 1}$
 $A \in \mathbb{R}^{n \times n}$

ans. $f(x) = x^T Ax + x^T z$

Shape of $f(x)$ = $(1 \times n)(n \times n)(n \times 1) + (1 \times n)(n \times 1)$
 $= (1 \times 1) + (1 \times 1)$
 $\rightarrow (1 \times 1) \rightarrow \text{scalar}$

$$\begin{aligned} \frac{df(x)}{dx} &= \frac{d[x^T Ax + x^T z]}{dx} \\ &= \frac{d(x^T Ax)}{dx} + \frac{d(x^T z)}{dx} \\ &= \frac{d(x^T A)(x)}{dx} + \frac{(x^T A)d(x)}{dx} + z \\ &= Ax + (x^T A)^T + z \\ &= Ax + A^T x + z = \underline{\underline{(A + A^T)x + z}} \end{aligned}$$

Shape check for $f'(x)$: $(n \times n)(n \times 1) + (n \times 1)$
 $= (n \times 1) + (n \times 1)$
 $= (n \times 1) \rightarrow \text{column vector.}$
 $\rightarrow \text{same shape as } \underline{\underline{x}}.$

(2) $L(w) = \frac{1}{2} (w^T x - y)^2$, $w, x \in \mathbb{R}^{n \times 1}$
 $y \in \mathbb{R}^{1 \times 1}$

ans. Shape of $[w^T x - y]$
 $= (1 \times n)(n \times 1) - (1 \times 1)$
 $= (1 \times 1) - (1 \times 1) = (1 \times 1) \rightarrow \text{scalar.}$

\therefore Shape of $\underline{\underline{L(w)}}$ = scalar.

Let $z = w^T x - y$

$$\therefore L(w) = \frac{1}{2} z^2$$

~~$$\frac{dL(w)}{dz}$$~~

$$\begin{aligned} \therefore \frac{dL(w)}{dw} &= \frac{dL(w)}{dz} \cdot \frac{dz}{dw} \\ &= \frac{1}{2} \times \cancel{z} \times \frac{d(w^T x - y)}{dw} \\ &= z \cdot \left[\frac{d(w^T x)}{dw} - \underbrace{\frac{dy}{dw}}_0 \right] \\ &= z \cdot x = \underline{\underline{(w^T x - y) \cdot x}} \end{aligned}$$

Shape check for $f'(w) =$
 $(1 \times 1)(n \times 1)$

→ since z is a scalar, $zx = xz$

$\therefore (n \times 1)(1 \times 1) \rightarrow (n \times 1) \rightarrow$ column vector

(3) $L(w) = \frac{1}{2m} \|xw - y\|^2$, $w, y \in \mathbb{R}^{n \times 1}$
 $X \in \mathbb{R}^{n \times n}$

ans. Shape of $L(w) = (1 \times 1)$ $[\because \text{norm of a vector is taken}]$.

$$L(w) = \frac{1}{2m} (xw - y)^T \cdot (xw - y)$$

Let $u = (xw - y)$
 $\Rightarrow L(w) = \frac{1}{2m} u^T \cdot u$

$$\begin{aligned} \frac{dL(w)}{dw} &= \frac{1}{2m} \left[\frac{du^T}{dw} \cdot u + u^T \cdot \frac{du}{dw} \right] \\ &= \frac{1}{2m} \left[\frac{du^T}{dw} \cdot u + \frac{du^T}{dw} \cdot u \right] = \frac{1}{2m} \left[\cancel{2} \cdot \frac{du^T}{dw} \cdot u \right] \end{aligned}$$

$$= \frac{1}{m} \left[\frac{d(xw - y)}{dw} \right] \cdot u$$

$$= \frac{1}{m} \left[\frac{d(xw)}{dw} - \frac{dy}{dw} \right] \cdot u$$

$$= \frac{1}{m} [x^T] \cdot u$$

$$\Rightarrow \frac{dL(w)}{dw} = \frac{1}{m} [x^T] \cdot [xw - y]$$

Shape check for $L'(w) = (n \times n) [(n \times n)(n \times 1) - (n \times 1)]$
 $= (n \times n) [(n \times 1)]$
 $= (n \times 1) \rightarrow \text{column vector}$

Q3. $z = wx + b$
 $L = \|z - y\|^2$
 $w \in \mathbb{R}^{n \times 1}$
 $x, b, y \in \mathbb{R}^{n \times 1}, \quad dL/dw = ?$

ans. Shape of $L \rightarrow \text{scalar}$ since norm is taken.

$$L = (z - y)^T \cdot (z - y)$$

Let $u = (z - y)$

$$\Rightarrow L = u^T \cdot u$$

$$\therefore \frac{dL}{dw} = \frac{d(u^T \cdot u)}{dw} = \frac{d(u^T)}{dw} \cdot u + u^T \cdot \frac{d(u)}{dw}$$

$$= \frac{d(u^T)}{dw} \cdot u + \frac{d(u)}{dw} \cdot u$$

$$= \frac{2d(u^T) \cdot u}{dw} \rightarrow \textcircled{1}$$

$$\text{Now, } \frac{d(u^T)}{dW} = \frac{d[Wx + b - y]^T}{dW} = \frac{d[x^T W^T]}{dW} + \frac{db^T}{dW} - \frac{dy^T}{dW}$$

→ (2)

We vectorize matrix $W \in \mathbb{R}^{n \times n}$ into a row vector of n features, of length (n) each.

$$\Rightarrow W = [w_1 \ w_2 \ \dots \ w_n]$$

$$\therefore \dim(W) = (n \times n)$$

Using this in eq (2),

$$\frac{d(u^T)}{dW} = x^T + 0 - 0 = x^T$$

Combining (1) & (3),

$$\frac{dL}{dW} = \underline{2(Wx + b - y) \cdot x^T}$$

$$\begin{aligned} \text{Shape check of } \frac{dL}{dW} &= [(n \times n)(n \times 1) + (n \times 1) - (n \times 1)](1 \times n) \\ &= (n \times 1) \cdot (1 \times n) \\ &= \underline{(n \times n)} \end{aligned}$$

Q4. (i) Gradient descent for $\alpha = 0.5, 1.5, 2.5$

Given :-

$$\hat{y} = xw, \quad x \in \mathbb{R}, \quad w \in \mathbb{R}$$

$$L2 \text{ loss} = \frac{1}{2} (xw - y)^2 = J$$

$$w_0 = 0$$

$$[x = 1, y = 100]$$

~~Q4~~ ~~Q4~~

$$\frac{dJ}{dw} = \frac{1}{2} \times 2 (xw - y) \cdot x = \boxed{(xw - y)x = \frac{dJ}{dw}} \rightarrow (1)$$

$$= (w - 100)$$

(i) $\alpha = 0.5$

ans. Iteration 1

$$w_1 = w_0 - \alpha \frac{dJ}{dw} \bigg|_{w_0}$$

$$= 0 - 0.5 (1 \times 0 - 100) (1)$$

$$= 50$$

$$L_1 = \frac{1}{2} (1 \times 50 - 100)^2 = \frac{2500}{2} = 1250$$

Iteration 2

$$w_2 = 50 - 0.5 (50 - 100)$$

$$= 50 + 25 = 75$$

$$L_2 = \frac{1}{2} (75 - 100)^2 = \frac{625}{2} = 312.5$$

Iteration 3

$$w_3 = 75 - 0.5 (75 - 100)$$

$$= 75 + 0.5 \times 25 = 87.5$$

$$L_3 = \frac{1}{2} (87.5 - 100)^2 = 78.125$$

Iteration 4

$$w_4 = 87.5 - 0.5 (87.5 - 100) \\ = 93.75$$

$$L_4 = \frac{1}{2} (93.75 - 100)^2 = 19.53$$

Iteration 5

$$w_5 = 93.75 - 0.5 (93.75 - 100) \\ = 96.875$$

$$L_5 = \frac{1}{2} (96.875 - 100)^2 = 4.88$$

Iteration 6

$$w_6 = 96.875 - 0.5 (96.875 - 100) \\ = 98.4375$$

$$L_6 = \frac{1}{2} (1.5625)^2 = 1.22$$

Iteration 7

$$w_7 = 98.4375 - 0.5 (98.4375 - 100) \\ = 99.21875$$

$$L_7 = \frac{1}{2} (99.21875 - 100)^2 = 0.305$$

\therefore Loss is < 1 , we can stop.

optimal $w = 99.21875$

(ii) $\alpha = 1.5$

Iteration 1

$$w_1 = 0 - 1.5(0 - 100) = 150.0$$

$$J_1 = \frac{1}{2}(150 - 100)^2 = 1250$$

Iteration 2

$$w_2 = 150 - 1.5(150 - 100) = 75$$

$$J_2 = \frac{(75 - 100)^2}{2} = 312.5$$

Iteration 3

$$w_3 = 75 - 1.5(75 - 100) = 112.5$$

$$J_3 = \frac{(112.5 - 100)^2}{2} = 78.125$$

Iteration 4

$$w_4 = 112.5 - 1.5(112.5 - 100) = 93.75$$

$$J_4 = 19.53$$

Iteration 5

$$w_5 = 93.75 - 1.5(93.75 - 100) = 103.125$$

$$J_5 = 4.88$$

Iteration 6

$$w_6 = 103.125 - 1.5(103.125 - 100) = 98.43$$

$$J_6 = 1.22$$

★★ w is oscillating around 100. Stopping for $w_6 = 98.43$.

(ii) $\alpha = 2.5$

Iteration 1

$$w_1 = 0 - 2.5(0 - 100) \\ = 250$$

$$J_1 = \frac{(250 - 100)^2}{2} = 11250$$

Iteration 2

$$w_2 = 250 - 2.5(250 - 100) \\ = -125$$

$$J_2 = \frac{(-125 - 100)^2}{2} = 25312.5$$

Iteration 3

$$w_3 = -125 - 2.5(-125 - 100) \\ = 437.5$$

$$J_3 = 56953.125$$

\therefore Since learning rate is too high, Gradient descent is showing divergent behaviour.
 \Rightarrow It will take a greater no. of iterations to converge.

(2) For optimal position, loss should be 0.
 $\Rightarrow J = \frac{1}{2} (wx - y)^2 = 0$ [can also be derived by finding minima]

$$\Rightarrow |wx - y| = 0$$

Since there is only 1 data point,
 $|w(1) - 100| = 0$

$$\Rightarrow |w - 100| = 0 \rightarrow (1)$$

for optimal w .

(i) Condition for α when G.D starts oscillating around optimal position.

ans. Given:

$$\text{OPTIMAL POINT} = 100$$

* Oscillation around optimal point means that simultaneous w_s are on either side of 100.

$$\Rightarrow \text{Case 1: } \begin{cases} w_n > 100 \\ w_{n-1} < 100 \end{cases}$$

We know,

$$w_n = w_{n-1} - \alpha(w_{n-1} - 100) > 100$$

$$\Rightarrow (w_{n-1} - 100) > \alpha(w_{n-1} - 100) \rightarrow (1)$$

$$\text{Given that } (w_{n-1} - 100 < 0)$$

$$\frac{(w_{n-1} - 100)}{(w_{n-1} - 100)} < \alpha$$

$$\Rightarrow \boxed{\alpha > 1} \rightarrow (2)$$

$$\text{Case 2: } \begin{cases} w_n < 100 \\ w_{n-1} > 100 \end{cases}$$

lly,

$$w_n = w_{n-1} - \alpha(w_{n-1} - 100) < 100$$

$$\Rightarrow (w_{n-1} - 100) < \alpha(w_{n-1} - 100) \rightarrow (3)$$

$$\text{Given that } (w_{n-1} - 100 > 0)$$

$$\frac{(w_{n-1} - 100)}{(w_{n-1} - 100)} < \alpha$$

$$\Rightarrow \boxed{\alpha > 1} \rightarrow (4)$$

Combining (2) & (4),

$$\boxed{\alpha > 1} \text{ for oscillation}$$

(ii) Condition for α that GD can converge.

ans. For convergence, assume stopping criteria to be s.t $|w_t - 100| < 1$

$$\Rightarrow -1 < (w_t - 100) < 1$$

$$\Rightarrow \boxed{99 < w_t < 101} \rightarrow (1)$$

We know,

$$w_t = w_{t-1} - \alpha(w_{t-1} - 100) \rightarrow (2)$$

$$\Rightarrow 99 < w_{t-1} - \alpha(w_{t-1} - 100) < 101$$

$$\Rightarrow (w_{t-1} - 99) > \alpha(w_{t-1} - 100) > (w_{t-1} - 101) \rightarrow (3)$$

~~w_{t-1} is constant~~

Since (w_{t-1}) is the weight of a previous iteration

$$(w_{t-1}) \leq 99 \quad \text{or} \quad (w_{t-1}) \geq 101$$

★ Case 1: $(w_{t-1}) \leq 99$

\Rightarrow From (3),

$$0 \leq \frac{(w_{t-1} - 99)}{(w_{t-1} - 100)} < \alpha < \frac{(w_{t-1} - 101)}{(w_{t-1} - 100)} \rightarrow \text{scribbles}$$

$$\Rightarrow \boxed{0 < \alpha < 2} \rightarrow (4)$$

★ Case 2: $(w_{t-1}) \geq 101$

$$2 \geq \left(\frac{w_{t-1} - 99}{w_{t-1} - 100} \right) > \alpha > \left(\frac{w_{t-1} - 101}{w_{t-1} - 100} \right) \geq 0$$

$$\Rightarrow \boxed{2 > \alpha > 0} \rightarrow (5)$$

∴ Combining case ① & ②,

$0 < \alpha < 2$ → for convergence.

③. (i) For $\alpha = 0.5$, it converges in iteration 7, for $w = 99.21875$

(ii) For $\alpha = 1.5$, it oscillates around the optimal point $[100]$ starting from iteration ①

(iii) For $\alpha = 2.5$, loss is monotonically increasing & never converges.