1) a) $g(x) := \text{GELU}(x) = x\,\sigma(1.702 \cdot x)$

$$\frac{dg(x)}{dx} = \sigma(1.702x) + x\left(\sigma(1.702x)(1-\sigma(1.702x))\right) \times 1.702$$

For G.D.

$$x_{t+1} = x_t - \eta \nabla g(x_t)$$

$\underline{t=0}$

$\quad x_0 = 0, \qquad g(0) = 0$

$\underline{t=1}$

$\Rightarrow\; x_1 = x_0 - g(0.1)\,\nabla g(0)$

$$= 0 - 0.1\left[\sigma(0) + 0\right] = -0.1\left(\frac{1}{2}\right)$$

$$\Rightarrow \boxed{x_1 = -0.05}$$

$\underline{t=2}$

$x_2 = x_1 - (0.1)\,\nabla g(-0.05)$

$$= -0.05 - (0.1)\left[\sigma(-0.05) - 0.05\left(\sigma(0.05)\right)(1-\sigma(\right.$$

$$= -0.05 - 0.1\left[\sigma(-0.085) - 0.05\left(\sigma(-0.085)\right) \cdot (1-\sigma(-0.085))\times 1.702\right]$$

$$= -0.05 - 0.1\left[0.4787 - 0.05\left(0.4787\,(1-0.4787)\right)\times 1.702\right]$$

$$\boxed{x_2 = -0.0957}$$

$t=3$

$$x_3 = x_2 - (0.1) \nabla g(-0.0957)$$

$$\Rightarrow x_3 = -0.0957 - 0.1\left[\sigma(-0.1629)(1 - \sigma(\right.$$

$$x_3 = -0.0957 - 0.1\left[\sigma(-0.1629) - 0.0957\left(\sigma(0.1629)(1-\sigma(0.1629))\right)\times 1.702\right]$$

$$\boxed{x_3 = -0.1376}$$

function values,

$$\boxed{\begin{array}{l} g(x_1) = -0.0239 \\ g(x_2) = -0.0439 \\ g(x_3) = -0.0608 \end{array}}$$

b)   with $\eta = 1.0$ , $x_0 = 0$

$t=1$

$$x_1 = 0 - (1.0)\nabla g(0) = -0.5$$

$$\Rightarrow \boxed{x_1 = -0.5} \quad \text{and} \quad \boxed{g(x_1) = -0.1496}$$

$t=2$

$$x_2 = -0.5 - (1.0)\left[\sigma\left(\overset{-0.85}{-0.5}\right) - 0.5\left(\sigma(-0.851)(1-\sigma(-0.851))\right)\times 1.702\right]$$

$$\boxed{x_2 = -0.6208} \quad \Rightarrow \boxed{g(x_2) = -0.1601}$$

$t = 3$

$$x_3 = (-0.6208) - (1.0)\left[\sigma(-1.057) - (0.6208)(\sigma(-1.057)(1-\sigma(-1.057))) \atop \times 1.702\right]$$

$$\boxed{x_3 = -0.6765} \quad \Rightarrow \boxed{g(x_3) = -0.1625}$$

with $\eta = 1.0$, the function value decreased by almost an order of magnitude faster than with $\eta = 0.1$

c) i)  $x_0 = -3$, $\eta = 0.1$

$t = 1$

$$x_1 = (-3) - (0.1)\,\nabla g(0.1)$$

$$\boxed{x_1 = -2.9975} \quad \Rightarrow \boxed{g(x_1) = -0.0181}$$

$t = 2$

$$x_2 = (-2.9975) - (0.1)\, \nabla g(-2.9975)$$

$$\Rightarrow \boxed{x_2 = -2.9951} \quad \Rightarrow \boxed{g(x_2) = -0.0182}$$

$t = 3$

$$x_3 = (-2.9951) - (0.1)\,\nabla g(-2.9951)$$

$$\Rightarrow \boxed{x_3 = -2.9926} \quad \Rightarrow \boxed{g(x_3) = -0.0183}$$

ii) GD with momentum, $\beta = 0.9$, $\eta = 0.1$, $x_0 = -3$

$$v_0 = \nabla g(x_0) = -0.0245848$$

t = 1

$$v_1 = (0.9) v_0 + (0.1) \nabla g(x_0) =$$

$$\Rightarrow \boxed{v_1 = -0.024548}$$

$$x_1 = x_0 - (0.1) v_1$$

$$\Rightarrow \boxed{x_1 = -2.99755} \quad , \quad \boxed{g(x_1) = -0.0181}$$

t = 2

$$v_2 = (0.9) v_1 + (0.1) \nabla g(x_1)$$

$$v_2 = -0.24556$$

$$x_2 = x_1 - (0.1) v_2$$

$$\Rightarrow \boxed{x_2 = -2.995089} \quad , \quad \boxed{g(x_2) = -0.018192}$$

t = 3

$$v_3 = (0.9) v_2 + (0.1) \nabla g(x_2)$$

$$\boxed{v_3 = -0.02457}$$

$$x_3 = x_2 - (0.1) v_3$$

$$\Rightarrow \boxed{x_3 = -2.992632} \Rightarrow \boxed{g(x_3) = -0.018253}$$

iii) In this case both the methods perform almost the same as which could be seen by in the $x_t$ and $g(x_t)$ values.
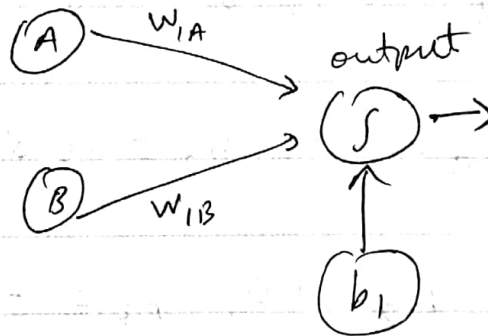
2)

AND gate



{ $S$ : signoid }

$$W_{1A} = \cancel{12} \; 1.0 \quad, \quad W_{1B} = 1.0 \quad, \quad b_1 = -1.5$$

OR Gate



{ $S$ : sigmoid }

$$W_{1A} = 1.0 \quad, \quad W_{1B} = 1.0, \quad b_1 = -0.5$$

XOR gate



{ $S$ : sigmoid }

$$W_{1A} = -0.4 \quad, \quad W_{2A}$$

# XOR gate



$$W_{1A} = -0.4, \quad W_{1B} = -0.4, \quad b_1 = 0.6$$

$$W_{2A} = 1.0, \quad W_{2B} = 1.0, \quad b_2 = -0.5$$

$$W_{31} = 1.0, \quad W_{32} = 1.0, \quad b_3 = -1.5$$

# XNOR gate



$$W_{1A} = 1.0, \quad W_{1B} = 1.0, \quad b_1 = -1.5$$
$$W_{2A} = -0.7, \quad W_{2B} = -0.7, \quad b_2 = 0.6$$
$$W_{31} = 1.0, \quad W_{32} = 1.0, \quad b_3 = -0.5$$

3) a)

$$E(\theta) = -\sum_{i=1}^{k} y_i \log o_i$$

$$\nabla_{f_2} E = \left(\nabla_{f_2} o\right)\left(\nabla_\theta E\right)$$

Now, $\nabla_\theta E = \begin{bmatrix} \partial E/\partial o_1 \\ \vdots \\ \partial E/\partial o_i \\ \vdots \\ \partial E/\partial o_k \end{bmatrix} = \begin{bmatrix} -y_1/o_1 \\ \vdots \\ \vdots \\ -y_k/o_k \end{bmatrix}$

$$\Rightarrow \nabla_\theta E_i = -y_i/o_i$$

and $\nabla_{f_2} o = \begin{bmatrix} \dfrac{\partial S(f_{21})}{\partial f_{21}} & \cdots & \dfrac{\partial S(f_{2k})}{\partial f_{21}} \\ \vdots & & \\ \dfrac{\partial S(f_{21})}{\partial f_{2k}} & - - & \dfrac{\partial S(f_{2k})}{\partial S(f_{2k})} \end{bmatrix}$

$$\nabla_{f_{2j}} o_i = \frac{\partial S(f_{2i})}{\partial f_{2j}} = \frac{e^{f_{2i}}\left(\frac{\partial z_i}{\partial x_j}\right)e^{\dagger}}{\left(\sum_{p=1}^{k} e^{x_p}\right)^2}$$

$$= \frac{e^{f_{2i}}}{\sum_{p=1}^{k} e^{f_{2p}}}\left(\frac{\partial f_{2i}}{\partial f_{2j}}\right) - \frac{e^{f_{2i}}\, e^{f_{2j}}}{\left(\sum_{p=1}^{k} e^{f_{2p}}\right)^2}$$

if $i = j$, $\Rightarrow \dfrac{\partial f_{2i}}{\partial f_{2j}} = 1$ $\Rightarrow \nabla_{f_{2i}} o_i = \dfrac{e^{f_{2i}}}{\sum_{p=1}^{k} e^{f_{2p}}}\left[1 - \left(\dfrac{e^{f_{2i}}}{\sum_{p=1}^{k} e^{f_{2p}}}\right)\right]$

$$\Rightarrow \frac{\partial f_{2i}}{\partial f_{2i}} =$$

$$\Rightarrow \quad \nabla_{f_{2i}} O_i = O_i(1-O_i)$$

if $i \neq j$ , $\Rightarrow \dfrac{\partial f_{2i}}{\partial f_{2j}} = 0 \Rightarrow \nabla_{f_{2j}} O_i = -O_i O_j$

$$\Rightarrow \quad \nabla_{f_2} O = \begin{bmatrix} O_1(1-O_1) & \cdots & -O_1 O_k \\ \vdots & & \\ -O_1 O_k & \cdots & O_k(1-O_k) \end{bmatrix}$$

or $\quad \nabla_{f_{2j}} O_i = \left\{ \begin{array}{l} O_i(1-O_i) \;\; ; \; \text{if } i=j \\[2mm] -O_i O_j \;\;\; ; \; \text{if } i \neq j \end{array} \right\}$

$$\Rightarrow \quad \nabla_{f_2} \oslash E = \begin{bmatrix} O_1(1-O_1) & \cdots & O_1 O_k \\ \vdots & & \vdots \\ -O_1 O_k & & O_k(1-O_k) \end{bmatrix} \begin{bmatrix} -y_1/O_1 \\ \vdots \\ -y_k/O_k \end{bmatrix}$$

In classification, Let $y_q = 1$ and $y_{i \neq q} = 0$

$$\Rrightarrow \quad \nabla_{f_2} E =$$

$$\Rightarrow \quad \nabla_{f_{2j}} E_{\ell} = \begin{bmatrix} \sum_{\substack{P=1 \\ P \neq j}}^{K} O_p(1-O_p)\left(-\dfrac{y_p}{O_p}\right) \end{bmatrix} \neq O_j \oslash$$

$$\Rightarrow \nabla_{f_{2i}} E_{\frac{\xi}{2}} = \left[ \sum_{\substack{j=1, \\ j \neq i}}^{K} (-O_i O_j)\left(-\frac{y_j}{O_j}\right) \right] + O_i(1-O_i)\left(-\frac{y_i}{O_i}\right)$$

$$\Rightarrow \nabla_{f_{2i}} E = \sum_{\substack{j=1, \\ j \neq i}}^{K} y_j O_i + y_i O_i - y_i$$

$$\Rightarrow \boxed{\nabla_{f_{2i}} E = \left[ \sum_{j=1}^{K} y_j O_i \right] - y_i}$$

b) $\quad \nabla_x E = \left(\nabla_x f_1\right)\left(\nabla_{f_1} a\right)\left(\nabla_a f_2\right)\left(\nabla_{f_2} E\right)$

Now $\rightarrow \nabla_x f_1 = \nabla_x (x W_1 + b_1) = W_1 \qquad -①$

$\rightarrow \nabla_{f_1} a = \nabla_{f_1} \sigma(f_1) = \cancel{\nabla(f_1)\left(\frac{1}{1} \sigma(f_1)\right)} \sigma(f_1) I \left(I - \frac{\sigma(f_1)}{I}\right)$

$\Rightarrow \nabla_{f_1} a = a I\left(\frac{I}{1} - a I\right) \qquad -②\qquad \left\{ \overline{I = [1 \cdots 1]}^{T} \right\}_{M \text{ times}}$

$\rightarrow \nabla_a f_2 = \nabla_a (a W_2 + b_2)$

$\Rightarrow \nabla_a f_2 = W_2$

$$\Rightarrow \boxed{\nabla_x E = W_1 \, a I(I_M - a I) \, W_2 \, \nabla_{f_2} E}$$

4) $\cancel{a}$  Given the feature map has same Height and width $(H_{in})$

$$H_{out} = \frac{H_{in} + \cancel{2xps} \; 2p - d(k_{\cancel{x}} - 1) - 1}{s} + 1 \quad - \text{①}$$

a)  given $s = 1$, $d = $ dialation $= 1$, and $H_{out} = H_{in}$,
and $k = 3$

$$\cancel{s}(H_{in}) = \frac{H_{in} + 2p - (3-1) - 1 + 1}{1}$$

$$2p = \cancel{x} \; 2 \quad \Rightarrow \quad \boxed{P = 1}$$

$$F_{out} = \begin{bmatrix} 13.5 & -18.5 & 9 & -12.0 \\ 4 & 21.5 & 5 & 11.5 \\ 20.5 & 10.5 & 24.5 & 17 \\ 6 & 24 & 17 & 13.5 \end{bmatrix}$$

b)  for filter 2, $k = 2$
Using ① with $H_{out} = H_{in} = 4$, $p = 1$, $d = 1$

$$\overset{4}{4 \cancel{H_{out}}} = \frac{\overset{4}{\cancel{H_{in}}} + 2 - 1(2-1) - 1}{s} + 1$$

$$3 = \frac{4}{s} \quad \Rightarrow \quad \boxed{s = \frac{4}{3}} \quad \Rightarrow \quad s \text{ is not an integer}$$
so it is not possible to get a feature map with $H_{out} = H_{in}$

c) $\quad F' = \begin{bmatrix} 13.5 & -18.5 & 9 & -12 \\ 4 & 21.5 & 5 & 11.5 \\ 20.5 & 10.5 & 24.5 & 17 \\ 6 & 24 & 17.0 & 13.5 \end{bmatrix}$

$\text{Avgpool}(F') = \begin{bmatrix} 5.125 & 3.375 \\ 15.25 & 18.0 \end{bmatrix}$

d) $\quad$ convolution filter for average pooling,

$\qquad \text{Filter} = \begin{bmatrix} 0.25 & 0.25 \\ 0.25 & 0.25 \end{bmatrix}$

$\qquad \text{padding} = 0 \quad, \quad \text{stride} = 2$