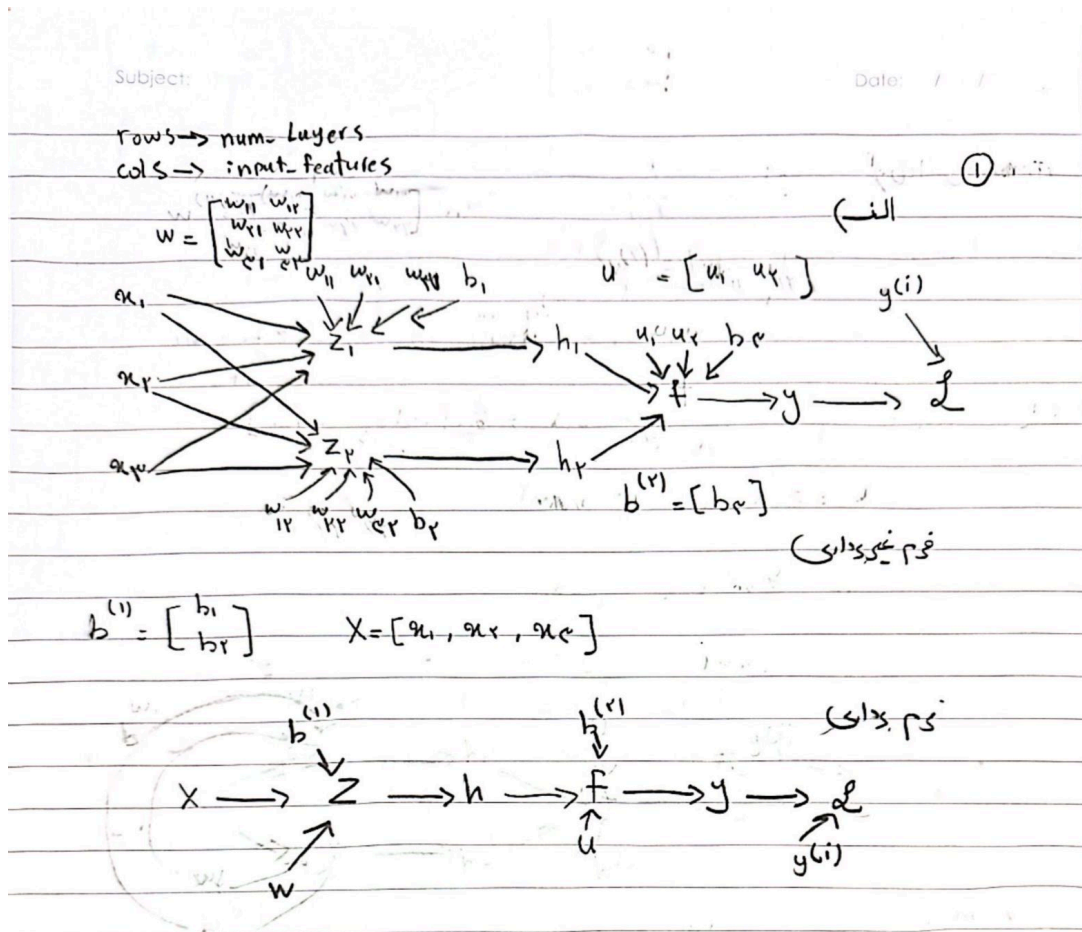


## تمرین دوم - بخش تئوری

(matrix calculation notebook)

سوال اول



\*  $\text{sigmoid}(x) = \frac{e^x}{1+e^x}$  (.)

- Feed Forward formulas

first layer:  $z_1 = \begin{bmatrix} w_{11} \\ w_{12} \\ w_{13} \end{bmatrix} [x_1, x_2, x_3] + b_1 \quad h_1 = \frac{z_1}{1+e^{-z_1}}$

second layer:  $z_2 = \begin{bmatrix} w_{21} \\ w_{22} \\ w_{23} \end{bmatrix} [h_1, h_2, h_3] + b_2 \quad h_2 = \frac{z_2}{1+e^{-z_2}}$

$f = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} [h_1, h_2] + b_3 \quad y = \frac{1}{1+e^{-f}} \quad L = y(i) \log(1+e^{-f}) + (1-y(i)) \log(1+e^f)$

Feed forward formulas (matrix form):

(-)

first layer:  $Z = X W = [z_1 \quad z_r]$

$$h = z \circ \sigma(z) = [z_1 \sigma(z_1) \quad z_r \sigma(z_r)]$$

second layer:

$$\begin{bmatrix} u_1 & u_r \end{bmatrix} \begin{bmatrix} h_1 \\ h_r \end{bmatrix} \quad f = u \cdot h \quad y = \frac{1}{1 + e^{-f}} \quad \mathcal{L} = -y^{(i)} \log(y) - (1 - y^{(i)}) \log(1 - y)$$

Backward pass (matrix form)

$$\bar{\mathcal{L}} = 1 \quad * \quad \bar{W} = \bar{z} \eta^T$$

$$\bar{y} = \bar{\mathcal{L}} \left[ -\frac{y^{(i)}}{y} + \frac{1 - y^{(i)}}{1 - y} \right] \quad * \quad \bar{b}^{(1)} = \bar{z}$$

$$\bar{f} = \bar{y} \frac{dy}{df} = \bar{y} \circ \sigma'(f) \quad * \quad \bar{u} = \bar{f} h^T$$

$$\bar{h} = \bar{f} \frac{df}{dh} = \bar{f} u^T \quad * \quad \bar{b}^{(2)} = \bar{f}$$

$$\bar{z} = \bar{h} \frac{dh}{dz} = \bar{h} \circ (z \circ \sigma(z))' = \bar{h} \circ (\sigma(z) + z \sigma'(z))$$

iteration 1

$$\begin{bmatrix} 0.8 & -0.3 \\ 0.2 & 0.2 \\ -0.5 & 0.1 \end{bmatrix} \leftarrow \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{c1} & w_{c2} \end{bmatrix} + \begin{bmatrix} b_1 & b_2 \end{bmatrix} \rightarrow b^{(1)}$$

forward pass:

$$\text{First layer: } z_1 = x^{(1)} W + b = \begin{bmatrix} 0.15 & -0.45 \end{bmatrix} \rightarrow z_r$$

$$x^{(1)} = [1 \ 0.2]$$

$$z_r = x^{(1)} W + b = \begin{bmatrix} -0.24 & -0.22 \end{bmatrix}$$

$$x^{(2)} = [0.1 \ 0.8 \ 1]$$

$$x^{(c)} = [0.1 \ 0.4]$$

$$z_c = x^{(c)} W + b = \begin{bmatrix} -0.15 & -0.12 \end{bmatrix}$$

$$\begin{aligned} h_1 &= z_1 \sigma(z_1) = \begin{bmatrix} 0.08 & -0.22 \end{bmatrix} \\ h_r &= z_r \sigma(z_r) = \begin{bmatrix} -0.15 & -0.10 \end{bmatrix} \\ h_c &= z_c \sigma(z_c) = \begin{bmatrix} -0.04 & -0.04 \end{bmatrix} \end{aligned} \left\{ \begin{bmatrix} \sigma(z_1) & \sigma(z_r) \end{bmatrix} \right.$$

second layer:

$$\begin{aligned} f_1 &= u h_1 + b_c = 0.15 \\ f_r &= u h_r + b_c = 0.06 \\ f_c &= u h_c + b_c = 0.09 \end{aligned} \left\{ \begin{aligned} f &= [u_1 \ u_r] \begin{bmatrix} h_1 \\ h_r \end{bmatrix} + b_r \\ &\rightarrow \begin{bmatrix} 0.05 & -0.1 \end{bmatrix} \end{aligned} \right.$$

$$\begin{aligned} y_1 &= \frac{1}{1+e^{-f_1}} = 0.537 \\ y_r &= \frac{1}{1+e^{-f_r}} = 0.517 \\ y_c &= \frac{1}{1+e^{-f_c}} = 0.521 \end{aligned} \left\{ \rightarrow y = \text{sigmoid}(f) \right.$$

$$L_{\text{-iter1}} = \frac{L_1 + L_r + L_c}{3} \left\{ \begin{aligned} L_1 &= 0.422 \\ L_r &= 0.422 \\ L_c &= 0.422 \end{aligned} \right\} \rightarrow -y^{(1)} \log(y) - (1-y) \log(1-y)$$

$$L_{\text{-iter1}} = 0.695$$



Backward pass:

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial f} \frac{\partial f}{\partial h} \frac{\partial h}{\partial z} \frac{\partial z}{\partial w} \quad \frac{dh}{dz}$$

$$\frac{\partial \mathcal{L}}{\partial w} = \left[ \frac{-y^{(i)}}{y} + \frac{1-y^{(i)}}{1-y} \right] \cdot \sigma'(f) U_0'(z \sigma(z))' X^T$$

$$= \left[ \frac{-y^{(i)}}{y} + \frac{1-y^{(i)}}{1-y} \right] \cdot (y(1-y)) U_0'(\sigma(z) + z \sigma(z)(1-\sigma(z)))$$

$$= \left[ \frac{-y^{(i)}}{y} + \frac{1-y^{(i)}}{1-y} \right] \cdot (y(1-y)) \cdot U_{1 \times 2} \left[ \begin{array}{c} \sigma(z_1) + z_1 \sigma(z_1)(1-\sigma(z_1)) \\ \sigma(z_2) + z_2 \sigma(z_2)(1-\sigma(z_2)) \end{array} \right]^T \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix}^T$$

scalar ← element wise

طبقی تعریف انجام شده

$$u \rightarrow 1 \times 2 \quad \frac{dh}{dz} \rightarrow 1 \times 2 \quad \rightarrow 0 \rightarrow 1 \times 2 \rightarrow [1 \times 2]^T \rightarrow 2 \times 1 \quad 1 \times 2 \rightarrow 2 \times 2$$

element wise

$$\frac{\partial \mathcal{L}}{\partial w} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial w_{11}} & \frac{\partial \mathcal{L}}{\partial w_{12}} & \frac{\partial \mathcal{L}}{\partial w_{13}} \\ \frac{\partial \mathcal{L}}{\partial w_{21}} & \frac{\partial \mathcal{L}}{\partial w_{22}} & \frac{\partial \mathcal{L}}{\partial w_{23}} \end{bmatrix}$$

for  $x^{(1)} = \begin{bmatrix} -0.008 & 0 & -0.005 \\ 0.009 & 0 & 0.005 \end{bmatrix}$   $a^{(1)} = [1 \ 0 \ 0.8] \quad y^{(1)} = 1$

for  $x^{(2)} = \begin{bmatrix} 0.005 & 0.002 & 0.005 \\ 0.005 & -0.001 & -0.002 \end{bmatrix}$   $a^{(2)} = [0.1 \ 0.8 \ 1] \quad y^{(2)} = 0$

for  $x^{(3)} = \begin{bmatrix} 0 & 0.002 & 0.005 \\ 0 & -0.002 & -0.002 \end{bmatrix}$   $a^{(3)} = [0 \ 1 \ 0.6] \quad y^{(3)} = 0$

average:

$$n=3 \rightarrow \frac{1}{n} \sum_{i=1}^n \frac{\partial \mathcal{L}_i}{\partial w} = \begin{bmatrix} -0.004 & 0.003 & 0.01 \\ 0.002 & -0.001 & -0.001 \end{bmatrix}$$

Gradients of  $b^{(1)}$  - first layer

$$\frac{\partial \mathcal{L}}{\partial b^{(1)}} = \frac{\partial \mathcal{L}}{\partial y} \cdot \frac{\partial y}{\partial f} \cdot \frac{\partial f}{\partial h} \cdot \frac{\partial h}{\partial z}$$

$$= \left[ \left( \frac{-y^{(1)}}{y} + \frac{1-y^{(1)}}{y} \right) \cdot y(1-y) \right] \cdot \underbrace{1 \times 1}_{\text{element-wise}} \cdot \underbrace{1 \times 2}_{\text{1x2}}$$

scalar ←
↓ 1x1
→ 1x2

$$= \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial b_1} & \frac{\partial \mathcal{L}}{\partial b_2} \end{bmatrix} = \left[ \left( \frac{-y^{(1)}}{y} + \frac{1-y^{(1)}}{y} \right) \cdot y(1-y) \right] \cdot \begin{bmatrix} \delta(z) + z\delta(z_1)(1-\delta(z_1)) \\ \delta(z_1) + z_1\delta(z_1)(1-\delta(z_1)) \end{bmatrix}^T$$

For  $X^{(1)} = \begin{bmatrix} -0.01 & 0.009 \end{bmatrix}$   $y^{(1)} = 1$

For  $X^{(2)} = \begin{bmatrix} 0.01 & -0.02 \end{bmatrix}$   $y^{(2)} = 0$

For  $X^{(3)} = \begin{bmatrix} 0.02 & -0.02 \end{bmatrix}$   $y^{(3)} = 0$

average

$$n=3 \rightarrow \frac{1}{n} \sum_{i=1}^n \frac{\partial \mathcal{L}_i}{\partial b^{(1)}} = \begin{bmatrix} 0.012 & -0.01 \end{bmatrix}$$

Gradients of  $U$ 

$$\frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial f} h^T$$

$$= \left[ \frac{-y^{(i)}}{y} + \frac{1-y^{(i)}}{y} \right] \cdot y(1-y) h_i$$

1x1

scalar

$$\frac{\partial \mathcal{L}}{\partial u} = \left[ \frac{\partial \mathcal{L}}{\partial u_1} \quad \frac{\partial \mathcal{L}}{\partial u_r} \right] = \left[ \left( \frac{-y^{(i)}}{y} + \frac{1-y^{(i)}}{y} \right) y(1-y) \right] [h_1 \quad h_r]$$

$$\text{for } x^{(1)} = [-0.98 \quad 0.10]$$

$$\text{for } x^{(2)} = [-0.25 \quad -0.05]$$

$$\text{for } x^{(3)} = [-0.10 \quad -0.35]$$

$$n=3 \rightarrow \frac{1}{n} \sum_{i=1}^n \frac{\partial \mathcal{L}}{\partial u_i} = [-0.05 \quad 0.004]$$

Gradients of  $b^{(2)}$  - second layer

$$\frac{\partial \mathcal{L}}{\partial b^{(2)}} = \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial f}$$

$$= \left[ \frac{-y^{(i)}}{y} + \frac{1-y^{(i)}}{y} \right] \cdot y(1-y) \rightarrow \text{scalar}$$

$$\frac{\partial \mathcal{L}}{\partial b^{(2)}} = \left[ \frac{\partial \mathcal{L}}{\partial b_e} \right]$$

$$\text{for } x^{(1)} = -0.98$$

$$\text{for } x^{(2)} = 0.25$$

$$\text{for } x^{(3)} = 0.25$$

$$n=3 \quad \frac{1}{n} \sum_{i=1}^n \frac{\partial \mathcal{L}_i}{\partial b^{(2)}} = [0.19]$$





First iteration weight updates

GID with momentum  $\rightarrow$  update velocity  $v_0 = \vec{0}$   
 $\rightarrow$  update weights  $\delta = 0.9$   
 $\eta = 0.1$

$$W \Rightarrow \begin{cases} v_i \leftarrow \delta v_0 - \frac{\eta}{n} \left[ \sum_{i=1}^n \frac{\partial \mathcal{L}_i}{\partial W} \right] \\ W_{\text{new}} \leftarrow W_{\text{old}} + v_i \end{cases}$$

$$\Rightarrow \begin{cases} v_i \leftarrow \vec{0} - 0.1 \begin{bmatrix} -0.015 & 0.02 & 0.019 \\ 0.002 & -0.01 & -0.01 \end{bmatrix} \\ W_{\text{new}} \leftarrow \begin{bmatrix} 0.1 & -0.1 & 0.1 \\ -0.1 & 0.1 & 0.1 \end{bmatrix} + \begin{bmatrix} 0.002 & -0.002 & -0.001 \\ -0.001 & 0.001 & 0.001 \end{bmatrix} \\ = \begin{bmatrix} 0.102 & -0.102 & 0.099 \\ -0.101 & 0.101 & 0.101 \end{bmatrix} \end{cases}$$

$$b^{(1)} \Rightarrow \begin{cases} v_i \leftarrow \delta v_0 - \frac{\eta}{n} \left[ \sum_{i=1}^n \frac{\partial \mathcal{L}_i}{\partial b^{(1)}} \right] \\ b_{\text{new}} \leftarrow b_{\text{old}} + v_i \end{cases}$$

$$\Rightarrow \begin{cases} v_i \leftarrow \vec{0} - 0.1 \begin{bmatrix} 0.012 & -0.01 \end{bmatrix} \\ b_{\text{new}} \leftarrow \begin{bmatrix} 0 & -0.1 \end{bmatrix} + \begin{bmatrix} -0.01 & 0.009 \end{bmatrix} \\ = \begin{bmatrix} -0.01 & -0.091 \end{bmatrix} \end{cases}$$

$$u \Rightarrow \begin{cases} v_i \leftarrow \delta - \frac{\eta}{n} \left[ \sum_{i=1}^n \frac{\partial \mathcal{L}_i}{\partial u} \right] \\ u_{\text{new}} \leftarrow u_{\text{old}} + v_i \end{cases}$$

$$\Rightarrow \begin{cases} v_i \leftarrow \vec{0} - 0.1 \begin{bmatrix} -0.005 & 0.004 \end{bmatrix} \\ u_{\text{new}} \leftarrow \begin{bmatrix} 0.1 & -0.1 \end{bmatrix} + \begin{bmatrix} 0.005 & -0.004 \end{bmatrix} \end{cases}$$

$$= \begin{bmatrix} 0,14 & -0,10 \end{bmatrix}$$

$$b^{(r)} \Rightarrow \begin{cases} v_i \leftarrow \delta v_i - \frac{\eta}{n} \left[ \sum_{j=1}^n \frac{\partial L_j}{\partial b^{(r)}} \right] \\ b_{new}^{(r)} \leftarrow b_{old}^{(r)} + v_i \end{cases}$$

$$\Rightarrow \begin{cases} v_i \leftarrow 0 - 0,14 \begin{bmatrix} 0,19 \\ -0,11 \end{bmatrix} \\ b_{new}^{(r)} \leftarrow [-1] + [-0,10] \end{cases}$$

$$= \begin{bmatrix} -0,06 \\ -0,10 \end{bmatrix}$$

weights after first iteration

$$W_1 = \begin{bmatrix} 0,81 & 0,14 & -0,81 \\ -0,90 & 0,20 & 0,10 \end{bmatrix} \quad U = \begin{bmatrix} 0,85 & -0,10 \end{bmatrix}$$

$$b^{(1)} = \begin{bmatrix} -0,01 & -0,09 \end{bmatrix} \quad b^{(r)} = \begin{bmatrix} -0,06 \\ -0,10 \end{bmatrix}$$

iteration 2

forward pass: first-layer:  $z_1 = X^{(1)} W^T + b^{(1)} = \begin{bmatrix} 0,18 & -0,14 \end{bmatrix}$   
 $z_2 = X^{(r)} W^T + b^{(1)} = \begin{bmatrix} -0,90 & -0,20 \end{bmatrix}$   
 $z_3 = X^{(e)} W^T + b^{(1)} = \begin{bmatrix} -0,19 & -0,11 \end{bmatrix}$   
 $h_1 = z_1 \sigma(z_1) = \begin{bmatrix} 0,01 & -0,12 \end{bmatrix}$   
 $h_2 = z_2 \sigma(z_2) = \begin{bmatrix} -0,12 & -0,09 \end{bmatrix}$   
 $h_3 = z_3 \sigma(z_3) = \begin{bmatrix} -0,09 & -0,06 \end{bmatrix}$

second layer:

$$\begin{aligned} f_1 &= U h_1 = -0,008 \\ f_2 &= U h_2 = -0,094 \\ f_3 &= U h_3 = -0,017 \end{aligned} \quad \xrightarrow{\text{sigmoid}} \quad \begin{aligned} y_1 &= 0,999 \\ y_2 &= 0,904 \\ y_3 &= 0,981 \end{aligned}$$



$$\Rightarrow \begin{aligned} L_1 &= 0.494 \\ L_2 &= 0.442 \\ L_c &= 0.455 \end{aligned} \rightarrow \boxed{L_{iter 2} = 0.445}$$

Backward pass:

در این جا باید فرمول ها را باز نویسی کنیم  
iteration اول برای محاسبه گرادیان ها نوشته شده است

Gradients of  $w$ :

$$\text{For } X^{(1)} \rightarrow \begin{bmatrix} -0.11 & 0 & -0.05 \\ 0.01 & 0 & 0.05 \end{bmatrix}$$

$$\text{For } X^{(2)} \rightarrow \begin{bmatrix} 0.005 & 0.02 & 0.05 \\ -0.002 & -0.01 & -0.02 \end{bmatrix} \xrightarrow{\text{avg}} \begin{bmatrix} -0.03 & 0.03 & 0.02 \\ 0.002 & -0.01 & -0.01 \end{bmatrix}$$

$$\text{For } X^{(c)} \rightarrow \begin{bmatrix} 0 & 0.42 & 0.05 \\ 0 & -0.42 & -0.02 \end{bmatrix}$$

Gradients of  $b^{(1)}$  - First layer

$$\text{For } X^{(1)} \rightarrow \begin{bmatrix} -0.16 & 0.15 \end{bmatrix}$$

$$\text{For } X^{(2)} \rightarrow \begin{bmatrix} 0.05 & -0.02 \end{bmatrix} \xrightarrow{\text{avg}} \begin{bmatrix} 0.004 & -0.104 \end{bmatrix}$$

$$\text{For } X^{(c)} \rightarrow \begin{bmatrix} 0.07 & -0.02 \end{bmatrix}$$

Gradients of  $u$

$$\text{For } X^{(1)} \rightarrow \begin{bmatrix} -0.04 & 0.11 \end{bmatrix}$$

$$\text{For } X^{(2)} \rightarrow \begin{bmatrix} -0.08 & 0.04 \end{bmatrix} \xrightarrow{\text{avg}} \begin{bmatrix} -0.05 & 0.01 \end{bmatrix}$$

$$\text{For } X^{(c)} \rightarrow \begin{bmatrix} -0.04 & -0.02 \end{bmatrix}$$

Subject:

Date:

Gradient of  $b^{(1)}$  - second layer

$$\text{for } x^{(1)} \rightarrow [-0.1, 0.0]$$

$$\text{for } x^{(2)} \rightarrow [0.4, 1] \xrightarrow{\text{avg}} [0.15]$$

$$\text{for } x^{(3)} \rightarrow [0.4, 1]$$

second iteration weight updates

$$W \Rightarrow \begin{cases} V_r \leftarrow \delta V_i - \frac{1}{n} \left[ \sum \frac{\partial \mathcal{L}_i}{\partial W} \right] \\ W_{\text{new}} \leftarrow W_{\text{old}} + V_r \end{cases}$$

$$\Rightarrow \begin{cases} V_r \leftarrow 0.9 \begin{bmatrix} 0.02 & -0.04 & -0.15 \\ -0.01 & 0.00 & 0.11 \end{bmatrix} - 0.1 \begin{bmatrix} -0.02 & 0.00 & 0.02 \\ 0.00 & -0.01 & -0.01 \end{bmatrix} \\ W_{\text{new}} \leftarrow \begin{bmatrix} -0.02 & 0.00 & 0.02 \\ 0.00 & -0.01 & -0.01 \end{bmatrix} + \begin{bmatrix} 0.004 & -0.04 & -0.04 \\ -0.004 & 0.00 & 0.02 \end{bmatrix} \end{cases}$$

$$W_{\text{new}} = \begin{bmatrix} 0.04 & 0.12 & -0.15 \\ -0.02 & 0.12 & 0.12 \end{bmatrix}$$

$$b^{(1)} \Rightarrow \begin{cases} V_r \leftarrow \delta V_i - \frac{1}{n} \left[ \sum \frac{\partial \mathcal{L}_i}{\partial b^{(1)}} \right] \\ b_{\text{new}}^{(1)} \leftarrow b_{\text{old}}^{(1)} + V_r \end{cases}$$

$$\Rightarrow \begin{cases} V_r \leftarrow 0.9 \begin{bmatrix} -0.01 & 0.01 \end{bmatrix} - 0.1 \begin{bmatrix} 0.004 & -0.004 \end{bmatrix} \\ b_{\text{new}}^{(1)} \leftarrow \begin{bmatrix} -0.01 & -0.00 \end{bmatrix} + \begin{bmatrix} -0.01 & 0.1 \end{bmatrix} \end{cases}$$

$$b_{\text{new}} = \begin{bmatrix} -0.02 & -0.00 \end{bmatrix}$$



Scanned with CamScanner

sam

$$u \Rightarrow \begin{cases} v_i \leftarrow \delta v_i - \frac{12}{n} \left[ \sum \frac{\partial \mathcal{L}_i}{\partial u} \right] \\ u_{\text{new}} \leftarrow u_{\text{old}} + v_i \end{cases}$$

$$\begin{cases} v_i \leftarrow 0.9 [0.04 \quad -0.10 \Delta] - 0.1 \Delta [-0.08 \quad 0.01] \\ u_{\text{new}} \leftarrow [0.04 \quad -0.10] + [0.04 \quad -0.02] \end{cases}$$

$$u_{\text{new}} = [0.08 \quad -0.12]$$

$$b \Rightarrow \begin{cases} v_i \leftarrow \delta v_i \leftarrow \frac{12}{n} \left[ \sum \frac{\partial \mathcal{L}_i}{\partial b^{(r)}} \right] \\ b_{\text{new}}^{(r)} \leftarrow b_{\text{old}}^{(r)} + v_i \end{cases}$$

$$\Rightarrow \begin{cases} v_i \leftarrow 0.9 [-0.1 \Delta] + 0.1 \Delta [0.1 \Delta] \\ b_{\text{new}}^{(r)} \leftarrow [-0.08] + [-0.04] \end{cases}$$

$$b_{\text{new}}^{(r)} = [-0.12]$$

Final weights after 2 iterations

$$W = \begin{bmatrix} 0.04 & 0.12 & -0.04 \\ -0.08 & 0.12 & 0.12 \end{bmatrix}$$

$$u = [0.08 \quad -0.12]$$

$$b^{(1)} = [-0.08 \quad -0.04]$$

$$b^{(r)} = [-0.12]$$



سوال دوم)

(الف)

- می‌دانیم در GD با momentum پارامتر سرعت محاسبه می‌شود که جهت گرادیان‌های قبلی را در نظر می‌گیرد و درواقع exponentially decaying average گرادیان‌های قبلی را حساب می‌کند، این کار باعث می‌شود oscillation ها کمتر شود زیرا گرادیان‌های قبلی را در نظر می‌گیرد (هرچه گرادیان ها برای پیمایش‌های قدیمی‌تر باشند تاثیرشان کمتر است). با توجه به اینکه گفته شده نمودار آبی رنگ GD است، می‌توان فهمید که نمودار سبز رنگ مربوط به momentum است زیرا بسیار شبیه به GD است ولی با stability کمتر. از طرفی در nesterov-momentum روی momentum کار شده تا آنرا از momentum هم بهتر کند و این کار را با تغییر نحوه بروزرسانی پارامترهای مدل و محاسبه گرادیان پارامترها در نقطه‌ای کمی جلوتر از حالت فعلی (lookahead) با کمک جهت سرعت محاسبه شده مانند حالت قبل می‌کند. در نتیجه با این نگاه به جلو باز هم از oscillation ها کاسته می‌شود. این تغییر نحوه بروزرسانی پارامترها و خود سرعت را هوشمندتر می‌کند. درواقع تفاوت اصلی آن با momentum در محاسبه گرادیان در هر پیمایش است که در nesterov-momentum عنصر آینده‌نگری به آن اضافه شده است و گرادیان پارامترها با فرض جهت حرکت کنونی حساب می‌شود.

SGD with momentum

$$v \leftarrow \beta v - \alpha \nabla L(w)$$

$$w \leftarrow w + v$$

SGD with nesterov-momentum

$$w_{lookahead} = w + \beta v$$

$$v \leftarrow \beta v - \alpha \nabla L(w)$$

$$W \leftarrow W + v$$

بنابراین نمودار قرمز رنگ برای nesterov-momentum است زیرا دارای oscillation های بسیار کمتر است.

دو روش بهینه‌سازی بر اساس momentum دارای نرخ یادگیری ثابت هستند اما RMSProp در حوزه الگوریتم‌های adaptive learning است و نرخ یادگیری آن با توجه به پیشروی گرادیان‌ها کاهش پیدا می‌کند. در واقع نرخ یادگیری هر پارامتر به صورت جداگانه با توجه به exponential moving average گرادیان‌های آن پارامتر از اول تا آن پیمایش به صورت معکوس کاهش می‌یابد. در نتیجه با نزدیک شدن به نقطه optimal نرخ یادگیری نیز کاهش پیدا می‌کند و پرش‌ها بسیار کمتر می‌شوند. در نتیجه RMSProp در این مثال از بروزرسانی‌های اضافی و با گام‌های بزرگ پیشگیری می‌کند و باعث می‌شود به صورت مستقیم به سمت نقطه بهینه حرکت کند. در واقع از حرکات زیگزاگی و پرش‌های بلند که جلوگیری می‌شود مسیر بسیار هموار و direct تر می‌شود زیرا همواره در جهت و مسیر درست می‌ماند.

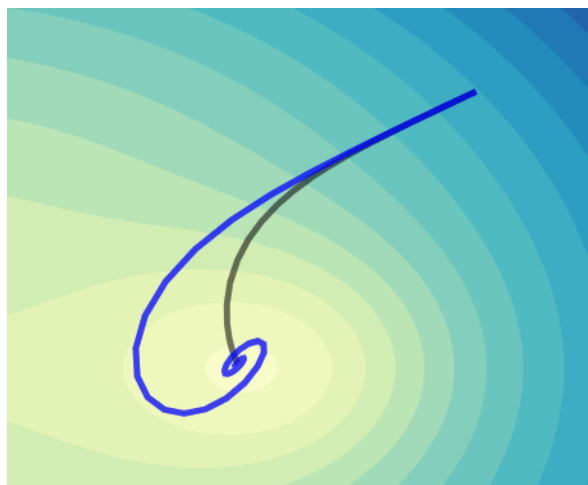
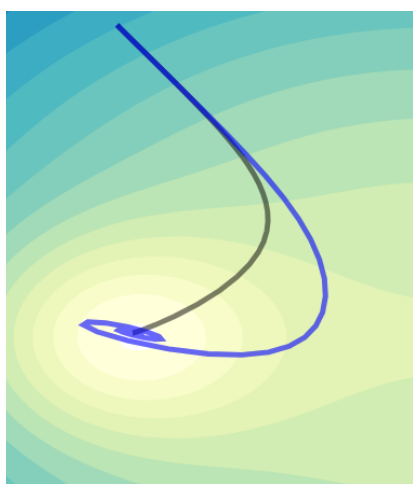
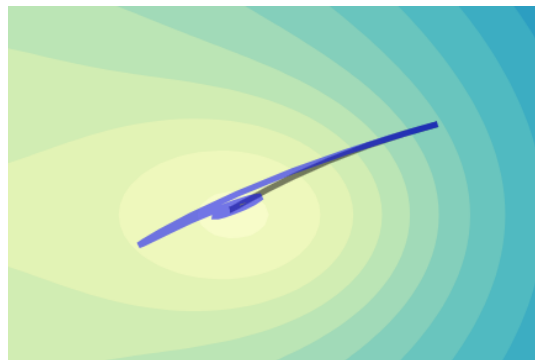
## RMSProp

$$r_t \leftarrow \rho r_{t-1} + (1 - \rho)g^2$$

$$w \leftarrow w - \alpha / (\sqrt{r_t} + \delta)g$$

- **مزایا و معایب روش momentum:** در این روش سعی شده تا با تعریف پارامتر سرعت که با توجه به گرادیان‌های پیمایش‌های قبلی محاسبه می‌شود فرایند بروزرسانی در هر پیمایش هموارتر و پایدارتر شود تا از oscillation ها تا حدی جلوگیری شود و نقطه بهینه overshoot نشود. در نتیجه باعث می‌شود نسبت به GD معمولی نرخ همگرایی بالاتری داشته باشد. از معایب این روش می‌توان به نیازمندی دو ابرپارامتر نرخ یادگیری و ضریب momentum اشاره کرد. اما مورد مهم تر در این روش که باعث شده تا در شکل داده شده عملکرد آن از GD بدون momentum بدتر شود، وابستگی و حساسیت زیاد این روش نسبت به نقطه اولیه گرادیان‌ها و چند گام اولیه می‌باشد با توجه به اینکه در روش momentum پارامتر سرعت نیز در نظر گرفته می‌شود در نتیجه هر

بروزرسانی گرادین‌ها بایاسی به سمت گام‌های قبلی دارد و همین مورد باعث می‌شود بعضی‌جاها با اینکه نباید جلو برود می‌رود و پارامترها را به اشتباه بروزرسانی می‌کند. در پایین این وابستگی به شرایط اولیه نشان داده شده است. (آبی مربوط به momentum است و مشکی SGD معمولی)



مزایا و معایب روش **nestrov-momentum**: این روش دارای مزایای روش momentum نیز می‌باشد اما پیش از بروزرسانی گرادین‌ها یک **further adjustment** روی پارامتر سرعت انجام می‌شود به این صورت که در جهت سرعت گرادین نقطه‌ای کمی جلوتر محاسبه می‌شود (به جای گرادین نقطه کنونی) و درواقع جلوی پیش را نگاه



می‌کند، این تکنیک دقیقا مشکل آخری که در momentum در بخش معایب آن مطرح شد را نشانه قرار می‌دهد و آن را حل می‌کند. حاصل این تغییر بروزرسانی گرادیان با توجه به نقطه‌ای کمی جلوتر بروزرسانی هوشمندانه‌تر پارامتر هاست، مخصوصا در محیط‌هایی که گرادیان‌ها با کمی تغییر بسیار oscillate می‌کنند. در نتیجه این روش نسبت به روش momentum دارای پایداری و در نتیجه نرخ یادگیری بالاتری می‌باشد. از معایب آن ثابت بودن نرخ یادگیری است که می‌تواند منجر به overshoot کردن نقطه بهینه شود زیرا به طور کلی ما با نزدیکتر شدن به نقطه بهینه به گام‌های کوچک‌تری در بروزرسانی پارامترها نیاز داریم همینطور همگرایی این روش بسیار حساس به نرخ یادگیری می‌باشد و نیاز به تنظیم دقیق دارد. از دیگر معایب این روش که شامل روش قبل هم می‌شود، می‌توان به گیر کردن در نقطه بهینه محلی اشاره کرد زیرا سعی می‌کند جهت بروزرسانی پارامترها را حفظ کند و این مهم می‌تواند باعث oscillation حول نقطه بهینه شود زیرا نرخ یادگیری نیز ثابت است.

مزایا و معایب روش RMSProp: هدف اصلی این روش پیاده سازی adaptive learning است و این کار را با کوچک کردن نرخ یادگیری هر پارامتر به صورت جداگانه با توجه به میانگین نمایی مربع گرادیان‌های گذشته آن انجام می‌دهد. این کار باعث می‌شود تا ما در نقاط نزدیک به بهینه آنرا overshoot نکنیم زیرا گام‌های بروزرسانی ما به مرور زمان کوچک‌تر می‌شود. همینطور RMSProp اهمیت کمتری به گرادیان‌های محاسبه‌شده در گام‌های قدیمی می‌دهد که باعث بروزرسانی هوشمندانه‌تر می‌شود و پایداری و همواری بیشتر می‌شود در نتیجه پارامترهایی که دارای گرادیان‌های کوچک‌تری بوده‌اند بروزرسانی‌های قوی‌تری روی آن‌ها انجام می‌شود و برعکس. از معایب این روش می‌توان به حساس بودن به نرخ یادگیری اولیه و نبود momentum و پارامتر سرعت و در نتیجه در نظر گرفتن جهت تغییر گرادیان‌ها اشاره کرد. همچنان تنظیم غیر مناسب پارامتر وزن در میانگین نمایی نیز می‌تواند نرخ همگرایی را به شدت پایین آورد. بار محاسباتی آن نیز کمی بالاتر است زیرا با میانگین گرادیان‌ها را در هر مرحله برای همه پارامترها حساب کند.

- الگوریتم AdaGrad نیز مانند RMSProp از روش‌های adaptive learning می‌باشد ولی نحوه کوچک‌تر کردن پارامتری‌های متفاوت است. در این روش کوچک کردن با توجه به

به حاصلی جمع مربع‌های گرادیان‌های گذشته رخ می‌دهد و میانگینی گرفته نمی‌شود. همین مورد باعث می‌شود این الگوریتم برای فرایندهای یادگیری طولانی مناسب نباشد زیرا نرخ یادگیری بدون هیچ  $\text{limit}$ ی می‌تواند کوچک شود که موجب جلوگیری از پیشروی پارامترها می‌شود یا سرعت آن را بسیار کند می‌کند. این مهم در RMSProp با میانگین گیری مورد توجه قرار گرفته است. با توجه به نکات گفته شده در بخش‌های قبل به نظر می‌رسد مسیر آن بسیار شبیه به RMSProp باشد با این تفاوت که نرخ همگرایی آن کند تر می‌شود و حتی ممکن است به نقطه بهینه نرسد زیرا ممکن است گام‌ها بسیار کوچک شوند، البته این موضوع به نرخ یادگیری اولیه ارتباط دارد همینطور با توجه به مشخص بودن تعداد استپ‌ها در مسیرهای رسم شده در GD به نظر می‌رسد فرایند یادگیری طولانی نباشد و AdaGrad نیز به نقطه بهینه برسد. (قبل از بسیار کوچک‌تر شدن استپ‌ها)

Subject:

Date:

$$s_t = \beta_r s_{t-1} + (1-\beta_r) g_t^r$$

$$s_{t-1} = \beta_r s_{t-2} + (1-\beta_r) g_{t-1}^r$$

$$* \quad \beta_r s_{t-1} = \beta_r (\beta_r s_{t-2} + (1-\beta_r) g_{t-1}^r)$$

$$* * \quad (\beta_r)^2 s_{t-2} = (\beta_r)^2 s_{t-3} + (\beta_r)^2 (1-\beta_r) g_{t-2}^r$$

$$* \quad s_t = (\beta_r)^r s_{t-r} + (1-\beta_r) g_{t-r+1}^r + (1-\beta_r) g_t^r$$

$$* * \quad s_t = (\beta_r)^r s_{t-r} + (\beta_r)^r (1-\beta_r) g_{t-r+1}^r + (1-\beta_r) g_t^r$$

$$s_t = (\beta_r)^n s_{t-n} + (1-\beta_r) g_t^r + \beta_r (1-\beta_r) g_{t-1}^r + \beta_r^2 (1-\beta_r) g_{t-2}^r + \dots + \beta_r^{n-1} (1-\beta_r) g_{t-n+1}^r$$

$$s_t = 0 \Rightarrow s_t = \sum_{i=0}^{t-1} \beta_r^i (1-\beta_r) g_{t-i}^r = (1-\beta_r) \sum_{i=0}^{t-1} \beta_r^i g_{t-i}^r$$

$$s_t = (1-\beta_r) \sum_{i=0}^{t-1} \beta_r^i g_{t-i}^r$$

$$E[s_t] = (1-\beta_r) \sum_{i=0}^{t-1} \beta_r^i E[g_{t-i}^r]$$

$$= (1-\beta_r) \mu \sum_{i=0}^{t-1} \beta_r^i$$

$$\sum_{i=0}^{n-1} ar^i = a \frac{r^n - 1}{r - 1}$$

$$E[s_t] = (1-\beta_r) \mu \frac{\beta_r^{t-1} - 1}{\beta_r - 1} = \mu (1-\beta_r^{t-1})$$

$$\mu = E[g_{t-i}^r]$$

average of  $g_{t-i}^r$ 's



- در الگوریتم Adam تصحیح بایاس روی هردوی گشتاورهای اول و دوم انجام می‌شود. دلیل اصلی این کار این است که هر دوی این گشتاورها با  $\beta$  در پیمایش اولیه مقداردهی می‌شوند. بدون تصحیح بایاس در پیمایش‌های اول الگوریتم مقادیر محاسبه شده برای این گشتاورها حول  $\beta$  بایاس هستند. در نظر داریم که ما در هر مرحله تخمینی از این گشتاورها را محاسبه می‌کنیم، پس برای محاسبه یک unbiased estimator از این گشتاورها از فرمول زیر استفاده می‌کنیم.

$$\hat{s} = s / (1 - \beta_2)$$

با توجه به نتیجه حاصل‌شده در بخش قبل در شروع فرایند یادگیری و در پیمایش‌های اولیه  $s_t$  و  $v_t$  هر دو بسیار به سمت  $\beta$  هستند زیرا  $\beta^t \approx 1$ .

اما با پیشروی یادگیری و بزرگ شدن  $t$ ،  $\beta^t \rightarrow 0$  در نتیجه دیگر بایاسی نداریم.

$$E[s_t] = \mu(1 - \beta_2^{t-1})$$

در نتیجه بدون این تصحیح بایاس الگوریتم ما در پیمایش‌های اولیه به اشتباه گرادیان‌ها را بسیار کوچک و نزدیک به  $\beta$  فرض می‌کند که منجر به گام‌های بزرگ در این مراحل می‌شود. این مهم باعث ناپایداری و غیرهموار بودن الگوریتم در مراحل اولیه می‌شود که می‌تواند باعث کاهش نرخ همگرایی و ضعف الگوریتم در پیدا کردن مسیر بهینه به سمت نقطه بهینه شود.

- می‌دانیم که Adam ایده‌های هردوی RMSProp و momentum را با هم ادغام کرده است. گشتاور دوم از RMSProp و گشتاور اول از momentum. از مزایای Adam نسبت به RMSProp به در نظر گرفتن جهت بروزرسانی گرادیان‌ها با بردار سرعت محاسبه شده اشاره می‌کنیم که منجر به بروزرسانی هوشمندانه تر می‌شود. (جزئیات در توضیحات مربوط به GD با مومنتوم ذکر شده.) همچنین در Adam بر خلاف RMSProp تصحیح بایاس وجود دارد که جزئیات آن را نیز گفتیم. اما از طرفی از مزایای RMSProp می‌توان به سریع‌تر بودن آن اشاره کرد زیرا گشتاور اول را در نظر نمی‌گیرد و در نتیجه محاسبات کمتری دارد. همچنین در محیط‌های convex بررسی شده است که RMSProp به خوبی Adam عمل می‌کند که با توجه به سریع‌تر بودن، استفاده کمتر از حافظه و ابرپارامتر کمتر از Adam می‌تواند انتخاب بهتری باشد.

سوال سوم)

(الف)

اگر به جای  $\hat{x}, y$  بگذاریم درواقع مثل این می ماند که در تمام شرایط داده های ورودی هر لایه شبکه ما دارای میانگین ۰ و واریانس ۱ است زیرا در هر مرحله صرفاً روی هر دسته z-score normalization انجام می دهیم و پارامترهای نرمالیزیشن  $\gamma$  و  $\beta$  آموزشی نمی بینند و shift و scale رخ نمی دهد در نتیجه داده در هر لایه محدود در بازه  $[-1, 1]$  یا  $[0, 1]$  می شود (بستگی به تابع فعال ساز دارد) که اتفاق خوشایندی نیست زیرا flexibility مورد نیاز را برای مسائل غیر خطی را از شبکه گرفته و ممکن است جلوی این را بگیرد که شبکه نحوه نمایش داده ها را در مراحل مختلف به صورت بهینه یاد بگیرد. همچنین در بسیاری از مسائل نیاز به این است که داده در range زیادی کشیده شده باشد تا فرایند یادگیری بتواند انجام شود زیرا تاثیر یک پارامتر را نمی توانیم به اندازه کافی کم و زیاد کنیم و فیدبک بگیریم. پس در نهایت با اعمال تغییر گفته شده قدرت یادگیری مدل کاهش می یابد. در نتیجه باید پارامترهای  $\gamma$  و  $\beta$  را یاد بگیریم تا به صورت بهینه از BN بتوانیم استفاده کنیم.

(ب)

- اگر در BN سائز هر دسته از حدی کوچکتر باشد برآورد و تخمین هایی که ما از میانگین و واریانس داده داریم به unrepresentative خواهد بود و درست نخواهد بود. به همین دلیل این تخمین ها در هر دسته بسیار با هم متفاوت خواهند بود و درواقع این نرمال سازی با یک سری آماره noisy انجام شده است. با وجود این نویز خروجی توابع فعال ساز می توانند بسیار غیر قابل پیش بینی باشند و شیفت زیادی در هر دسته نسبت به هم داشته باشند. این موجب می شود که تا گرادیان ها پایدار نباشند و فرایند یادگیری پارامترها مختل می شود و یا نرخ همگرایی آن پایین می آید. همچنین می تواند موجب overfit روی این نویزها هم بشود.
- اگر اندازه هر دسته برابر ۱ باشد و مقدار پارامتر اولیه شیفت ۰ باشد، طبق الگوریتم در هر مرحله  $y_i$  برابر صفر می شود زیرا میانگین  $x_i$  با  $x$  برابر خواهد بود و واریانس دسته نیز ۰ خواهد بود (البته تقسیم بر ۰ رخ نمی دهد زیرا e در نظر گرفته شده است) و همچنین چون  $\beta = 0$  در نتیجه شیفتی هم نداریم. در نتیجه همه  $y_i$  ها صفر می شوند

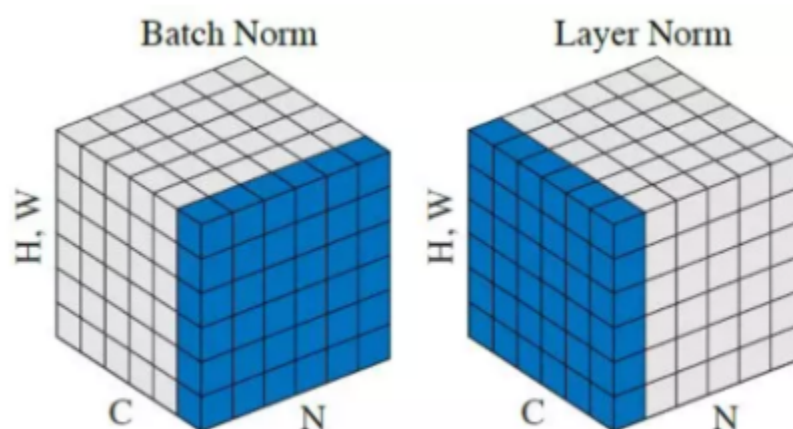
و مقادیر توابع فعال‌ساز به ازای تمام ورودی‌ها با هم برابر می‌شوند و عملاً فرایند یادگیری‌ای نمیتواند رخ دهد زیرا تمام اطلاعات قابل یادگیری از داده از بین می‌رود.

(ج)

با توجه به اینکه در این لایه ۲۰ نورون داریم و ورودی ما ۱۰ بعدی است و لایه ما تماماً متصل است، هر نورون تمام ابعاد ورودی را دریافت می‌کند در نتیجه هر نورون به ۱+۲۰ پارامتر وزن نیاز دارد. با احتساب بایاس. از طرفی BN در سطح هر نورون (ویژگی) و قبل از اعمال تابع فعال‌ساز و مستقل از دیگر نورون‌ها انجام می‌شود تا به ازای ویژگی‌های مختلف فرایندهای نرمالیزشن بهینه را داشته باشیم پس هر نورون ۲ پارامتر BN را نیز دارد، یعنی هر نورون ۲۳ پارامتر قابل یادگیری دارد. ما در کل ۲۰ نورون داریم پس ۲۰\*۲۳ یعنی ۴۶۰ پارامتر آموزشی داریم.

(د)

بزرگ‌ترین تفاوت بین LN و BN این است که در اولی نرمال‌سازی در سطح یک نمونه ولی در بازه تمام نورون‌های یک لایه و در دومی در سطح چند نمونه (سایز دسته) و برای هر ویژگی (نورون) به صورت جداگانه انجام می‌شود. در نتیجه پارامترهای  $\gamma$  و  $\beta$  در LN در هر لایه مشترک هستند ولی در BN هر نورون به صورت جداگانه  $\gamma$  و  $\beta$  دارد. در LN میانگین و واریانس روی ویژگی‌های نمونه انجام می‌شود. در نتیجه دیگر وابستگی به سایز دسته در این روش وجود ندارد.





در شرایطی BN بهتر است که اندازه دسته‌ها بزرگ باشد زیرا اینگونه برآورد ما از میانگین و واریانس داده واقعی‌تر خواهد بود و همینطور نشان داده شده است [\[لینک مقاله\]](#) که در شبکه‌های عصبی feedforward, و به ویژه CNNها نرمال‌سازی در ابعاد دسته بهتر است زیرا نرمال‌سازی دسته‌ای ورودی‌های هر لایه را به گونه‌ای نرمال می‌کند که توزیع ثابتی داشته باشند، و این باعث کاهش (Internal Covariate Shift) می‌شود، جایی که ورودی‌های لایه‌ها در طول آموزش تغییر می‌کنند. این پایداری امکان استفاده از نرخ یادگیری بالاتر و همگرایی سریع‌تر را فراهم می‌کند.

از طرفی زمانی LN بهتر است که شبکه ما از نوع RNN باشد زیرا LN آماره‌های نرمال‌سازی را بر روی ویژگی‌های یک نمونه تکی محاسبه می‌کند، که آن را برای RNNها مناسب می‌سازد زیرا طول sequence ممکن است متغیر باشد و وابستگی‌ها بین گام‌های زمانی و نمونه‌ها اهمیت دارند. استفاده از نرمال‌سازی دسته‌ای (Batch Normalization یا BN) در RNNها به دلیل طول متغیر توالی‌ها و نیاز به آماره‌های یکنواخت در تمام گام‌های زمانی، چالش‌برانگیز است. در کل در زمانی که داده ما temporal است یعنی وابسته به زمان، بهتر است از LN استفاده کنیم. همچنین در سناریوهایی با اندازه‌های کوچک دسته اتکای BN به آماره‌های دسته‌ای باعث ایجاد واریانس بالا و ناپایداری می‌شود. LN که مستقل از اندازه‌ی دسته است، در این شرایط پایدار باقی می‌ماند. LN در حین آموزش و استنتاج همان محاسبات را انجام می‌دهد، زیرا وابسته به آماره‌های دسته‌ای نیست. این سازگاری می‌تواند در سناریوهای پیاده‌سازی که حفظ ساختار یکسان بین آموزش و inference مطلوب است، مزیت داشته باشد.

سوال چهارم)

(الف)

درآپ اوت و منظم‌ساز هر دو به منظور جلوگیری از Overfitting استفاده می‌شوند، اما از روش‌های متفاوتی برای دستیابی به این هدف بهره می‌برند.

- تضعیف وابستگی به نورون‌های خاص: درآپ اوت به صورت تصادفی تعدادی از نورون‌ها را در طول آموزش غیرفعال می‌کند. این باعث می‌شود که مدل به ترکیب خاصی از

وزن‌ها یا نورون‌ها وابسته نباشد. منظم‌ساز نیز وابستگی مدل به مقادیر بزرگ وزن‌ها را کاهش می‌دهند.

- ایجاد یکپارچگی در یادگیری: دراپوت عملاً ترکیبات مختلفی از نورون‌ها را در هر مرحله از آموزش ایجاد می‌کند این کار به نوعی نقش عبارت تنظیمی را ایفا می‌کند زیرا پیچیدگی مدل را کاهش می‌دهد.
- کاهش پیچیدگی مدل: منظم‌سازها وزن‌ها را کنترل می‌کنند تا از رشد بیش از حد پارامترها جلوگیری کنند. دراپوت هم از طریق غیرفعال کردن نورون‌ها به نوعی تعداد پارامترهای فعال را کاهش می‌دهد، که اثری مشابه دارد.

این روش از co-adapting بیش از حد نورون‌ها جلوگیری می‌کند و شبکه را مجبور می‌سازد تا ویژگی‌های مقاوم‌تر و مستقل‌تری را یاد بگیرد. همچنین یک اثر مجموعه‌ای (Ensemble Effect) دارد، زیرا هر تکرار آموزشی را می‌توان به عنوان آموزش یک زیرشبکه متفاوت در نظر گرفت.

بنابراین، دراپوت به طور غیرمستقیم شبیه منظم‌ساز عمل می‌کند، زیرا هر دو به کاهش وابستگی مدل به ویژگی‌ها یا پارامترهای خاص کمک کرده و احتمال بیش‌برازش را کاهش می‌دهند.

(ب)

بله، می‌توان گفت دراپوت عملکردی شبیه یادگیری جمعی (Ensemble Learning) دارد. در زیر توضیح این شباهت را ارائه می‌کنم:

- ایجاد زیرشبکه‌های متفاوت: در هر بار استفاده از دراپوت در طول آموزش، به صورت تصادفی بخشی از نورون‌ها غیرفعال می‌شوند. این کار باعث می‌شود که هر بار یک زیرشبکه متفاوت از شبکه اصلی آموزش ببیند. در نتیجه، فرآیند آموزش مشابه آموزش چندین مدل (یا مجموعه‌ای از مدل‌ها) است.
- تاثیر ترکیب مدل‌ها: در یادگیری جمعی، ترکیب خروجی‌های چندین مدل برای بهبود عملکرد استفاده می‌شود. در دراپوت نیز، هنگام پیش‌بینی، وزن‌های نورون‌های

غیرفعال شده به صورت میانگین گیری مقیاس بندی می شوند. این کار باعث می شود شبکه نهایی رفتاری مشابه ترکیب چندین مدل داشته باشد.

• افزایش تعمیم پذیری: همانطور که در بخش قبل توضیح داده شده این روش تعمیم پذیری مدل را افزایش می دهد که باز هم شبیه به خواص یادگیری جمعی است.

در نتیجه، می توان گفت که دراپ اوت به طور ضمنی مانند یادگیری جمعی عمل می کند، زیرا عملکرد کلی شبکه حاصل از ترکیب زیر شبکه های مختلفی است که در طی فرآیند آموزش ایجاد می شوند.

(ج)

با توجه به اینکه گفته شده هر نود به صورت مستقل می تواند روشن یا خاموش باشد به صورت زیر عمل می کنیم.

یکبار فرض می کنیم از  $N$  گره هیچ کدام خاموش نشده،

یکبار فرض می کنیم از  $N$  گره 1 گره خاموش شده،

یکبار فرض می کنیم از  $N$  گره 2 گره خاموش شده،

و ...

در نتیجه داریم:

$$\sum_{k=0}^n \binom{n}{k} = 2^n$$

(د)

در زمان آزمایش، دراپ اوت معمولاً اعمال نمی شود و خروجی شبکه با احتمال دراپ اوت (p) مقیاس بندی می شود تا سازگاری با مقدار مورد انتظار در زمان آموزش حفظ شود.

دلیل این مقیاس بندی: در طول آموزش، نورون ها به صورت تصادفی غیرفعال می شوند و شبکه فقط بخشی از وزن ها را استفاده می کند. در زمان تست، تمام نورون ها فعال هستند. اگر

خروجی شبکه بدون مقیاس‌بندی استفاده شود، مقدار خروجی به‌طور قابل‌توجهی بزرگ‌تر از مقدار مورد انتظار در زمان آموزش خواهد شد.

برای جلوگیری از این مشکل، وزن‌های نوروها در زمان آزمایش در مقدار  $1 - p$  ضرب می‌شوند. این کار باعث می‌شود مقدار خروجی شبکه با مقدار میانگین محاسبه‌شده در زمان آموزش برابر باشد.