200100154

1E643
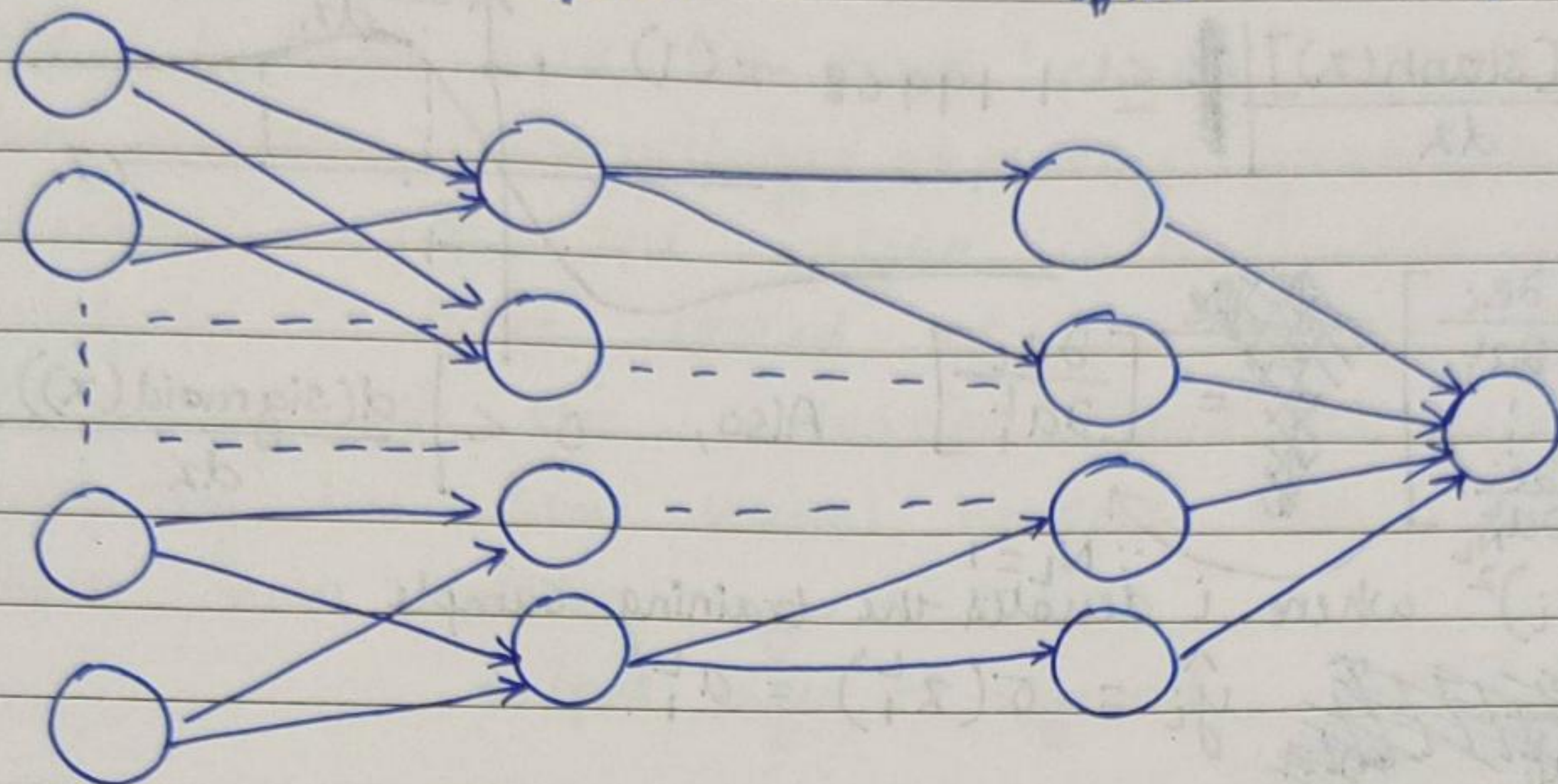
Assignment 2

RANKA
DATE   /   /
PAGE

(d) each hidden layer → stanh activation functions
output layer → logistic sigmoid activation function
784          200          100          1



Input layer    1st Hidden    2nd Hidden    Output
               Layer         Layer         Layer

In a generalized setting, we have shown in class that →

$$\therefore \nabla_{W^l} e = \text{Diag}(\phi^{l'}) V^{l+1} \cdots V^L \delta^L (a^{l-1})^T$$

where $V^{l+1} = \begin{bmatrix} w_{11}^{l+1} & \cdots & w_{N_{l+1}1}^{l+1} \\ \vdots & & \vdots \\ w_{1N_l}^{l+1} & \cdots & w_{N_{l+1}N_l}^{l+1} \end{bmatrix} \begin{bmatrix} \phi'(z_1^{l+1}) & & \\ & \ddots & \\ & & \phi'(z_{N_{l+1}}^{l+1}) \end{bmatrix}$
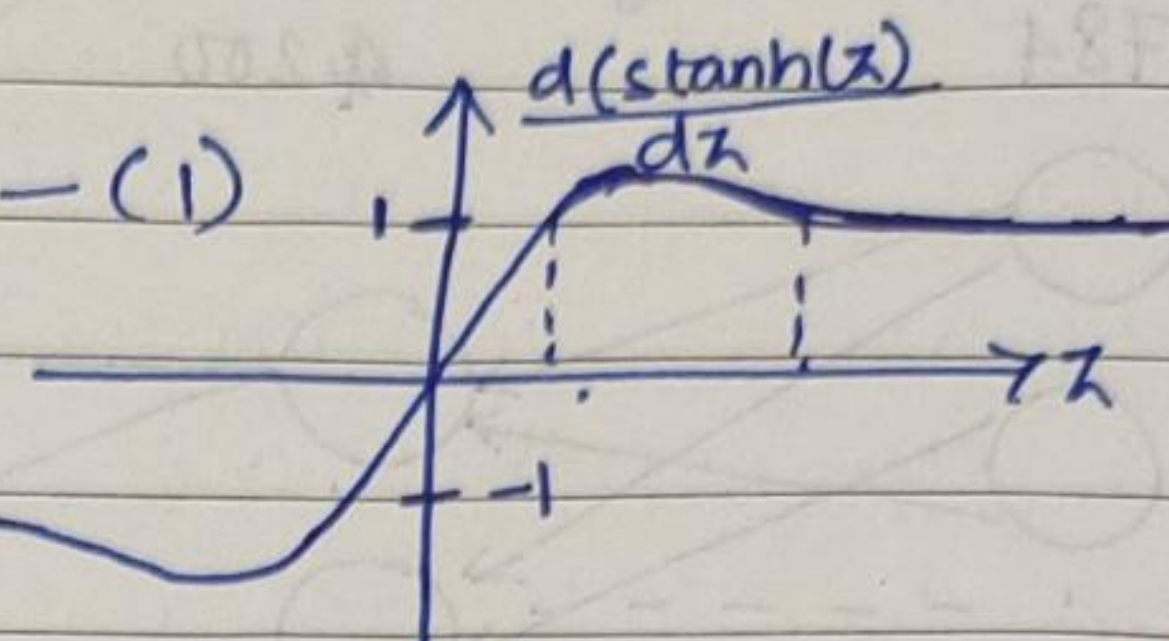
Note that the error gradients at the last layer flow
back into the previous layers

· exploding gradient problem $\equiv \nabla_{W^l} e$ gradients become very
large in magnitude i.e. approach $\infty$

· vanishing gradient problem $\equiv \nabla_{W^l} e$ gradients become very
small, i.e. they approach 0
in magnitude

* The magnitude of $\nabla_{W^l} e$ gradients
depends upon $V^{l+1} V^{l+2} \cdots V^L \delta^L$

* For our analysis, it suffices to only look at the
derivatives of the activation functions and the derivative
of the error function wrt the activation of the output
layer.

$$\text{stanh}(z) = z\,\sigma_{tan}(z) \quad \text{where} \quad \sigma_{tan}(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\text{So,}\ \frac{d[\text{stanh}(z)]}{dz} = \tanh(z) + z\,\text{sech}^2(z) \quad \text{zzzz}$$

$$0.19968 \le \left\|\frac{d[\text{stanh}(z)]}{dz}\right\| \le 1.19968 \quad —(1)$$



$$\text{Now,}\ \delta_L = \begin{bmatrix} \frac{\partial e_i}{\partial a_1^L} \\ \vdots \\ \frac{\partial e_i}{\partial a_{N_L}^L} \end{bmatrix} = \begin{bmatrix} \frac{\partial e_i}{\partial a_1^L} \end{bmatrix} \qquad \text{Also,}\quad 0 < \left|\frac{d(\text{sigmoid}(z))}{dz}\right| \le \frac{1}{4}$$

$$\because N_L = 1$$

- $e_i = (\hat{y}_i - y_i)^2$ where $i$ denotes the training sample
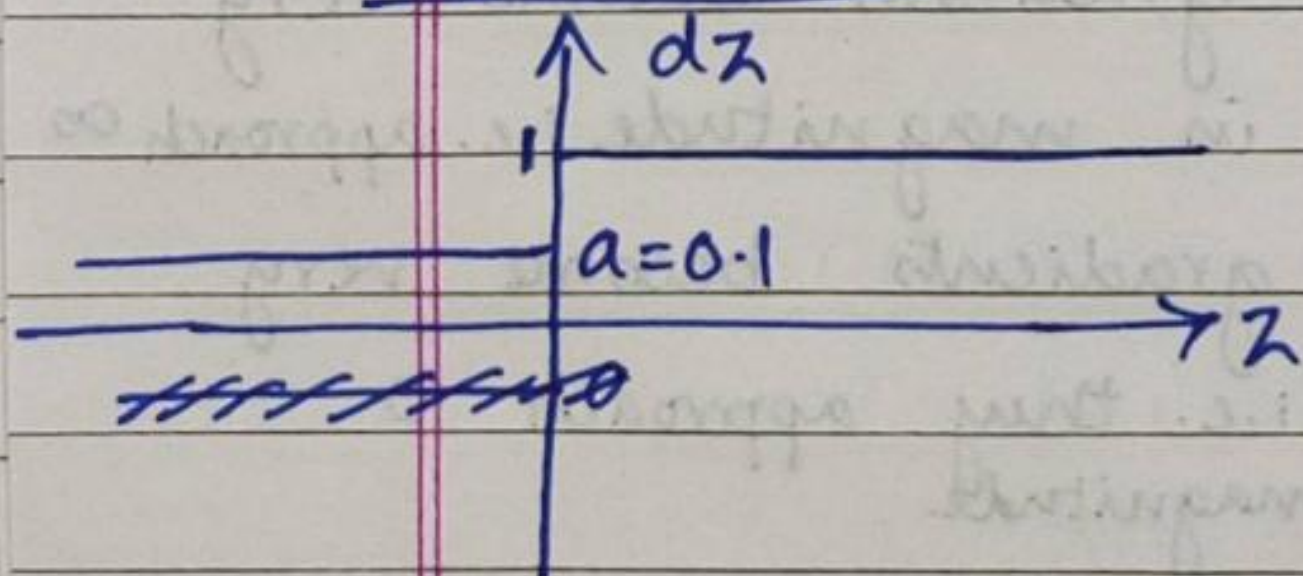- We have $\hat{y}_i = \sigma(z_1^L) = a_1^L$

$$\text{So,}\ \frac{\partial e_i}{\partial a_1^L} = \frac{\partial e_i}{\partial \hat{y}_i} \times \frac{\partial \hat{y}_i}{\partial a_1^L} = 2(\hat{y}_i - y_i)\cdot 1 = 2(\hat{y}_i - y_i) \quad \text{where}\ |\hat{y}_i| \le 1$$

$$\underbrace{\qquad}_{\text{finite}} \qquad \text{and}\ y_i = \{+1, -1\}$$

Due to (1), we can say that as the derivatives or slope gets larger and larger as we go backward with every layer during backprop $\Rightarrow$ Exploding gradients

(e) Considering the same neural network architecture before where now hidden layer $\rightarrow$ a-relu
output layer $\rightarrow$ sigmoid

$$\text{Now,}\quad a\text{-relu}(z) = \begin{cases} az & \text{if } z < 0 \\ z & \text{else} \end{cases} \quad \text{where } a = 0.1$$

$$\text{So,}\quad \frac{d(a\text{-relu}(z))}{dz} = \begin{cases} a & \text{if } z < 0 \\ 1 & \text{else} \end{cases}$$



$$\text{Also,}\quad 0 < \left|\frac{d(\text{sigmoid}(z))}{dz}\right| \le \frac{1}{4}$$

$$0.1 \le \left\|\frac{d(a\text{-relu}(z))}{dz}\right\| \le 1$$

From this bound and the previous discussion on $\frac{\partial e_i}{\partial a^L}$, we can safely interpret that a-ReLU causes vanishing gradients for inputs smaller than zero while for other inputs they allow the

gradient to pass through. The chance for a vanishing gradient can be minimized by taking a larger value of 'a'. Ideally, we would want $a = \lim_{x \to 1^-} x$ i.e. $a = 0.99, 0.999 \cdots$
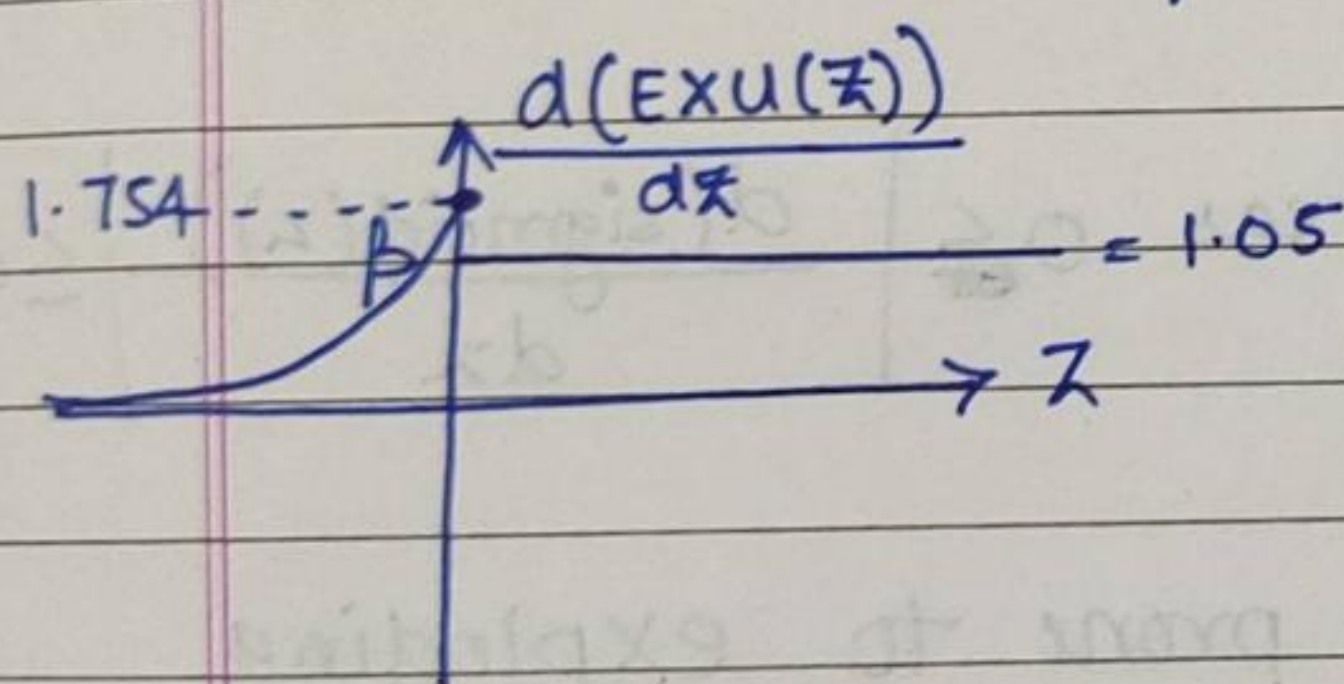
* Note :- I have used the term 'chance' above because this problem of vanishing / exploding gradients ~~allow~~ depends upon the weight initialization. also. For my analysis, I have looked at the bounds of the derivatives of the activation functions used.

(f) Considering the same neural network architecture before where now hidden layer → EXU
output layer → sigmoid

Now, $EXU(z) = \begin{cases} \beta z & \text{if } z \geq 0 \\ \beta y (\exp(z)-1) & \text{else} \end{cases}$

where $\beta = 1.05$ and $y = 1.67$

$\dfrac{d(EXU(z))}{dz} = \begin{cases} \beta & \text{if } z \geq 0 \\ \beta u \exp(z) & \text{else} \end{cases}$



$1.754$

$\dfrac{d(EXU(z))}{dz}$

$= 1.05$

So, $0 \leq \left| \dfrac{d(EXU(z))}{dz} \right| \leq 1.754$

Also, $0 < \left| \dfrac{d(\text{sigmoid}(z))}{dz} \right| \leq \dfrac{1}{4}$

Note, here we may be tempted to say that this model would suffer from exploding gradients since. $\left| \dfrac{d(EXU(z))}{dz} \right|_{max} = 1.754$. However, if we closely examine the graph of $\dfrac{d(EXU(z))}{dz}$ v/s $z$, we will notice that the range of $z$ for ~~ever~~ which $\dfrac{d(EXU(z))}{dz} \geq 1.05$ is extremely small. Rather, this model ~~should~~ should more likely to suffer from vanishing gradients since $\lim_{z \to -\infty} \dfrac{d(EXU(z))}{dz} = 0$

• Networks that were more prone to
exploding

(g) d → stanh, stanh, sigmoid    ~~Explosiong~~ gradient issue → 'd' mode
      with stanh

     e → a-ReLU, a-ReLU, sigmoid

     f → EXU, EXU, sigmoid   Reasons →

                 ○ large weight ~~we~~ initialization

• Networks that were        ○ accumulation of large error

more prone to vanishing     gradients which result in

gradient issue → 'f' model    extremely large updates to

        with EXU           neural network model weights

Reason → • many times           during training.

multiplication of smaller      • |mag. of derivative of stanh|$_{max}$ ⩾1

number (magnitude of activation function's' derivative)

~~function~~ results in a very

small number (close to 0)

• derivative of sigmoid is always below 0.25

• derivative of EXU over a longer range

    of $z$ gives a value close to 0.

(h)    d → stanh, stanh, linear

     e → a-ReLU, a-ReLU, linear

     f → EXU, EXU, linear

$$\left| \frac{d(\text{linear}(z))}{dz} \right| = 1 \qquad \text{whereas} \quad 0 \leq \left| \frac{d(\text{sigmoid}(z))}{dz} \right| \leq 0.25$$

⇒ The model becomes more prone to exploding
gradient problem and becomes less prone to
vanishing gradient problem. due to an increase
in the maximum, value of derivative of activation function
used in the output layer.
               absolute

$$0 \leq \left| \frac{d(\text{ReLU}(z))}{dz} \right| \leq 1 \qquad \Rightarrow \text{The model should be more}$$
                                 stable wrt to both exploding
and vanishing gradient problems due to the
range of $\frac{d(\text{ReLU}(z))}{dz} \in [0,1]$.