

ASSIGNMENT → 2

NAME :- PUSHPENDRA SINGH CHAUHAN

STUDENT ID :- A20472647

COURSE NO. :- CS577

SEMESTER :- SPRING 2021

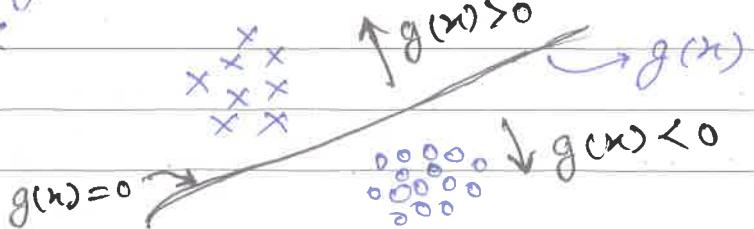
X X X X XArtificial Neurons

Q.1). Given a neuron with linear activation and weights $(0.1, 0.2, 0.3)$. Compute the output for the input vector $(1, 1)$.

Solu: Output for the input vector $(1, 1)$ = $0.2 \times 1 + 0.3 \times 1 + 0.1$
 $= 0.2 + 0.3 + 0.1$
 $= \underline{\underline{0.6}}$

Q.2). Given a linear discriminant function $g(x)$ explain what should be its value on one side of the decision boundary, on the other side of the decision boundary, and on the decision boundary.

Solu: Let $g(x)$ be a linear discriminant function, then:



So,

 $\rightarrow g(x) > 0$ on one side of the decision

boundary.

$\rightarrow g(x) < 0$ on the other side of the decision boundary.

$\rightarrow g(x) = 0$ on the decision boundary.

Q.3) Given a linear discriminant function $g(x_1, x_2) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$. Explain the meaning of the coefficients $\theta_0, \theta_1, \theta_2$.

Solu: (a). Here, for geometric interpretation:

(b). For template matching interpretation:

$$g(x_1, x_2) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

bias or threshold
for strength of match
Template

Q.4>. Let $g(x_1, x_2) = 1 + 2x_1 + 3x_2$. Find the normal to the decision boundary defined by $g(x_1, x_2) = 0$ and the distance of the decision boundary from the origin (note that for the distance computation the normal has to be normalized).

Solu: :- Normal to the decision boundary defined by $\{ g(x_1, x_2) \} = (2, 3)$ Ans.

$$\begin{aligned} \text{Distance of the decision boundary from the origin} &= \frac{1}{\sqrt{(2)^2 + (3)^2}} \\ &= \frac{1}{\sqrt{13}} \cdot \underline{\underline{\text{Ans}}} \end{aligned}$$

Q.5), When writing a discriminant function as $g(x) = \theta^T x$. where x and θ are vectors, explain where the bias coefficient in θ is and what is the required change needed to the feature vector x to allow writing the discriminant in this way.

Solu: The bias coefficient in θ is at the beginning or at 1st position as shown below :

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \xrightarrow{\text{bias coefficient}}$$

In feature vector x , we have to add 1 at the beginning or at 1st position as shown below :

$$x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \xrightarrow{\text{added 1}} \text{to the feature vector } x.$$

Q. 6). Write the expressions for the step and logistic (sigmoid) activation functions. Explain the advantage of the sigmoid activation over step activation. What happens to the sigmoid when θ is small?

Soln:- Expression for step activation function :-

$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x > 0 \\ 0 & \text{otherwise.} \end{cases}$$

→ Expression for sigmoid activation function :-

$$h_{\theta}(x) = \frac{1}{1 + \exp(-\theta^T x)}$$

→ The advantage of the sigmoid activation over step activation is that we can also get intermediate values between 0 and 1 which can give more accurate information probability of an object belonging to a class.

→ The value of sigmoid becomes approximately equal to 0.5 when θ is small.

→ Sigmoid activation function behaves like a step activation function when θ is small.

Q.7). Explain how the sigmoid is obtained using a linear function to model the log-likelihood ratio.

Soln:- Using a linear function to model the log-likelihood ratio:

$$\Rightarrow \log \left(\frac{P(y=1|x)}{P(y=0|x)} \right) = \theta^T x$$

On taking anti-log both sides, we get

$$\Rightarrow \frac{P(y=1|x)}{P(y=0|x)} = \exp(\theta^T x)$$

$$\Rightarrow P(y=1|x) = P(y=0|x) \exp(\theta^T x)$$

On replacing $P(y=0|x)$ with $(1 - P(y=1|x))$, we get

$$\Rightarrow P(y=1|x) = (1 - P(y=1|x)) \exp(\theta^T x)$$

$$\Rightarrow P(y=1|x) \left(1 + \exp(\theta^T x) \right) = \exp(\theta^T x)$$

$$\Rightarrow P(y=1|x) = \frac{\exp(\theta^T x)}{1 + \exp(\theta^T x)}$$

$$\Rightarrow P(y=1|x) = \frac{1}{1 + \exp(-\theta^T x)} = h_{\theta}(x)$$

activation .

sigmoid

Q.87. Write the derivative of a sigmoid and use it to compute the derivative of log-sigmoid.

Soln :- Sigmoid (x) = $\frac{1}{1 + \exp(-x)}$

$$\frac{d}{dx} \text{sigmoid}(x) = \text{sigmoid}(x) (1 - \text{sigmoid}(x))$$

$$= \frac{1}{1 + \exp(-x)} \left\{ 1 - \frac{1}{1 + \exp(-x)} \right\}$$

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

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

$$\frac{d}{dx} \text{log sigmoid}(x) = \frac{1}{\text{sigmoid}(x)} \cdot \text{sigmoid}'(x)$$

$$= 1 - \frac{1}{1 + \exp(-x)}$$

$$= \frac{1 + \exp(-x) - x}{1 + \exp(-x)}$$

$$= \frac{\exp(-x)}{1 + \exp(-x)}$$

Q.9) Explain how to compute the direction used for updating parameters in gradient descent. How is the size of the update controlled?

Solu:- We update the parameters in the opposite direction of the gradient of the objective function with respect to the parameters where the gradient gives the direction of the steepest ascent.

The size of the update is controlled by a hyperparameter called "learning rate".

Q.10) Explain the stop condition for gradient descent. Should the condition use the loss change or parameter change? Explain your answer.

Solu:- When the change in loss is less than a threshold in consecutive iteration, then we have to stop the process of gradient descent.

Date should Condition should use the loss change as we are minimizing the loss function and if loss doesn't change much then we can say that ~~it~~ it is converged but we don't know that smaller change in parameter will result in small change in loss. So, it should use loss change.

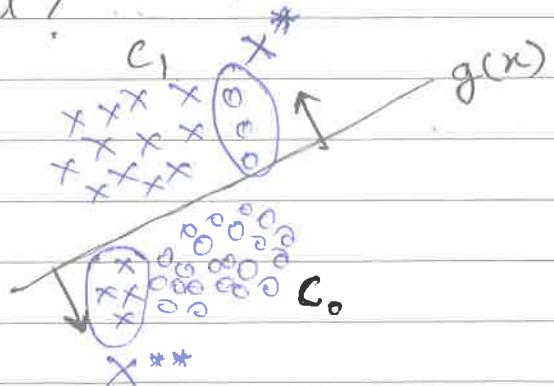
Q.11). Explain what happens when using a learning rate that is too large or too small.

Soln :- • If learning rate is too large, it may fail to converge and overshoot the minimum.

• If learning rate is too small, it would take long time to converge and become computationally expensive.

Q.12). Explain the empirical error loss is computed. What is the problem in using the empirical error loss in gradient descent?

Soln :-



$$x^* = \{x^{(i)} \mid \text{if } (y^{(i)} = 0) \wedge \theta^T x > 0\}$$

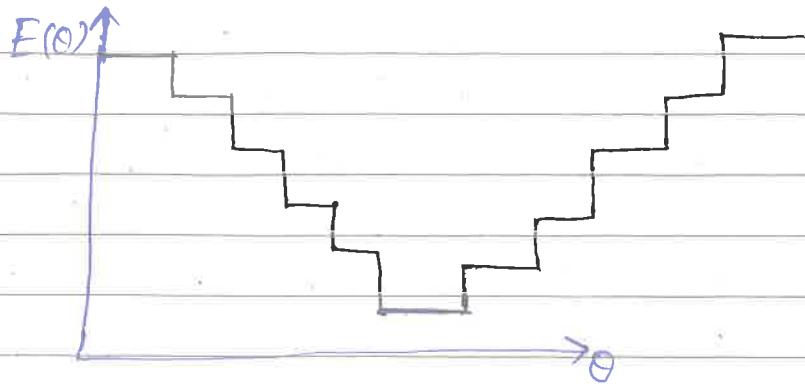
$$x^{**} = \{x^{(i)} \mid \text{if } (y^{(i)} = 1) \wedge \theta^T x < 0\}$$

$$E(\theta) = \# x^* + \# x^{**} = \text{Total no. of errors.}$$

$$\theta^* = \underset{\theta}{\operatorname{arg\,min}} E(\theta)$$

need

$$\nabla E(\theta)$$



The problem in using the empirical error loss in gradient descent is that there exist many values of ' θ ' for same value of ' $E(\theta)$ '.

Q.13) Write the log likelihood for a binary classifier and show how to develop it into binary cross-entropy.

$$\text{Solu: } l(\theta) = -\log \left(\prod_{x^{(i)} \in C_1} P(y=1|x^{(i)}) * \prod_{x^{(i)} \in C_0} P(y=0|x^{(i)}) \right)$$

negative log-likelihood.

$$= -\log \left(\prod_{i=1}^m P(y=1|x^{(i)})^{y^{(i)}} * P(y=0|x^{(i)})^{(1-y^{(i)})} \right)$$

$$= -\sum_{i=1}^m y^{(i)} \log \underbrace{P(y=1|x^{(i)})}_{h_\theta(x)} + (1-y^{(i)}) \log \underbrace{P(y=0|x^{(i)})}_{h_\theta(x)}$$

$$l(\theta) = -\sum_{i=1}^m y^{(i)} \log (h_\theta(x^{(i)})) + (1-y^{(i)}) \log (1-h_\theta(x^{(i)}))$$

binary cross-entropy

Q.12). Write the gradient of the binary cross-entropy loss function with respect to the parameter vector when using a sigmoid activation function. Write the gradient descent update rule when using this loss function. Explain the meaning of the update equation you wrote.

→ binary cross-entropy loss function

$$\text{Solu: } l(\theta) = -\sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))$$

here,

assume $h_{\theta}(x)$ = sigmoid activation function.

then,

$$\begin{aligned} \frac{d l(\theta)}{d \theta} &= -\sum_{i=1}^m y^{(i)} * \frac{1}{h_{\theta}(x^{(i)})} * h_{\theta}(x^{(i)}) (1-h_{\theta}(x^{(i)})) x^{(i)} \\ &+ \sum_{i=1}^m (1-y^{(i)}) * \frac{1}{(1-h_{\theta}(x^{(i)}))} * \left\{ -h_{\theta}(x^{(i)}) (1-h_{\theta}(x^{(i)})) \right\} x^{(i)} \end{aligned}$$

$$= -\sum_{i=1}^m y^{(i)} (1-h_{\theta}(x^{(i)})) x^{(i)} + (1-y^{(i)}) (-h_{\theta}(x^{(i)})) x^{(i)}$$

$$= -\sum_{i=1}^m \left\{ y^{(i)} - y^{(i)} h_{\theta}(x^{(i)}) - h_{\theta}(x^{(i)}) + y^{(i)} h_{\theta}(x^{(i)}) \right\} x^{(i)}$$

gradient of binary cross-entropy

$$\boxed{\frac{d}{d \theta} (-l(\theta)) = \sum_{i=1}^m \{ h_{\theta}(x^{(i)}) - y^{(i)} \} x^{(i)}}$$

Gradient descent update rule:

$$\theta_{\text{new}} = \theta_{\text{old}} - \eta \left\{ \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} \right\}$$

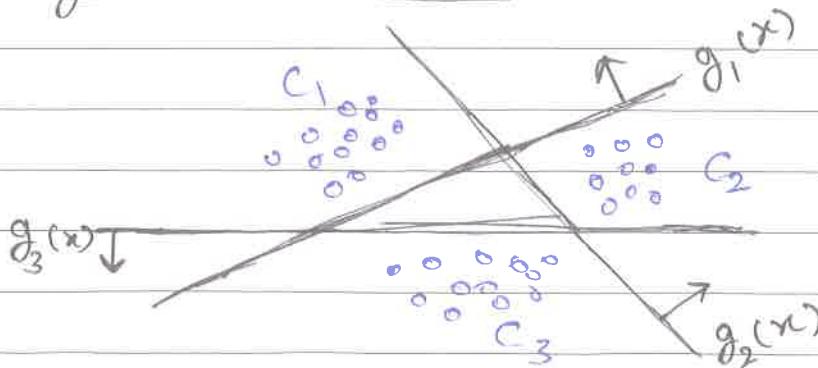
learning rate. $\frac{d}{d\theta} (-l(\theta))$

Here, we calculate the gradient in order to get the direction so that we can update parameters in this direction.

Learning rate decides the size of the update and then subtract them from old parameters to get new parameters.

Q.15). Explain the two strategies for multi-class classification (one against all others, one against each other). In which strategy is it easier to discriminate the data?

Solu:- • One against all others



Here,

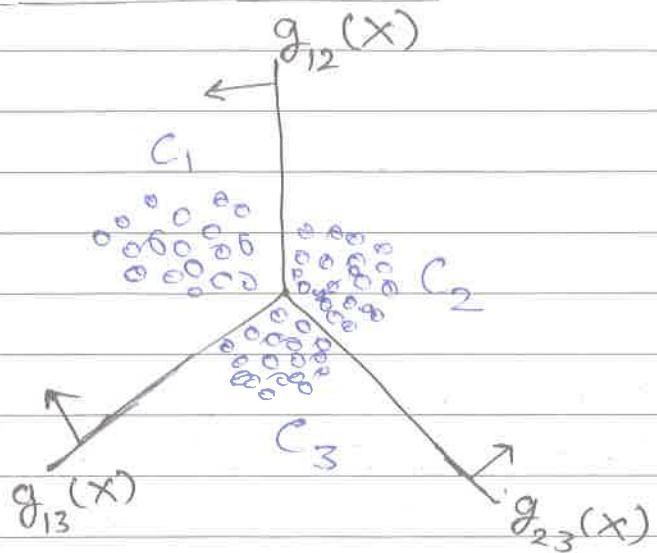
we have 'k' discriminant functions for 'k' classes.

$$g_1(x), \dots, g_k(x)$$

$$\hat{y} = \operatorname{argmax}_j g_j(x)$$

Here, we choose a discriminant function for which data has maximum value.

- One against each other



Hence, we have $\frac{k(k-1)}{2}$ discriminant functions for 'k' classes.

$$\hat{y} = \operatorname{argmax} \sum_{i \neq j} g_{ij}(x)$$

in

- The against each other strategy, it is easier to discriminate the data.

Q.16) Explain the template matching interpretation of linear discriminants. Using this interpretation what is the meaning of using multiple linear discriminant functions.

Solu:- Template Matching Interpretation

- Rows of Θ^T (columns of Θ) are templates.
- With K rows of Θ^T we have K templates (one template per class).
- $\Theta^T x$ measures how well x matches each of the K templates (dot product of x with rows of Θ^T).
- High similarity to a template of a particular class indicates high membership in this class.

The meaning of using multiple linear discriminant functions is to have multiple templates for multi-class classification problem.

Q.17) Explain the need for using softmax in the output layer of a multi-class classifier.

Solu:- We know that, logistic regression produces a decimal between 0 and 1. For example,

a logistic regression output of 0.8 from an email classifier suggests an 80% chance of an email being spam and a 20% chance of it being not spam. Clearly, the sum of the probabilities of an email being either spam or not spam is 1.

Softmax extends this idea into a multi-class world. That is, softmax assigns decimal probabilities to each class. in a multi-class problem. Those decimal probabilities must add upto 1.

Due to this, softmax is needed in the output layer of a multi-class classifier.

Q.18), Write the derivative of softmax and use it to compute the derivative of log-softmax.

Soln :- Given :-

$$\text{softmax}(z_i) = \frac{\exp(z_i)}{\sum_{i=1}^K \exp(z_i)}$$

$$\frac{\partial}{\partial z_i} \text{softmax}(z_i)$$

$$= \text{softmax}(z_i) (\delta_{ij} - \text{softmax}(z_i))$$

$$\frac{\partial}{\partial z_i} (\text{log-softmax}(z_i)) = \frac{1}{\text{softmax}(z_i)} * \{\text{softmax}(z_i) \\ (\delta_{ij} - \text{softmax}(z_i)) \\ = (\delta_{ij} - \text{softmax}(z_i))$$

Q.19) Write the log-likelihood for a multi-class classifier and show how to develop it into categorical cross-entropy.

$$\text{Soln: } l(\theta) = -\log \left(\prod_{i=1}^m \prod_{j=1}^k \underbrace{P(y^{(i)}=j | x^{(i)})}_{h_{\theta_j}(x^{(i)})} \mathbb{1}(y^{(i)}=j) \right)$$

negative log-likelihood

$$= -\sum_{i=1}^m \sum_{j=1}^k \mathbb{1}(y^{(i)}=j) \log (P(y^{(i)}=j | x^{(i)}))$$

$$l(\theta) = -\sum_{i=1}^m \sum_{j=1}^k \mathbb{1}(y^{(i)}=j) \log (h_{\theta_j}(x^{(i)}))$$

↳ Categorical cross-entropy

Q.20) Write the gradient of the categorical cross-entropy loss function with respect to the parameter vector when using a softmax activation function. Write the gradient descent update rule when using this loss function. Explain the meaning of

the update equation you wrote.

$$\text{Solu: } l(\theta) = - \sum_{i=1}^m \sum_{j=1}^k \mathbb{1}(y^{(i)}=j) \log(h_{\theta_j}(x^{(i)}))$$

→ Categorical cross-entropy loss of h .

Here,

assume $h_{\theta_j}(x^{(i)})$ = softmax activation function

Then,

$$\frac{\partial}{\partial \theta_j} (l(\theta)) = - \sum_{i=1}^m \mathbb{1}(y^{(i)}=j) * \frac{1}{h_{\theta_j}(x^{(i)})} \{ h_{\theta_j}(x^{(i)}) * (\mathbb{1}(y^{(i)}=j) - h_{\theta_j}(x^{(i)})) x^{(i)} \}$$

$$\frac{\partial}{\partial \theta_j} (l(\theta)) = \sum_{i=1}^m (h_{\theta_j}(x^{(i)}) - \mathbb{1}(y^{(i)}=j)) x^{(i)}$$

→ gradient of categorical cross-entropy

Gradient descent update rule:

$$\theta_j^{(\text{new})} = \theta_j^{(\text{old})} - \eta \sum_{i=1}^m (h_{\theta_j}(x^{(i)}) - \mathbb{1}(y^{(i)}=j)) x^{(i)}$$

learning rate $\frac{\partial}{\partial \theta_j} (l(\theta))$

Here, we calculate the gradient in order to get the direction so that we can update parameter in this direction.

Learning rate decides the size of the update and then subtract them from old parameter to get new parameter.

NEURAL NETWORKS

Q.1) Given a two layer network where the number of hidden units is larger than the number of inputs. explain the dimensionality increase interpretation of the hidden layer. Explain what is the possible benefit of such dimensionality increase.

Solu:- Dimensionality increase interpretation of the hidden layer:

no. of hidden units > no. of inputs.

$$\begin{matrix} 2D & & 5D \\ (\mathbf{x}_1, \mathbf{x}_2) & \longrightarrow & (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_1^2, \mathbf{x}_2^2, \mathbf{x}_1 \cdot \mathbf{x}_2) \\ \text{input} & & \text{hidden units} \end{matrix}$$

The possible benefit of such dimensionality increase is to establish relationship in higher dimensional space.

Q.2) Given a two layer network where the no. of hidden units is smaller than the no. of inputs explain the dimensionality decrease interpretation of the hidden layer. Explain what is the possible benefit of dimensionality decrease.

Soln :- Dimensionality reduction interpretation of the hidden layer:

no. of hidden units $<$ no. of inputs.

$$\begin{array}{ccc} 6D & & 4D \\ (x_1, x_2, x_3, x_4, x_5, x_6) \rightarrow & (x_1, x_3, x_4, x_6) \\ \text{input} & & \text{hidden units.} \end{array}$$

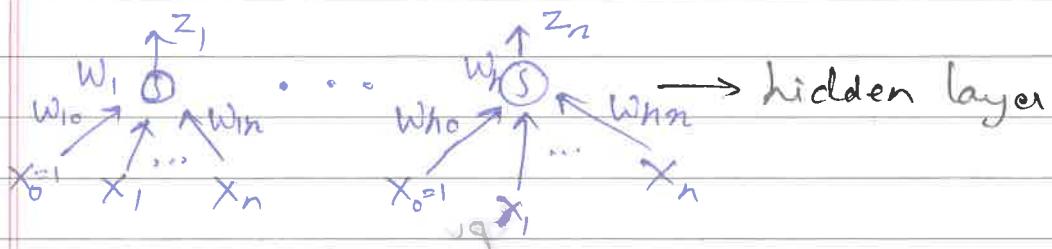
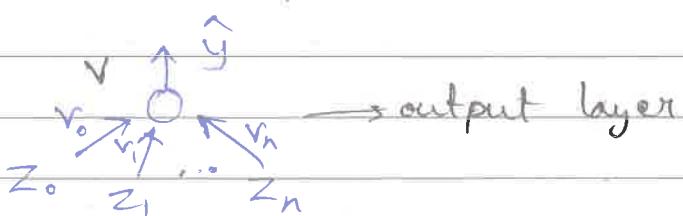
The possible benefit of dimensionality decrease is to extract important features and discard the features which are correlated with each other.

Q.3) Given the update equations for the output layer in a 2 layer network, explain why we cannot use directly the same equations for the hidden layer. Assume that both layers have the same activation function.

Soln :- We cannot use directly the same equations for the hidden layer because it will not take gradient of loss with respect to the weight parameters of hidden layer instead it will take gradient of loss with respect to the weight parameters of the output layer.

Q.4) Apply the chain rule to compute the gradients with respect to both hidden layer and output layer parameters of a 2-layer regression network with a single output where the hidden layers use sigmoid activation and the output layer uses linear activation.

Soln:-



$$\text{Loss: } E = \frac{1}{2} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

$$\text{Outputs: } \hat{y} = v^T z = v_0 + v_1 z_1 + \dots + v_n z_n$$

$$z_j = \text{sigmoid}(w_j^T x)$$

Chain rule for output layer:

$$\begin{aligned} \frac{\partial E}{\partial v} &= \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial v} = \frac{1}{2} \times 2 \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) z^{(i)} \\ &= \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) z^{(i)} \end{aligned}$$

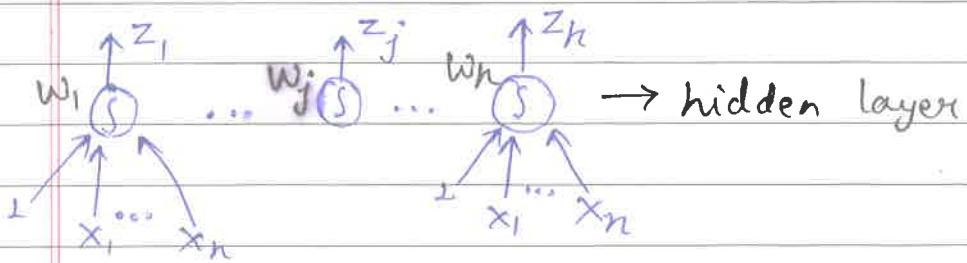
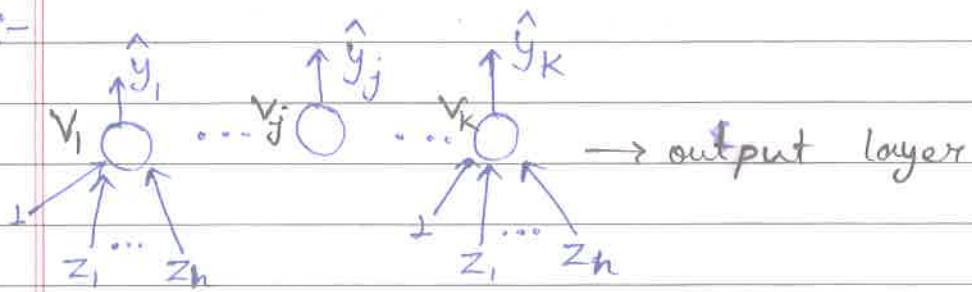
Chain rule for hidden layer:

$$\frac{\partial E}{\partial w_j} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_j} \frac{\partial z_j}{\partial w_j}$$

$$\frac{\partial E}{\partial w_j} = \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) v_j z_j^{(i)} (1 - z_j^{(i)}) x^{(i)}$$

Q.5). Apply the chain rule to compute the gradients with respect to both hidden layer and output layer parameters of a 2-layer regression network with multiple outputs where the hidden layers use sigmoid activation and the output layer uses linear activation

Soln:-



$$\text{Loss} : E = \frac{1}{2} \sum_{i=1}^m \sum_{l=1}^k (\hat{y}_l^{(i)} - y_l^{(i)})^2$$

Outputs : $\hat{y}_j = v_j^T z = v_{j0} + v_{j1} z_1 + \dots + v_{jn} z_n$

$$z_j = \text{sigmoid}(w_j^T x)$$

Chain rule for output layer:

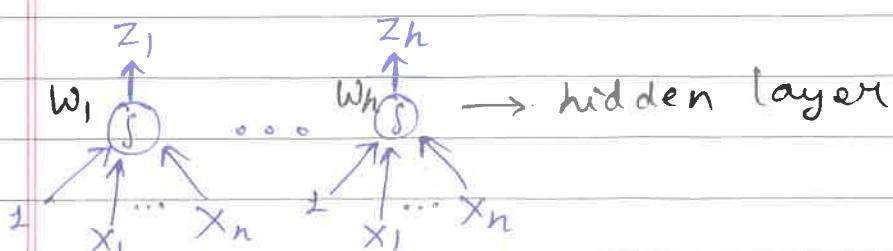
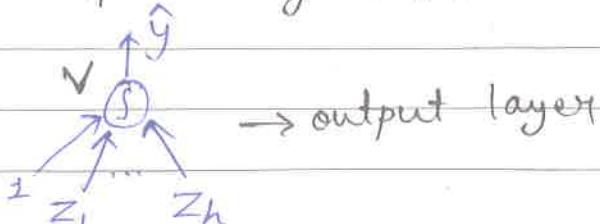
$$\frac{\partial E}{\partial v_j} = \frac{\partial E}{\partial \hat{y}_j} \frac{\partial \hat{y}_j}{\partial v_j} = \frac{1}{N} \times \sum_{i=1}^m (\hat{y}_j^{(i)} - y_j^{(i)}) z_j^{(i)}$$

Chain rule for hidden layer:

$$\begin{aligned} \frac{\partial E}{\partial w_j} &= \sum_{l=1}^k \frac{\partial E}{\partial \hat{y}_l} \frac{\partial \hat{y}_l}{\partial z_j} \frac{\partial z_j}{\partial w_j} \\ &= \sum_{i=1}^m \sum_{l=1}^k (\hat{y}_l^{(i)} - y_l^{(i)}) v_{lj} z_j^{(i)} (1 - z_j^{(i)}) x_j^{(i)} \end{aligned}$$

Q.6} Apply the chain rule to compute the gradients with respect to both hidden layer and output layer parameters of a 2-layer binary classification network where the hidden layers use sigmoid activation and the output layer uses sigmoid activation.

Soln :-



Loss :

$$L(0) = \prod_{i=1}^m P(y=1 | x^{(i)})^{\hat{y}^{(i)}} \underbrace{\left(1 - P(y=1 | x^{(i)})\right)^{1 - \hat{y}^{(i)}}}_{P(y=0 | x^{(i)})}$$

$$l(0) = -\log(L(0))$$

$$= -\sum_{i=1}^m y^{(i)} \log \underbrace{P(y=1/x^{(i)})}_{\hat{y}^{(i)}} + (1-y^{(i)}) \log \underbrace{1-P(y=1/x^{(i)})}_{1-\hat{y}^{(i)}}$$

$$l(0) = -\sum_{i=1}^m y^{(i)} \log (\hat{y}^{(i)}) + (1-y^{(i)}) \log (1-\hat{y}^{(i)})$$

Chain rule for o

Outputs :-

$$\hat{y} = \text{sigmoid}(v^T z)$$

$$z_j = \text{sigmoid}(w_j^T x)$$

Chain rule for output layer :-

$$\frac{\partial l(0)}{\partial v} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial v} =$$

$$= \sum_{i=1}^m \left(y^{(i)} \frac{1}{\hat{y}^{(i)}} - (1-y^{(i)}) \frac{1}{1-\hat{y}^{(i)}} \right) \hat{y}^{(i)} (1-\hat{y}^{(i)}) z^{(i)}$$

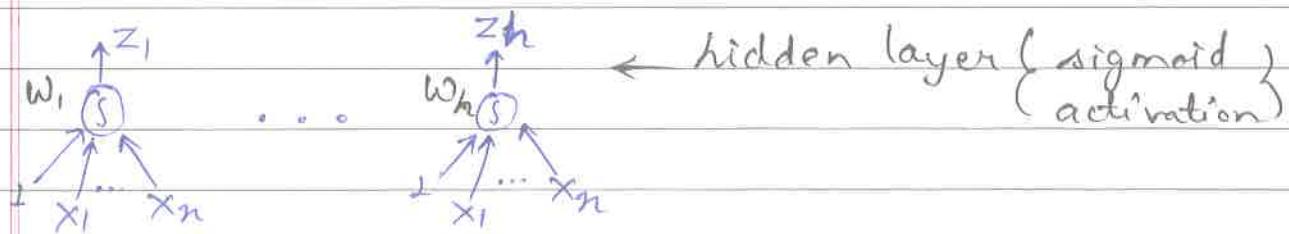
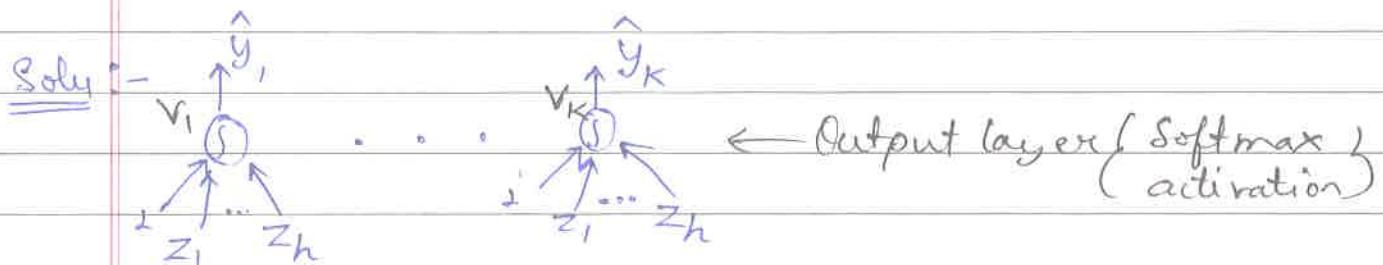
$$= \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) z^{(i)}$$

Chain rule for hidden layer :-

$$\frac{\partial l(0)}{\partial w_j} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_j} \frac{\partial z_j}{\partial w_j}$$

$$\frac{\partial L}{\partial w_j} = \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) v_j z_j^{(i)} (1 - z_j^{(i)}) x^{(i)}$$

Q.7) Apply the chain rule to compute the gradients with respect to both hidden layer and output layer parameters of a 2-layer multiclass classification network where the hidden layers use sigmoid activation and the output layer uses softmax activation.



Loss :-

$$L(\theta) = \prod_{i=1}^m \prod_{j=1}^k P(y=j | X^{(i)})^{I(y^{(i)}=j)}$$

$$L(\theta) = -\log(L(\theta))$$

$$= -\sum_{i=1}^m \sum_{j=1}^k I(y^{(i)}=j) \log(\hat{y}_j^{(i)})$$

$$L(\theta) = -\sum_{i=1}^m \sum_{j=1}^k I(y^{(i)}=j) \log(\hat{y}_j^{(i)})$$

Outputs:-

$$\hat{y}_j = \text{softmax} (v_j^T z)$$

$$z_j = \text{sigmoid} (w_j^T x)$$

Chain rule for output layer:

$$\begin{aligned} \frac{\partial l}{\partial v_j} &= \frac{\partial l}{\partial y_j} \frac{\partial y_j}{\partial v_j} \\ &= -\sum_{i=1}^m (\hat{y}_j^{(i)} - y_j^{(i)}) z_j^{(i)} \end{aligned}$$

Chain rule for hidden layer:

$$\begin{aligned} \frac{\partial l}{\partial w_j} &= \sum_{l=1}^k \frac{\partial l}{\partial g_l} \frac{\partial \hat{y}_l}{\partial z_j} \frac{\partial z_j}{\partial w_j} \\ &= \sum_{i=1}^m \left(\sum_{l=1}^k (\hat{y}_l^{(i)} - y_l^{(i)}) v_{lj} \right) z_j^{(i)} (1 - z_j^{(i)}) x^{(i)} \end{aligned}$$

Q.8) Describe how weights in the network should be initialized and explain why they are initialized in this way.

Soln:- Weights in the network should be initialized randomly from a normal distribution with mean $\mu=0$ and

$$\text{variance } \sigma^2 = \frac{1}{\text{no. of neurons in the previous layer}}$$

Ensuring zero-mean and maintaining the value of the variance of the input of every layer guarantees no exploding / vanishing gradient. Due to this reason, they are initialized in this way.

COMPUTATION GRAPHS

Q.1) Explain the advantage of using computation graphs. Describe what is computed during the forward pass and backward pass of the network during backpropagation. Explain what each node in the graph needs to be able to compute and what information it needs to store (and why).

Solu:- Advantage of using computation graphs is its efficiency in calculating output node values and also computing gradients for a node individually. Due to this property, we can do multiple operations in parallel.

All the intermediate node values is calculated in the forward pass. ~~and gradient~~ gradients is computed and pushes towards beginning from the end of the network in backward pass of the network during backpropagation.

Each node in the graph needs to be able to compute the gradient of the output values with respect to each input to the node. and it needs to store output value of each individual node which is computed in forward pass because this intermediate output value is used for computing gradients in backward pass.

Q.2) Given a two layer network where the hidden layer output is given by $z = f_1(W_1, x)$ and the output layer output is given by $\hat{y} = f_2(W_2, z)$ express the output \hat{y} as a function composition between $f_1()$ and $f_2()$. Assuming the output layer has a single node, express the L_2 loss for a batch $\{x^{(i)}, y^{(i)}\}_{i=1}^m$ in terms of the composite function you wrote and use the chain rule to write its derivatives with respect to W_1 and W_2 . You do not need to compute the actual derivatives of $f_1()$ and $f_2()$ as they are not specified. Repeat the question with cross entropy loss instead of L_2 .

Solu:-

$$\hat{y} = f_2(W_2, f_1(W_1, x))$$

L_2 loss for batch $\{x^{(i)}, y^{(i)}\}_{i=1}^m$:

$$\boxed{\text{L2 loss} = \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2}$$

$$\boxed{L(\text{loss}) = \sum_{i=1}^m (f_2(w_2, f_1(w_1, x^{(i)})) - y^{(i)})^2}$$

$$\frac{\partial L}{\partial w_1} = \cancel{\frac{\partial L}{\partial \hat{y}}} \frac{\partial \hat{y}}{\partial f_1}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f_2} \frac{\partial f_2}{\partial f_1} \frac{\partial f_1}{\partial w_1}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial f_2} \frac{\partial f_2}{\partial w_2}$$

Here,

$$\frac{\partial L}{\partial f_2} = 2(\hat{y} - y)$$

where,

$$\hat{y} = f_2(w_2, f_1(w_1, x))$$

With cross-entropy loss instead of L_2 :

$$\boxed{L(\text{loss}) = - \sum_{i=1}^m y^{(i)} \log(f_2(w_2, f_1(w_1, x^{(i)}))) + (1 - y^{(i)}) \log(1 - f_2(w_2, f_1(w_1, x^{(i)})))}$$

$$\boxed{\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f_2} \frac{\partial f_2}{\partial f_1} \frac{\partial f_1}{\partial w_1} ; \quad \frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial f_2} \frac{\partial f_2}{\partial w_2}}$$

Here,

$$\frac{\partial L}{\partial f_2} = \sum_{i=1}^m \frac{y^{(i)}}{f_2(w_2, f_1(w_1, x^{(i)}))} + \frac{(1-y^{(i)})}{1-f_2(w_2, f_1(w_1, x^{(i)}))}$$

Q.3). Repeat the previous question for a multiple output network.

Soln :-

$$\hat{y}_j = f_2(w_2, f_1(w_1, x))$$

L_2 loss for batch $\{x^{(i)}, y^{(i)}\}_{i=1}^m$:

$$L(\text{loss}) = \sum_{i=1}^m \sum_{l=1}^k (\hat{y}_l^{(i)} - y_l^{(i)})^2$$

$$\text{loss} = \sum_{i=1}^m \sum_{l=1}^k (f_2(w_2, f_1(w_1, x^{(i)})) - y_l^{(i)})^2$$

$$\frac{\partial L}{\partial w_{1j}} = \sum_{l=1}^k \frac{\partial L}{\partial f_{2l}} \frac{\partial f_{2l}}{\partial f_{1j}} \frac{\partial f_{1j}}{\partial w_{1j}}$$

$$\frac{\partial L}{\partial w_{2j}} = \frac{\partial L}{\partial f_{2j}} \frac{\partial f_{2j}}{\partial w_{2j}}$$

Here,

$$\frac{\partial L}{\partial f_{2j}} = 2 \sum_{i=1}^m (f_2(w_2, f_1(w_1, x^{(i)})) - y_j^{(i)})$$

With cross-entropy loss instead of L_2 :-

$$L(\text{loss}) = -\sum_{i=1}^m \sum_{j=1}^k \mathbb{1}(y^{(i)}=j) \log(\hat{y}_j^{(i)})$$

$$= -\sum_{i=1}^m \sum_{j=1}^k \mathbb{1}(y^{(i)}=j) \log(f_{2j}(w_2, f_1(w_1, x^{(i)})))$$

$$\frac{\partial L}{\partial w_{1j}} = \sum_{l=1}^k \frac{\partial L}{\partial f_{2l}} \frac{\partial f_{2l}}{\partial f_{1j}} \frac{\partial f_{1j}}{\partial w_{1j}}$$

$$\frac{\partial L}{\partial w_{2j}} = \frac{\partial L}{\partial f_{2j}} \frac{\partial f_{2j}}{\partial w_{2j}}$$

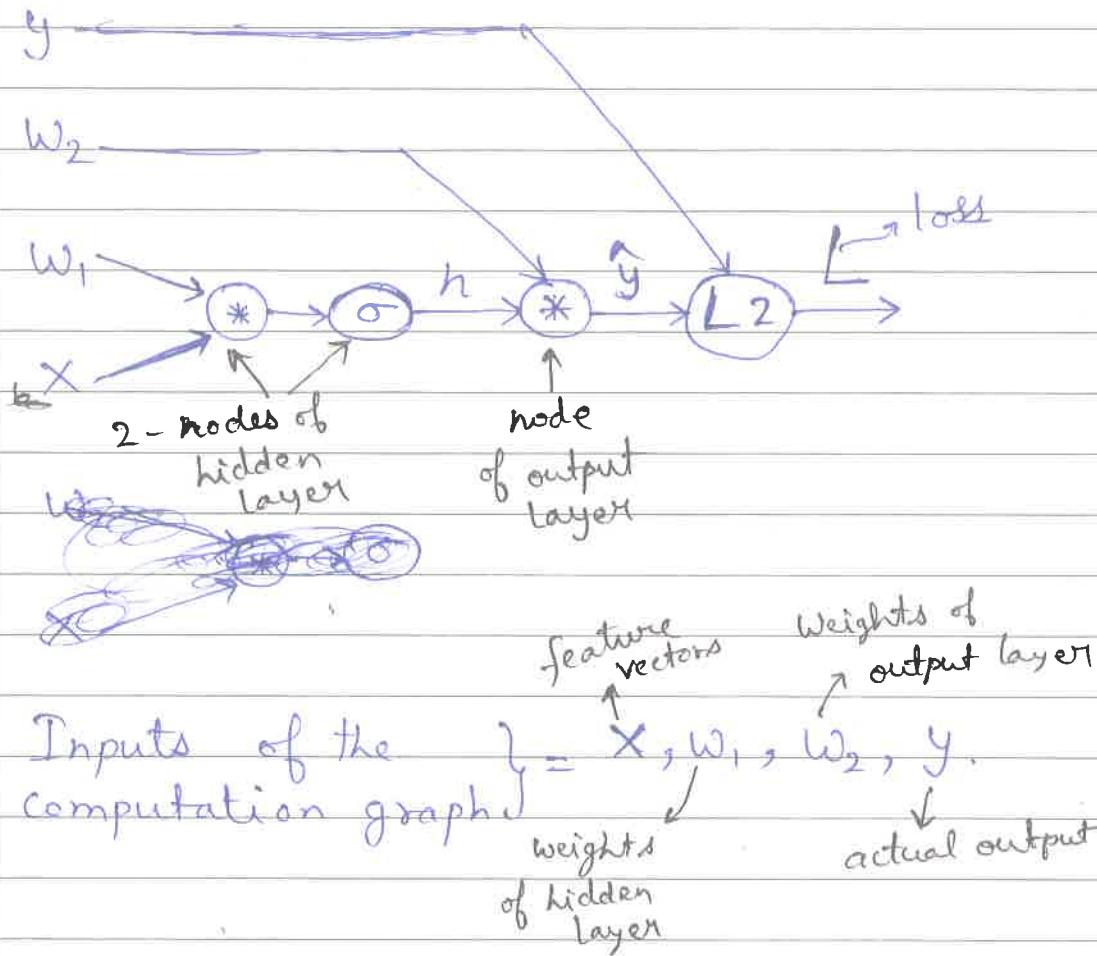
Here,

$$\frac{\partial L}{\partial f_{2j}} = -\sum_{i=1}^m \frac{\mathbb{1}(y^{(i)}=j)}{f_{2j}(w_2, f_1(w_1, x^{(i)}))}$$

- Q.4) Write the computation graph for a two layer regression network with 2 nodes in the hidden layer and 1 node in the output layer. Assume sigmoid activation in the hidden layer and linear activation in the output layer. Assume the feature vectors are 3-dimensional. List all the inputs of the computation graphs and write all the derivatives needed to compute the loss gradient with respect to

network parameters when using the backpropagation pass.

Soln:- Computation graph:



All gradients needed to compute are:

$$L = (\hat{y} - y)^2$$

$$h = \sigma(xw_1)$$

$$\hat{y} = h w_2$$

$$\frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y) \quad \left| \quad \frac{\partial \hat{y}}{\partial w_2} = h \right.$$

$$\frac{\partial \hat{y}}{\partial h} = w_2 \quad \left| \quad \frac{\partial h}{\partial w_1} = h(1-h)x \right.$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_2} = 2(\hat{y} - y) \cdot h$$

$$= 2(\hat{y} - y) \cdot \sigma(xw_1)$$

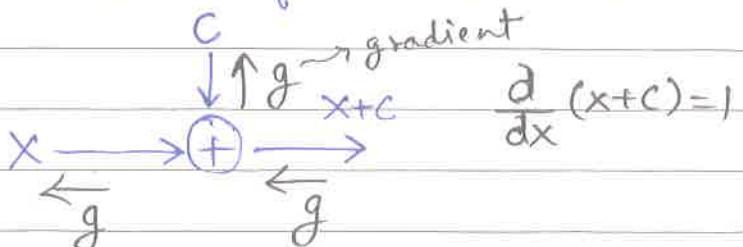
$$\frac{\partial L}{\partial h} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial h} = 2(\hat{y} - y) \cdot w_2$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h} \frac{\partial h}{\partial w_1} = 2(\hat{y} - y) \cdot w_2 \cdot h(1-h) \cdot x$$

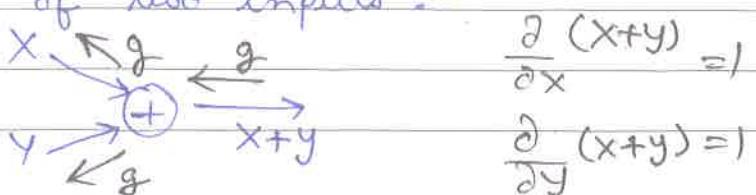
$$= 2(\hat{y} - y) \cdot w_2 \cdot \sigma(xw_1)(1 - \sigma(xw_1)) \cdot x$$

Q.5). Describe the back flow patterns for the following nodes: addition of a constant, addition of two inputs, multiplication by a constant, multiplication of two inputs, max of two inputs, min of two inputs.

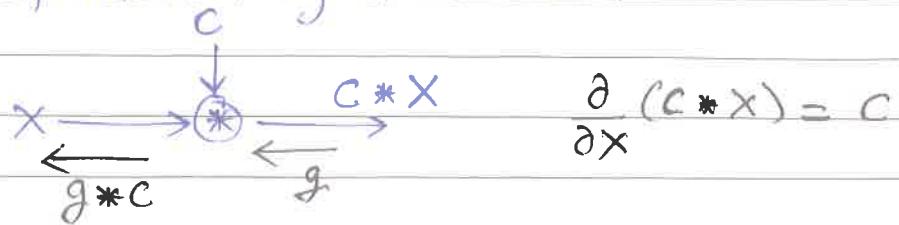
Solu:- • addition of a constant :-



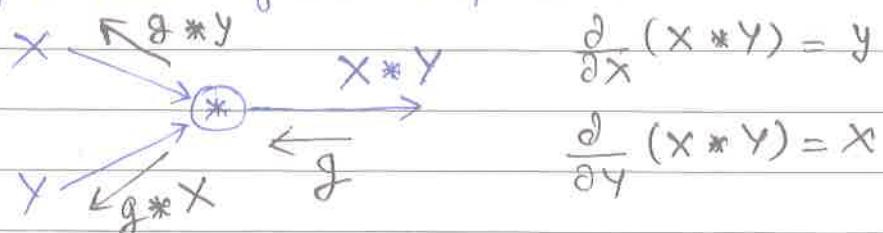
• addition of two inputs :-



• Multiplication by a constant:



• Multiplication of two inputs:

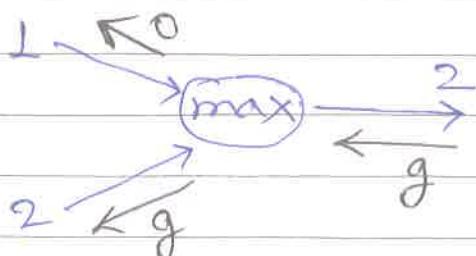


• Max of two inputs:

$$\max(x, y) = \begin{cases} x & \text{if } x \geq y \\ y & \text{otherwise} \end{cases}$$

$$\frac{\partial}{\partial x} \max(x, y) = \begin{cases} 1 & \text{if } x \geq y \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial}{\partial y} \max(x, y) = \begin{cases} 0 & \text{if } x \geq y \\ 1 & \text{otherwise} \end{cases}$$

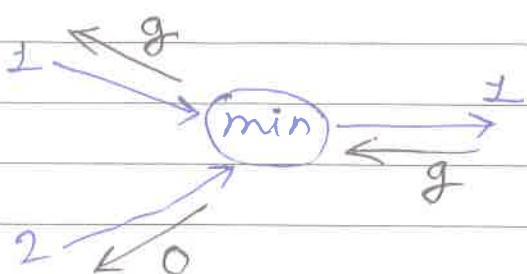


• Min of two inputs:

$$\min(x, y) = \begin{cases} x & \text{if } x \leq y \\ y & \text{otherwise} \end{cases}$$

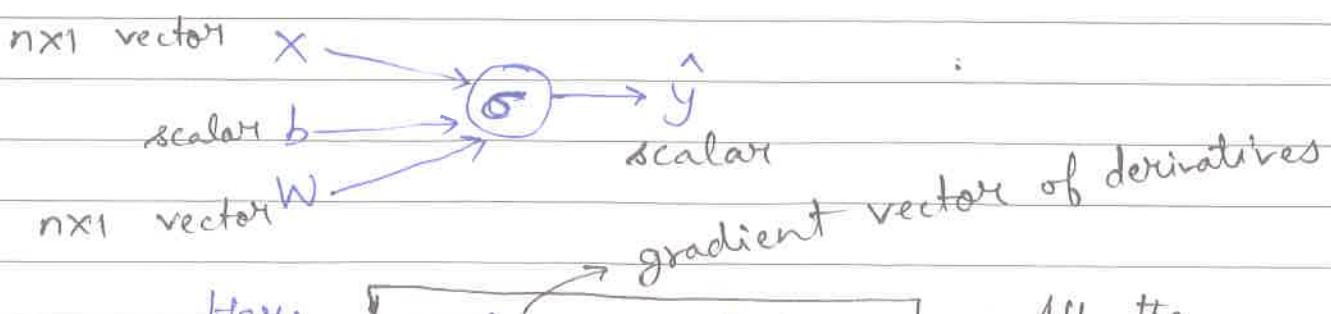
$$\frac{\partial}{\partial x} \min(x, y) = \begin{cases} 1 & \text{if } x \leq y \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial}{\partial y} \min(x, y) = \begin{cases} 0 & \text{if } x \leq y \\ 1 & \text{otherwise} \end{cases}$$



Q.6) Explain what derivative you expect when the input to a node is a vector and the output is a scalar. Write all the components of the derivative.

Solu:- When the input to a node is a vector and the output is a scalar then we will get a rank 1 tensor or gradient vector.



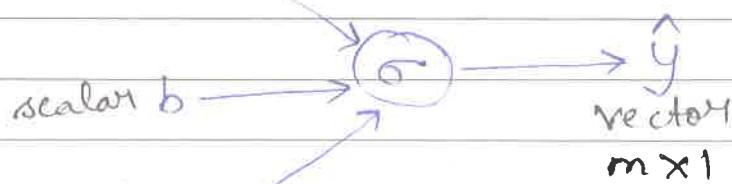
$$\frac{\partial \hat{y}}{\partial w} = \begin{bmatrix} \frac{\partial \hat{y}}{\partial w_1} \\ \vdots \\ \frac{\partial \hat{y}}{\partial w_n} \end{bmatrix}.$$

All the components of the derivative

Q.7) Explain what derivative you expect when the input to a node is a vector and the output is a vector. Write all the components of the derivative.

Solu :- When the input to a node is a vector and the output is a vector then we will get rank 2 tensor of derivatives.

$n \times 1$ vector X



$n \times 1$ vector W

$$W = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$$

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_m \end{bmatrix}$$

All the

components
of the
derivative.

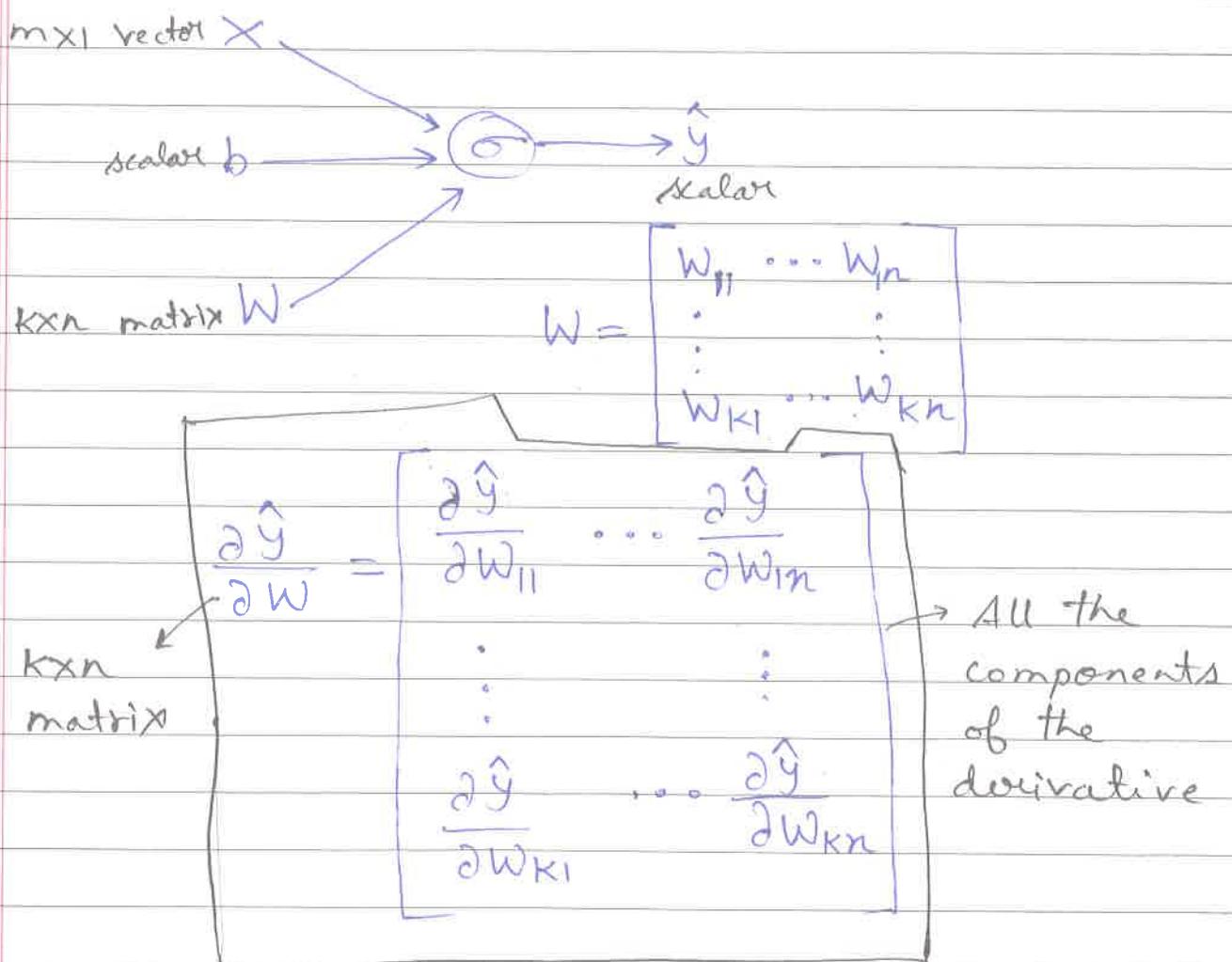
$$\frac{\partial \hat{y}}{\partial W} = \begin{bmatrix} \frac{\partial \hat{y}_1}{\partial w_1} & \dots & \frac{\partial \hat{y}_1}{\partial w_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial \hat{y}_m}{\partial w_1} & \dots & \frac{\partial \hat{y}_m}{\partial w_n} \end{bmatrix}_{(m \times n)}$$

It is a
rank 2
tensor.

Shape of this
matrix is $m \times n$.

Q.8). Explain what derivative you expect when the input to a node is a matrix and the output is a scalar. Write all the components of the derivative.

Soln:- When the input to a node is a matrix and the output is a scalar then we will get rank 2 tensor of derivatives.



Q.9). Write the chain rule for the composition of vector valued vector functions F and G . Pay attention to the order of the Jacobian matrices you write.

Solu:- Let $F: R^k \rightarrow R^m$ and $G: R^h \rightarrow R^k$ be vector valued functions.

→ The composition of F and G is the function:

$F \circ G: R^h \rightarrow R^m$ given by:

$$(F \circ G)(P) = F(G(P))$$

$\uparrow \quad \uparrow \quad \uparrow$
 $m \times 1 \quad k \times 1 \quad n \times 1$

Chain rule for the composition of vector valued functions:

$$D(F \circ G)(P) = \underbrace{D^m F(G(P))}_{\substack{\text{compute } DF \\ \text{then assign} \\ G(P) \text{ instead of } P}} \cdot \underbrace{D^k G(P)}_{\substack{\text{Product of matrices} \\ \text{rank 2} \\ \text{tensor}}}$$

$$\begin{matrix} m \times k \\ \text{rank 2} \\ \text{tensor} \end{matrix} * \begin{matrix} k \times n \\ \text{rank 2} \\ \text{tensor} \end{matrix} = \begin{matrix} m \times n \\ \text{rank 2} \\ \text{tensor} \end{matrix}$$

Example:-

$$F(x, y, z) = \begin{bmatrix} 2xy \\ y-z \end{bmatrix} \quad G(x, y) = \begin{bmatrix} x-4y \\ xy \\ x+y \end{bmatrix}$$

$$DF(x, y, z) = \begin{bmatrix} \frac{\partial}{\partial x}(2xy) & \frac{\partial}{\partial y}(2xy) & \frac{\partial}{\partial z}(2xy) \\ \frac{\partial}{\partial x}(y-z) & \frac{\partial}{\partial y}(y-z) & \frac{\partial}{\partial z}(y-z) \end{bmatrix} = \begin{bmatrix} 2y & 2x & 0 \\ 0 & 1 & -1 \end{bmatrix}$$

Using the chain rule

→ first compute DF and then plug in $G(x,y)$

$$DF(G(x,y)) = \begin{bmatrix} 2(xy) & 2(x-4y) & 0 \\ 0 & 1 & -1 \end{bmatrix}$$

$$DG(x,y) = \begin{bmatrix} \frac{\partial}{\partial x} (x-4y) & \frac{\partial}{\partial y} (x-4y) \\ \frac{\partial}{\partial x} (xy) & \frac{\partial}{\partial y} (xy) \\ \frac{\partial}{\partial x} (x+y) & \frac{\partial}{\partial y} (x+y) \end{bmatrix} = \begin{bmatrix} 1 & -4 \\ y & x \\ 1 & 1 \end{bmatrix}$$

Using the chain rule:

$$D(F \circ G)(x,y) = DF(G(x,y)) * DG(x,y)$$

$$= \begin{bmatrix} 2xy & 2x-8y & 0 \\ 0 & 1 & -1 \end{bmatrix}_{2 \times 3} \begin{bmatrix} 1 & -4 \\ y & x \\ 1 & 1 \end{bmatrix}_{3 \times 2}$$

$$= \begin{bmatrix} 2xy + 2xy - 8y^2 & -8xy + 2x^2 - 8xy \\ y-1 & x-1 \end{bmatrix}$$

$$= \begin{bmatrix} 4xy - 8y^2 & 2x^2 - 16xy \\ y-1 & x-1 \end{bmatrix}_{2 \times 2}$$

Q.10) Let $F(x, y, z) = [3xy \quad y-z]^T$ and

$G(x, y) = [x-5y \quad xy \quad x-y]^T$ write the

composition $(f \circ G)(x, y)$. Compute the Jacobian matrix of the composition in two ways: directly and by using the chain rule. Make sure the result you compute in both ways is identical.

Soln :-

$$F(x, y, z) = \begin{bmatrix} 3xy \\ y-z \end{bmatrix} \quad G(x, y) = \begin{bmatrix} x-5y \\ xy \\ x-y \end{bmatrix}$$

$$(F \circ G)(x, y) = \begin{bmatrix} 3(x-5y)(xy) \\ (xy) - (x-y) \end{bmatrix} = \begin{bmatrix} 3x^2y - 15xy^2 \\ xy - x + y \end{bmatrix}$$

Computing the Jacobian matrix directly :

$$D(F \circ G)(x, y) = \begin{bmatrix} \frac{\partial}{\partial x}(3x^2y - 15xy^2) & \frac{\partial}{\partial y}(3x^2y - 15xy^2) \\ \frac{\partial}{\partial x}(xy - x + y) & \frac{\partial}{\partial y}(xy - x + y) \end{bmatrix}$$

$$D(F \circ G)(x, y) = \begin{bmatrix} 6xy - 15y^2 & 3x^2 - 30xy \\ y-1+0 & x-0+1 \end{bmatrix}$$

$$= \begin{bmatrix} 6xy - 15y^2 & 3x^2 - 30xy \\ y-1 & x+1 \end{bmatrix} \quad \text{Ans.}$$

Computing the Jacobian matrix using chain rule:

$$DF(x, y, z) = \begin{bmatrix} \frac{\partial}{\partial x}(3xy) & \frac{\partial}{\partial y}(3xy) & \frac{\partial}{\partial z}(3xy) \\ \frac{\partial}{\partial x}(y-z) & \frac{\partial}{\partial y}(y-z) & \frac{\partial}{\partial z}(y-z) \end{bmatrix}$$

$$DF(x, y, z) = \begin{bmatrix} 3y & 3x & 0 \\ 0 & 1 & -1 \end{bmatrix}$$

$$DG(x, y) = \begin{bmatrix} \frac{\partial}{\partial x}(x-5y) & \frac{\partial}{\partial y}(x-5y) \\ \frac{\partial}{\partial x}(xy) & \frac{\partial}{\partial y}(xy) \\ \frac{\partial}{\partial x}(x-y) & \frac{\partial}{\partial y}(x-y) \end{bmatrix}$$

$$DG(x, y) = \begin{bmatrix} 1 & -5 \\ y & x \\ 1 & -1 \end{bmatrix}$$

Using the chain rule:-

$$DF(G(x, y)) = \begin{bmatrix} 3(xy) & 3(x-5y) & 0 \\ 0 & 1 & -1 \end{bmatrix}$$

Compute DF

then plug in

$$G(x, y)$$

$$= \begin{bmatrix} 3xy & 3x-15y & 0 \\ 0 & 1 & -1 \end{bmatrix}$$

Using the chain rule :-

$$D(F \circ G)(x, y) = DF(G(x, y)) \cdot DG(x, y)$$

$$= \begin{bmatrix} 3xy & 3x-15y & 0 \\ 0 & 1 & -1 \end{bmatrix} \begin{matrix} 1 & -5 \\ y & x \\ 1 & -1 \end{matrix} \begin{matrix} 2 \times 3 \\ 3 \times 2 \end{matrix}$$

$$= \begin{bmatrix} 3xy + 3xy - 15y^2 & -15xy + 3x^2 - 15xy \\ y-1 & x+1 \end{bmatrix}$$

$$= \begin{bmatrix} 6xy - 15y^2 & 3x^2 - 30xy \\ y-1 & x+1 \end{bmatrix} \begin{matrix} \text{Ans} \\ \underline{\underline{=}} \\ 2 \times 2 \end{matrix}$$

Q.11) Explain the need to order nodes when traversing the computation graph in forward and backward passes.

Solu :- Here, I need to compute the nodes preceding me before I move to the next node.

Due to this reason, it is needed to order nodes when traversing the computation graph in forward and backward passes.

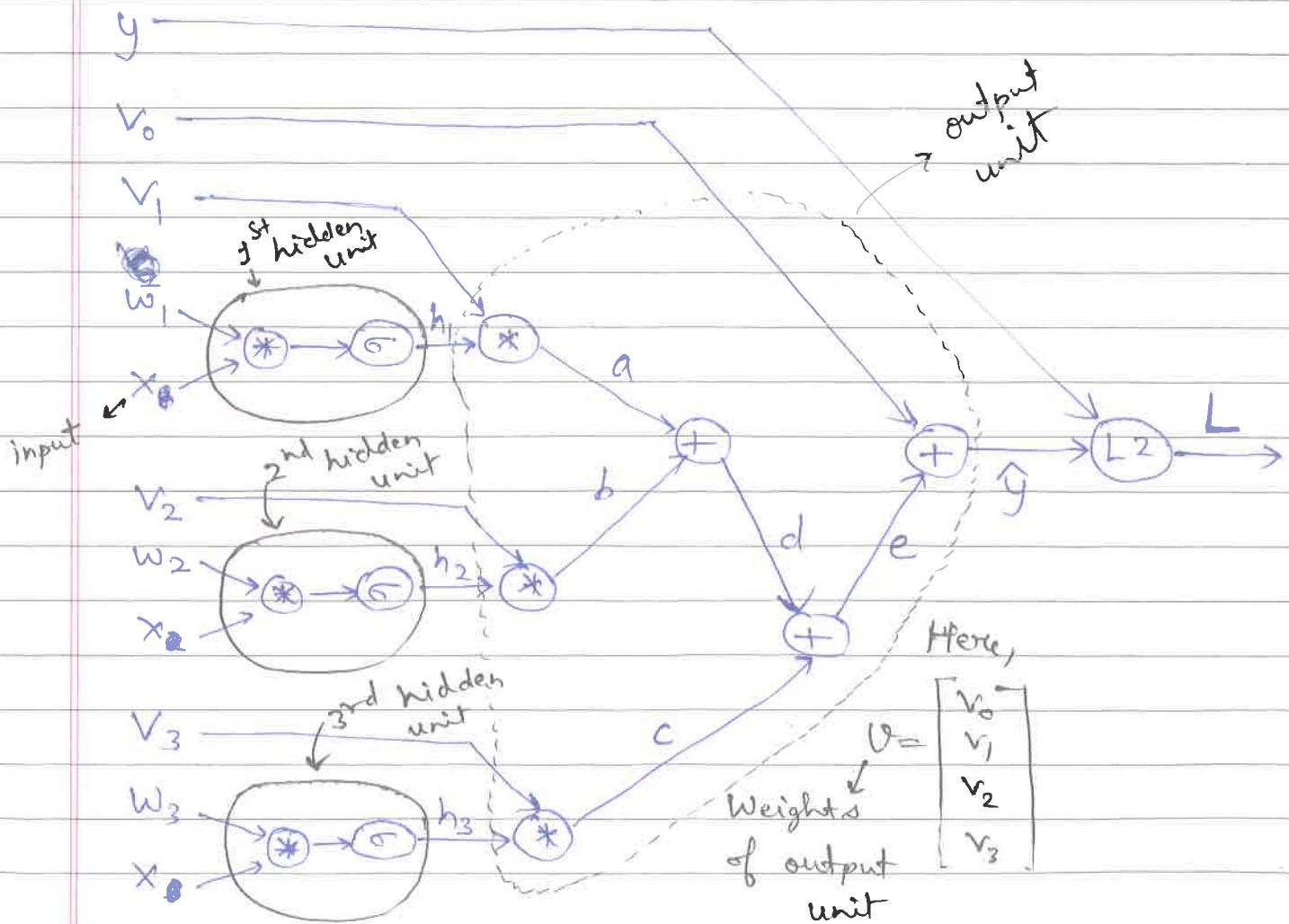
Q.12) Scalar computation graph:

(a). Draw the computation graph for a regression network with three hidden units and one output unit. Assume logistic activation at hidden layer nodes and linear activation at output. Assume L2 loss.

Soly:-

$$L = L_2 \left(V_0 + \sigma(xw_1) V_1 + \sigma(xw_2) V_2 + \sigma(xw_3) V_3 \rightarrow y \right)$$

$$L = \left((v_0 + \epsilon(xw_1)v_1 + \epsilon(xw_2)v_2 + \epsilon(xw_3)v_3) - y \right)^2.$$



(b). Assume the current weight parameters guess is given by:

$$w_1 = (0.01, 0.02, 0.03)$$

$$w_2 = (0.03, 0.01, 0.02)$$

$$w_3 = (0.02, 0.03, 0.01)$$

$$v = (0.01, 0.02, 0.03, 0.04)$$

where, w_1, w_2, w_3 are hidden unit weight vectors and v is the output unit weight vector. Assume the training data in the batch is given by:

$$x_1, y_1 = [(1, 2), 8]$$

$$x_2, y_2 = [(1, 3), 11]$$

$$x_3, y_3 = [(2, 2), 10]$$

Compute the values at the network for the forward pass for each data point.

Solu :- ~~for~~ For $x_1, y_1 = [(1, 2), 8]$

$$h_1 = \sigma(w_1 x_1)$$

$$= \sigma(0.01 + 0.02 \times 1 + 0.03 \times 2)$$

$$= \sigma(0.09)$$

$$= \underline{0.52248}$$

$$\begin{aligned}
 h_2 &= \sigma(w_2 x_1) \\
 &= \sigma(0.03 + 0.01x_1 + 0.02x_2) \\
 &= \sigma(0.08) \\
 &= \underline{0.5199} .
 \end{aligned}$$

$$\begin{aligned}
 h_3 &= \sigma(w_3 x_1) \\
 &= \sigma(0.02 + 0.03x_1 + 0.01x_2) \\
 &= \sigma(0.07) \\
 &= \underline{0.51749} .
 \end{aligned}$$

$$\begin{aligned}
 a &= v_1 \times h_1 \\
 &= 0.02 \times 0.52248 \\
 &= \underline{0.01049} .
 \end{aligned}$$

$$\begin{aligned}
 b &= v_2 \times h_2 \\
 &= 0.03 \times 0.5199 \\
 &= \underline{0.01559} .
 \end{aligned}$$

$$\begin{aligned}
 c &= v_3 \times h_3 \\
 &= 0.04 \times 0.51749 \\
 &= \underline{0.02069} .
 \end{aligned}$$

$$\begin{aligned}
 d &= a + b \\
 &= 0.01049 + 0.01559 \\
 &= \underline{0.02608} .
 \end{aligned}$$

$$\begin{aligned}
 e &= d + c \\
 &= 0.02069 + 0.02608 \\
 &= \underline{0.04677} .
 \end{aligned}$$

$$\begin{aligned}
 \hat{y} &= v_0 + e \\
 &= 0.01 + 0.04677 \\
 &= \underline{0.05677} .
 \end{aligned}$$

$$\begin{aligned}
 L &= (\hat{y} - y)^2 \\
 &= (0.05677 - 8)^2 \\
 &= \underline{\underline{63.09}}
 \end{aligned}$$

$$\Rightarrow \text{For } x_2, y_2 = [(1, 3), 11]$$

$$\begin{aligned}
 h_1 &= \sigma(w_1 x_1) \\
 &= \sigma(0.01 + 0.02 \times 1 + 0.03 \times 3) \\
 &= \sigma(0.12) \\
 &= \underline{\underline{0.52996}}
 \end{aligned}$$

$$\begin{aligned}
 h_2 &= \sigma(w_2 x_2) \\
 &= \sigma(0.03 + 0.01 \times 1 + 0.02 \times 3) \\
 &= \sigma(0.1) \\
 &= \underline{\underline{0.52497}}
 \end{aligned}$$

$$\begin{aligned}
 h_3 &= \sigma(w_3 x_3) \\
 &= \sigma(0.02 + 0.03 \times 1 + 0.01 \times 3) \\
 &= \sigma(0.08) \\
 &= \underline{\underline{0.5199}}
 \end{aligned}$$

$$\begin{aligned}
 a &= v_1 x h_1 \\
 &= 0.02 \times 0.52996 \\
 &= \underline{\underline{0.0106}}
 \end{aligned}$$

$$\begin{aligned}
 b &= v_2 x h_2 \\
 &= 0.03 \times 0.52497 \\
 &= \underline{\underline{0.0157491}}
 \end{aligned}$$

$$\begin{aligned}
 c &= v_3 \times h_3 \\
 &= 0.04 \times 0.5199 \\
 &= \underline{\underline{0.020796}}
 \end{aligned}$$

$$\begin{aligned}
 d &= a+b \\
 &= 0.0106 + 0.0157491 \\
 &= \underline{\underline{0.0263}}
 \end{aligned}$$

$$\begin{aligned}
 e &= d+c \\
 &= 0.0263 + 0.0208 \\
 &= \underline{\underline{0.0471}}
 \end{aligned}$$

$$\begin{aligned}
 \hat{y} &= v_0 + e \\
 &= 0.01 + 0.0471 \\
 &= \underline{\underline{0.0571}}
 \end{aligned}$$

$$\begin{aligned}
 L &= (\hat{y} - y)^2 \\
 &= (0.0571 - 11)^2 \\
 &= \underline{\underline{119.7}}
 \end{aligned}$$

For $x_3, y_3 = [(2, 2), 10]$

$$\begin{aligned}
 h_1 &= \sigma(w_1 x_3) \\
 &= \sigma(0.01 + 0.02 \times 2 + 0.03 \times 2) \\
 &= \sigma(0.11) \\
 &= \underline{\underline{0.5275}}
 \end{aligned}$$

$$\begin{aligned}
 h_2 &= \sigma(w_2 x_3) \\
 &= \sigma(0.03 + 0.01 \times 2 + 0.02 \times 2) \\
 &= \sigma(0.09) \\
 &= \underline{\underline{0.5225}}
 \end{aligned}$$

$$\begin{aligned}
 h_3 &= \sigma(w_3 x_3) \\
 &= \sigma(0.02 + 0.03 \times 2 + 0.01 \times 2) \\
 &= \sigma(0.1) \\
 &= \underline{0.525}
 \end{aligned}$$

$$\begin{aligned}
 a &= v_1 \times h_1 \\
 &= 0.02 \times 0.525 \\
 &= \underline{0.01055}
 \end{aligned}$$

$$\begin{aligned}
 b &= v_2 \times h_2 \\
 &= 0.03 \times 0.525 \\
 &= \underline{0.01568}
 \end{aligned}$$

$$\begin{aligned}
 c &= v_3 \times h_3 \\
 &= 0.04 \times 0.525 \\
 &= \underline{0.021}
 \end{aligned}$$

$$\begin{aligned}
 d &= a + b \\
 &= 0.01055 + 0.01568 \\
 &= \underline{0.02623}
 \end{aligned}$$

$$\begin{aligned}
 e &= d + c \\
 &= 0.02623 + 0.021 \\
 &= \underline{0.04723}
 \end{aligned}$$

$$\begin{aligned}
 \hat{y} &= v_0 + e \\
 &= 0.01 + 0.04723 \\
 &= \underline{0.05723}
 \end{aligned}$$

$$\begin{aligned}
 L &= (\hat{y} - y)^2 \\
 &= (0.05723 - 10)^2 \\
 &= \underline{98.8}
 \end{aligned}$$

(c). Compute the gradients using the formula of the loss function with respect to w_1, w_2, w_3, v due to all training examples using back-propagation.

Solu: ~~For $x_1, y_1 = [(1, 2), 8]$~~

$$\frac{\partial L}{\partial \hat{y}} = \frac{\partial}{\partial \hat{y}} (L) \quad \left| \begin{array}{l} \frac{\partial e}{\partial d} = 1 \\ \frac{\partial d}{\partial b} = 1 \end{array} \right.$$

$$\frac{\partial \hat{y}}{\partial e} = 1 \quad \left| \begin{array}{l} \frac{\partial d}{\partial a} = 1 \\ \frac{\partial e}{\partial c} = 1 \end{array} \right.$$

$$\left| \begin{array}{l} \frac{\partial a}{\partial h_1} = v_1 \\ \frac{\partial b}{\partial h_2} = v_2 \\ \frac{\partial c}{\partial h_3} = v_3 \end{array} \right.$$

$$\frac{\partial h_1}{\partial w_1} = \sigma(w_1 x) (1 - \sigma(w_1 x)) x$$

$$\frac{\partial h_2}{\partial w_2} = \sigma(w_2 x) (1 - \sigma(w_2 x)) x$$

$$\frac{\partial h_3}{\partial w_3} = \sigma(w_3 x) (1 - \sigma(w_3 x)) x$$

$$\text{So, } \frac{\partial \hat{y}}{\partial v_0} = 1 \quad \left| \begin{array}{l} \frac{\partial a}{\partial v_1} = h_1 = \sigma(w_1 x) \\ \frac{\partial b}{\partial v_2} = h_2 = \sigma(w_2 x) \\ \frac{\partial c}{\partial v_3} = h_3 = \sigma(w_3 x) \end{array} \right.$$

$$\frac{\partial b}{\partial v_2} = h_2 = \sigma(w_2 x)$$

$$\frac{\partial c}{\partial v_3} = h_3 = \sigma(w_3 x)$$

$$\frac{\partial L}{\partial v_0} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial v_0} = 2(\hat{y} - y) * 1$$

$$\frac{\partial L}{\partial v_1} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial e} \frac{\partial e}{\partial d} \frac{\partial d}{\partial a} \frac{\partial a}{\partial v_1},$$

$$= 2(\hat{y} - y) * 1 * 1 * \sigma(w, x)$$

$$= 2(\hat{y} - y) \sigma(w, x)$$

$$\frac{\partial L}{\partial v_2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial e} \frac{\partial e}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial v_2}$$

$$= 2(\hat{y} - y) * 1 * 1 * \sigma(w_2, x)$$

$$= 2(\hat{y} - y) \sigma(w_2, x)$$

$$\frac{\partial L}{\partial v_3} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial e} \frac{\partial e}{\partial c} \frac{\partial c}{\partial v_3}$$

$$= 2(\hat{y} - y) * 1 * 1 * \sigma(w_3, x)$$

$$= 2(\hat{y} - y) \sigma(w_3, x).$$

$$\frac{\partial L}{\partial v} = \begin{bmatrix} \frac{\partial L}{\partial v_0} \\ \frac{\partial L}{\partial v_1} \\ \frac{\partial L}{\partial v_2} \\ \frac{\partial L}{\partial v_3} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \\ \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \sigma(w, x) \\ \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \sigma(w_2, x) \\ \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \sigma(w_3, x) \end{bmatrix} \stackrel{\text{Ans}}{=} .$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial g} \frac{\partial \hat{y}}{\partial e} \frac{\partial e}{\partial d} \frac{\partial d}{\partial a} \frac{\partial a}{\partial h_1} \frac{\partial h_1}{\partial w_1}$$

$$= 2(\hat{y} - y) \cdot 1 \cdot 1 \cdot 1 \cdot v_1 \cdot \sigma(w_1 x) (1 - \sigma(w_1 x)) x$$

$$\boxed{\frac{\partial L}{\partial w_1} = \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \cdot v_1 \cdot \sigma(w_1 x^{(i)}) (1 - \sigma(w_1 x^{(i)})) x^{(i)}} \quad \text{Ans}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial g} \frac{\partial \hat{y}}{\partial e} \frac{\partial e}{\partial d} \frac{\partial d}{\partial b} \frac{\partial b}{\partial h_2} \frac{\partial h_2}{\partial w_2}$$

$$= 2(\hat{y} - y) \cdot 1 \cdot 1 \cdot 1 \cdot v_2 \cdot \sigma(w_2 x) (1 - \sigma(w_2 x)) x$$

$$\boxed{\frac{\partial L}{\partial w_2} = \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \cdot v_2 \cdot \sigma(w_2 x^{(i)}) (1 - \sigma(w_2 x^{(i)})) x^{(i)}} \quad \text{Ans}$$

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial e} \frac{\partial e}{\partial c} \frac{\partial c}{\partial h_3} \frac{\partial h_3}{\partial w_3}$$

$$= 2(\hat{y} - y) \cdot 1 \cdot 1 \cdot v_3 \cdot \sigma(w_3 x) (1 - \sigma(w_3 x)) x$$

$$\boxed{\frac{\partial L}{\partial w_3} = \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \cdot v_3 \cdot \sigma(w_3 x^{(i)}) (1 - \sigma(w_3 x^{(i)})) x^{(i)}} \quad \text{Ans}$$

$$\boxed{\frac{\partial L}{\partial v} = 2 \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) z^{(i)}} \quad \text{Here, } z^{(i)} =$$

$$\begin{bmatrix} 1 \\ \sigma(w_1 x^{(i)}) \\ \sigma(w_2 x^{(i)}) \\ \sigma(w_3 x^{(i)}) \end{bmatrix}$$

Q(d). Compute the gradients using the formula given in class and compare to your previous results.

$$\text{Solu: } \frac{\partial L}{\partial v} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial v}$$

$$\boxed{\frac{\partial L}{\partial v} = 2 \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) z^{(i)}}$$

$$\frac{\partial L}{\partial w_j} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_j} \frac{\partial z_j}{\partial w_j}$$

no. of units in hidden layer.

$$= 2 \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) v_j \cdot z_j^{(i)} (1 - z_j^{(i)}) x^{(i)}$$

So,

$$\boxed{\frac{\partial L}{\partial w_1} = 2 \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) v_1 \cdot z_1^{(i)} (1 - z_1^{(i)}) x^{(i)}}$$

$$\boxed{\frac{\partial L}{\partial w_2} = 2 \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) v_2 \cdot z_2^{(i)} (1 - z_2^{(i)}) x^{(i)}}$$

$$\boxed{\frac{\partial L}{\partial w_3} = 2 \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) v_3 \cdot z_3^{(i)} (1 - z_3^{(i)}) x^{(i)}}$$

Here, $z_1^{(i)} = \sigma(w_1 x^{(i)})$; $z_2^{(i)} = \sigma(w_2 x^{(i)})$

$$z_3^{(i)} = \leftarrow (w_3, x^{(i)})$$

On comparing with previous results, we came to know that both are same.

Q.13). Vector computation graph :-

(a). Let $f(x, y) = (2x+3y)^2$. Compute $\nabla f(x, y)$.

Soln:- $f(x, y) = (2x+3y)^2$

$$\nabla f(x, y) = \begin{bmatrix} \frac{\partial f(x, y)}{\partial x} \\ \frac{\partial f(x, y)}{\partial y} \end{bmatrix} = \begin{bmatrix} 2(2x+3y) \cdot 2 \\ 2(2x+3y) \cdot 3 \end{bmatrix}$$

$$\nabla f(x, y) = \begin{bmatrix} 8x+12y \\ 12x+18y \end{bmatrix} \quad \underline{\text{Ans}}$$

(b). Let $F(x, y) = \begin{bmatrix} x^2+2y \\ 3x+4y^2 \end{bmatrix}$. Compute the Jacobian

matrix $\nabla F(2, 2)$.

Soln:- $F(x, y) = \begin{bmatrix} x^2+2y \\ 3x+4y^2 \end{bmatrix}$

$$DF(x, y) = \begin{bmatrix} \frac{\partial}{\partial x} (x^2 + 2y) & \frac{\partial}{\partial y} (x^2 + 2y) \\ \frac{\partial}{\partial x} (3x + 4y^2) & \frac{\partial}{\partial y} (3x + 4y^2) \end{bmatrix}$$

$$DF(x, y) = \begin{bmatrix} 2x & 2 \\ 3 & 8y \end{bmatrix}.$$

Now,

$$DF(1, 2) = \begin{bmatrix} 2*1 & 2 \\ 3 & 8*2 \end{bmatrix}$$

$$DF(1, 2) = \begin{bmatrix} 2 & 2 \\ 3 & 16 \end{bmatrix} \quad \underline{\underline{\text{Ans}}}.$$

(c). Let $G_1(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$ and $F(x, y)$ as in the

previous question. Compute the Jacobian matrix $D(F \circ G_1)(2)$ both with and without the chain rule.

Solu :- $G_1(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}; F(x, y) = \begin{bmatrix} x^2 + 2y \\ 3x + 4y^2 \end{bmatrix}$

$$(F \circ G_1)(x) = \begin{bmatrix} x^2 + 2(x^2) \\ 3x + 4(x^2)^2 \end{bmatrix} = \begin{bmatrix} 3x^2 \\ 3x + 4x^4 \end{bmatrix}.$$

Computing the Jacobian matrix directly :-

$$D(F_0 G)(x) = \begin{bmatrix} \frac{\partial}{\partial x} (3x^2) \\ \frac{\partial}{\partial x} (3x + 4x^4) \end{bmatrix} = \begin{bmatrix} 6x \\ 3 + 16x^3 \end{bmatrix}$$

$$D(F_0 G)(2) = \begin{bmatrix} 6 \cdot 2 \\ 3 + 16(2)^3 \end{bmatrix} = \begin{bmatrix} 12 \\ 3 + 128 \end{bmatrix}$$

$$D(F_0 G)(2) = \begin{bmatrix} 12 \\ 131 \end{bmatrix} \stackrel{\text{Ans}}{=}$$

Computing the Jacobian matrix using chain rule :-

$$DF(x, y) = \begin{bmatrix} \frac{\partial}{\partial x} (x^2 + 2y) & \frac{\partial}{\partial y} (x^2 + 2y) \\ \frac{\partial}{\partial x} (3x + 4y^2) & \frac{\partial}{\partial y} (3x + 4y^2) \end{bmatrix}$$

$$DF(x, y) = \begin{bmatrix} 2x & 2 \\ 3 & 8y \end{bmatrix}$$

$$DG(x) = \begin{bmatrix} \frac{\partial}{\partial x} (x) \\ \frac{\partial}{\partial x} (x^2) \end{bmatrix} = \begin{bmatrix} 1 \\ 2x \end{bmatrix}$$

$$DF(G(x)) = \begin{bmatrix} 2(x) & 2 \\ 3 & 8(x^2) \end{bmatrix} = \begin{bmatrix} 2x & 2 \\ 3 & 8x^2 \end{bmatrix}$$

↓
first compute DF
and then plug in $G(x)$

Using the chain rule:-

$$D(F \circ G)(x) = DF(G(x)) * DG(x)$$

$$= \begin{bmatrix} 2x & 2 \\ 3 & 8x^2 \end{bmatrix}_{2 \times 2} * \begin{bmatrix} 1 \\ 2x \end{bmatrix}_{2 \times 1}$$

$$= \begin{bmatrix} 2x + 4x \\ 3 + 16x^3 \end{bmatrix}$$

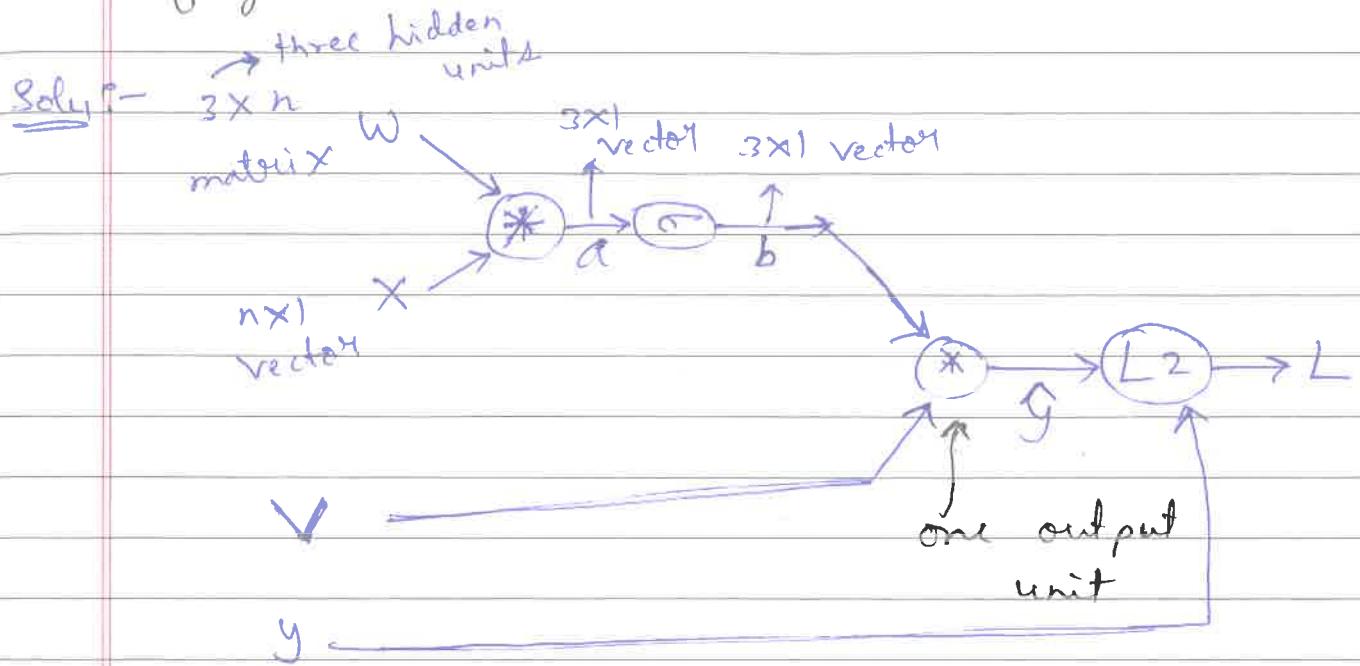
$$= \begin{bmatrix} 6x \\ 3 + 16x^3 \end{bmatrix}$$

$$D(F \circ G)(2) = \begin{bmatrix} 6 * 2 \\ 3 + 16(2)^3 \end{bmatrix} = \begin{bmatrix} 12 \\ 3 + 128 \end{bmatrix}$$

$$D(F \circ G)(2) = \begin{bmatrix} 12 \\ 131 \end{bmatrix} \text{ Ans.}$$

- (d). Repeat the scalar computation graph question (question 1) but this time using a vector notation. That is, draw the computation graph using vectors for input data and weights, perform a forward pass using vectors, and perform back propagation by creating

and multiplying Jacobian matrices instead of gradient vectors.



forward pass:

$$a = Wx$$

$$b = \sigma(Wx) \quad \sigma(a) = \sigma(Wx)$$

$$\hat{y} = Vb$$

$$L = L_2(\hat{y}, y)$$

$$L = (\hat{y} - y)^2$$

backward pass:

$$\frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)$$

$$\frac{\partial \hat{y}}{\partial b} = V$$

$$\frac{\partial b}{\partial a} = \sigma'(a)(1 - \sigma(a))$$

$$\frac{\partial a}{\partial W} = x$$

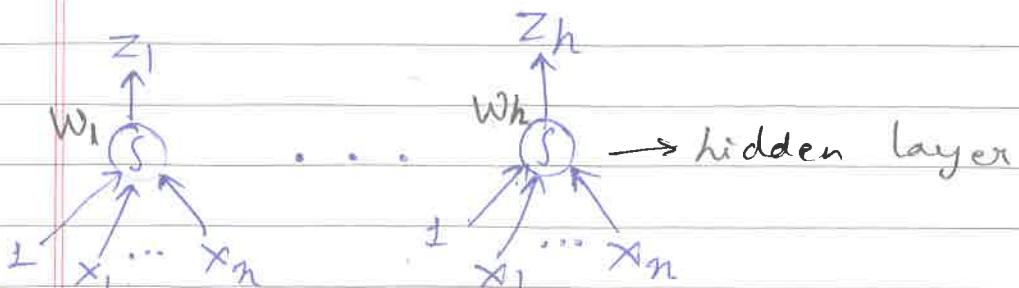
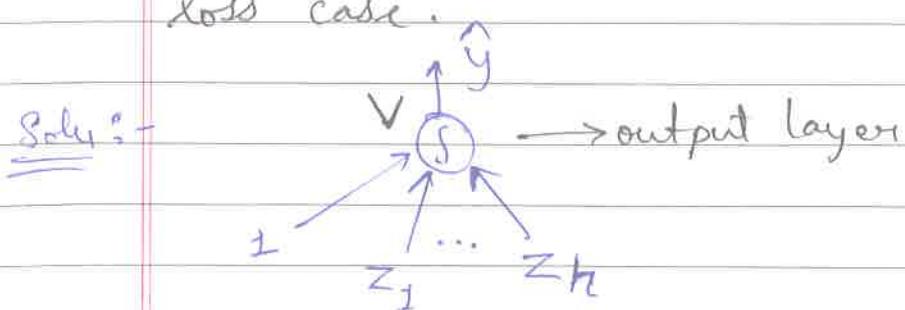
$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial b} * \frac{\partial b}{\partial a} * \frac{\partial a}{\partial w}$$

$$= 2(\hat{y} - y) * v * \sigma'(a) (1 - \sigma(a)) * x$$

Q.14) Assume a two layer classification network with a single output where the hidden units and output unit use a Sigmoid activation. Assume a training set $\{x^{(i)}, y^{(i)}\}_{i=1}^m$ where

$x^{(i)} \in \mathbb{R}^n$ and $y^{(i)} \in \{0, 1\}$. Assume L2 loss ($\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$ instead of binary

cross entropy) derive the back propagation update equations for the output weights in a similar way to what was done in class for the binary cross entropy loss case.



Loss :-

$$L(\theta) = \frac{1}{2} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

Outputs :-

$$\hat{y} = \text{sigmoid}(V^T z)$$

$$z_j = \text{sigmoid}(w_j^T x)$$

Chain rule for output layer :-

$$\frac{\partial L(\theta)}{\partial V} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial V}$$

$$= \sum_{i=1}^m \frac{1}{2} \times 2 (\hat{y}^{(i)} - y^{(i)}) \hat{y}^{(i)} (1 - \hat{y}^{(i)}) z^{(i)}$$

$$= \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \hat{y}^{(i)} (1 - \hat{y}^{(i)}) z^{(i)}$$

Chain rule for hidden layer :-

$$\frac{\partial L(\theta)}{\partial w_j} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_j} \frac{\partial z_j}{\partial w_j}$$

$$\frac{\partial L}{\partial w_j} = \sum_{i=1}^m \frac{1}{2} \times 2 (\hat{y}^{(i)} - y^{(i)}) v_j z_j^{(i)} (1 - z_j^{(i)}) x^{(i)}$$

$$\boxed{\frac{\partial l}{\partial w_j} = \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \cdot v_j \cdot z_j^{(i)} (1 - z_j^{(i)}) x^{(i)}}$$

Update equation for output weights :-

$$V^{(\text{new})} = V^{(\text{old})} - n * \frac{\partial l}{\partial V} \quad \text{learning rate}$$

$$\boxed{V^{(\text{new})} = V^{(\text{old})} - n * \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \hat{y}^{(i)} (1 - \hat{y}^{(i)}) z^{(i)}}$$

Ans

Update equation for hidden layer weights :-

$$w_{kj}^{(\text{new})} = w_{kj}^{(\text{old})} - n * \frac{\partial l}{\partial w_{kj}}$$

$$\boxed{w_{kj}^{(\text{new})} = w_{kj}^{(\text{old})} - n * \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \cdot v_j \cdot z_j^{(i)} (1 - z_j^{(i)}) x^{(i)}}$$

Ans