

# **CS 577 - F22 - Assignment 2**

## **Questions Portion (Part 1)**

Due date: 10/04/2022

Anas Puthawala - A20416308

Professor Gady Agam

## Artificial Neurons

1.

$$1. \quad w = [0.1, 0.2, 0.3] \quad b_0 + b_1 x_1 + b_2 x_2$$

$$x = [1, 1] \quad y = 0.1 + 0.2(1) + 0.3(1)$$

$$= 0.6$$

2. For a linear discriminant, on one side the value is greater than 0, the other side has values less than 0 and on the boundary itself the value is 0.

3.  $b_0$  is the bias term and the distance from the decision boundary to the origin, whereas  $b_1$  and  $b_2$  are the normal vectors for the decision boundary.

4.

$$4) \quad g(x_1, x_2) = 1 + 2x_1 + 3x_2 \quad w = [1, 2, 3]$$

$$x = (0, 0)$$

$$\frac{\partial g}{\partial x_1} = 2 \quad \frac{\partial g}{\partial x_2} = 3. \quad \text{normal is } (2, 3)$$

$$\text{normalized normal} = \left( \frac{1}{\sqrt{13}}, \frac{3}{\sqrt{13}} \right)$$

$$\text{distance} : \frac{1}{\sqrt{13}} + \frac{2}{\sqrt{13}}(0) + \frac{3}{\sqrt{13}}(0) = \frac{1}{\sqrt{13}}$$

5.

$$5) \quad \theta^T = [b, \dots, b_n]$$

$$\uparrow$$

$$\text{Bias} = \text{first term}$$

$$\text{features go from } \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$\text{features: } \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \text{ must be augmented to include 1}$$

6.

$$6) \quad \begin{array}{l} \text{step} \\ \text{sigmoid} \end{array}$$

$$\mathcal{L}(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z \leq 0 \end{cases} \quad \text{sigmoid} = \frac{1}{1 + e^{-z}}$$

The advantage of using sigmoid activation over the step function is that the sigmoid function you can take derivatives for easily and attain useful information but for the step function the

derivative is going to be 0. The sigmoid also can be used to output probabilities so that can come in handy when we deal with classification problems. The step function is way more decisive as-well which can be a bad thing if there's something like 0.00001, it'll instantly assign it to a 1 even though it's significantly closer to 0.

When  $\theta$  is too small then the slope is a lot more flatter and less decisive (difficult to decide class).

7.

$$\begin{aligned}
 7) & \frac{P(y=1|x)}{P(y=0|x)} = \exp(\theta^T x) \\
 & \log\left(\frac{P(y=1|x)}{P(y=0|x)}\right) = \theta^T x \quad | \exp(\dots) \\
 & P(y=1|x) \left[ 1 + \exp(\theta^T x) \right] = \exp(\theta^T x) \\
 & P(y=1|x) = \frac{\exp(\theta^T x)}{1 + \exp(\theta^T x)} \quad | \text{same as sigmoid} \\
 & P(y=1|x) = \frac{\exp(\theta^T x)}{1 - P(y=1|x)} = \exp(\theta^T x) \\
 & P(y=1|x) = \exp(\theta^T x)(1 - P(y=1|x)) \\
 & \rightarrow P(y=1|x) + \exp(\theta^T x)P(y=1|x) = \exp(\theta^T x)
 \end{aligned}$$

The linear discriminant  $\theta^T x$  is used to model the log-likelihood ratio and by rearranging and mathematical manipulation you can obtain the sigmoid function equation in the box.

8.

$$\begin{aligned}
 \frac{\partial \text{loss}}{\partial z} &= \sigma(z)(1 - \sigma(z)) \quad \frac{\partial \log(\sigma(z))}{\partial z} = \frac{\partial \sigma(z)}{\partial z} \cdot \frac{\partial \log(\sigma(z))}{\partial z} \\
 &\rightarrow \sigma(z)(1 - \sigma(z)) \cdot \frac{1}{\sigma(z)} = [1 - \sigma(z)] \quad | \text{chain rule}
 \end{aligned}$$

9. The direction is computed by calculating the gradient, the derivative of the loss function. The size of the update is controlled by the learning rate.

10. You can stop when you're at the global minimum of the function (it could also get stuck at local minima however). The condition could be numerically tracked by the difference between the current gradient and the previous gradient and see if it's less than a certain threshold (i.e. the change is very small). This will ensure that there isn't any significant change in the gradient and that it may have indeed reached a local or global minima.

11. If a learning rate is too large then the algorithm will diverge, if the learning rate is too small it will take too small of steps and a much longer time to converge (may not ever converge) and plus it can get stuck in a local minima.

12. The empirical error loss calculates the total number of errors, it's not differentiable as it is a piecewise function and hence it can't be used in gradient descent.

13.

$$BCE = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1-y_i) \cdot \log(1-p(y_i))$$

*(log likelihood)*

$$\log \left( \prod_{x^{(i)} \in C_1} p(y=1 | x^{(i)}) \right) + \log \left( \prod_{x^{(i)} \in C_0} p(y=0 | x^{(i)}) \right) \Rightarrow p(y=0) = 1 - p(y=1)$$

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

*becomes*  $\sum_{i=1}^m y^{(i)} \log(p(y=1 | x^{(i)})) + (1-y^{(i)}) \log(1-p(y=1 | x^{(i)}))$

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

14.

$$\frac{d l(\theta)}{d \theta} = \frac{d}{d \theta} \left( \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right) \quad \text{where } h_\theta(x) = \theta^\top x$$

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

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

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

*Update rule:*

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

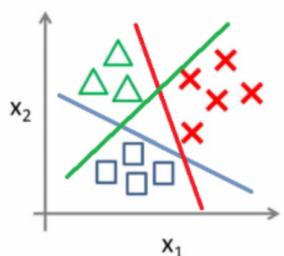
The meaning of the update equation is that it will advance the algorithm in search of the minima in the direction of the gradient which was computed by taking the derivative of the BCE loss function.

15.

One against all

- One against all tries to separate each individual class by having a discriminate function for each respective class.
- Therefore, there are 'k' discriminant functions where k is the number of classes
- You can imagine this by there being three clusters of data spaced out and there are exactly three lines that run and separate the three clusters of data, those three lines are the three respective discriminant functions:

Multi-class classification:



### One against each other

- One against each other aims to tackle one class at a time and derive a discriminant function to ensure that that singular class is best optimized and separated and that the discriminant function achieves high performance when separating the class
- There are a total of  $k(k-1)/2$  total models / discriminant functions that we must examine where  $k$  is the number of classes
- The primary dataset gets split into one binary classification set for each pair of classes. One against each other is easier to discriminate the data.

16. The rows of  $\theta^T$  are the templates and there is effectively one template per class.  $\theta^T x$  measures how similar the features match to each of the template and concurrently, high similarity can be indicative of high membership in the respective class. Multiple linear discriminant functions means you have multiple templates per each class.

17. Using softmax is an effective way to represent probabilities, all the class outputs will be in the interval 0,1 and the entirety of the class outputs would add up to 1. The larger the value in the range of 0,1 the higher the probability and certainty of it being that class

18.

$$(18) \text{Softmax}(z_j) = \frac{\exp(z_j)}{\sum_{i=1}^k \exp(z_i)}$$

$$\frac{\partial}{\partial \theta_i} \text{softmax}(z_j) = \text{softmax}(z_j) (\delta_{ij} - \text{softmax}(z_i))$$

$$\frac{\partial \text{softmax}(\theta_j^T x)}{\partial \theta_i} = \text{softmax}(\theta_j^T x) (\delta_{ij} - \text{softmax}(\theta_j^T x)) \cdot x$$

log sigmoid

$$\frac{\partial \log(\text{softmax}(\theta_j^T x))}{\partial \theta_i} = \frac{1}{\text{softmax}(\theta_j^T x)} \cdot \frac{\partial \text{softmax}(\theta_j^T x)}{\partial \theta_i}$$

$$= \frac{1}{\text{softmax}(\theta_j^T x)} \cdot \text{softmax}(\theta_j^T x) (\delta_{ij} - \text{softmax}(\theta_j^T x)) \cdot x$$

19.

$$(19) l(\theta_0, \dots, \theta_k) = \log \left( \prod_{i=1}^m \prod_{j=1}^k p(y^{(i)} = j | x^{(i)}) \right)$$

turns into

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

$$= - \sum_{i=1}^m \sum_{j=1}^k 1 \cdot (y^{(i)} = j) \log (\text{softmax}(x^{(i)}))$$

Categorical cross entropy.

20.

20)

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \alpha \nabla J(\theta^{(i)})$$

$$J(\theta) = -\sum_{i=1}^m \sum_{j=1}^k 1(y^{(i)}=j) \log(\text{logit}(x^{(i)}))$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_{i=1}^m (\text{logit}(x^{(i)}) - 1(y^{(i)}=j)) x^{(i)}$$

$$\theta_j \leftarrow \theta_j - \alpha \sum_{i=1}^m (\text{logit}(x^{(i)}) - 1(y^{(i)}=j)) x^{(i)}$$

update eqn.

The update rule shows how the new value of the parameter theta is going to be the old value plus an added factor adjusting it based on the partial derivative of the loss function for categorical cross entropy which in other words is also known as the gradient. The update equation is iteratively adjusted until it reaches either a local minima or a global minima. The parameter 'alpha' is the learning rate which is responsible for the size of the adjustment and is a hyper-parameter that gets tuned. The update equation shows that the softmax will output multiple classes as it looks at the 'k' classes.

## Neural Networks

1. It maps the inputs to a higher dimensional space in an effort to find finer relationships between the original inputs. The benefit of this process is that it is easier to find relationships between features when expanding the dimension of said features and analyzing between them in the higher dimensional space.
2. Dimensionality reduction, it compresses the input features into a lower dimensional space. The purpose and benefit of this process can be to extract solely the useful features or extract the features that have useful information encoded within them.
3. The gradients will differ between the two layers. The values z and Z\_hat are unknown and so we can't get the expected values we want from the hidden layer unless there are additional steps that we follow up with. You can see from the following picture that the chain rule for the two layers differs:

chain rule:

$$\begin{cases} \frac{\partial E}{\partial v} = \frac{\partial E}{\partial g} \frac{\partial g}{\partial v} \\ \frac{\partial E}{\partial w_j} = \frac{\partial E}{\partial g} \frac{\partial g}{\partial z_j} \frac{\partial z_j}{\partial w_j} \end{cases}$$

$$v_j \leftarrow v_j - \eta \sum_{i=1}^m (\hat{y}_j^{(i)} - 1(y^{(i)}=j)) z^{(i)}$$

$$w_j \leftarrow w_j - \eta \sum_{i=1}^m (\hat{z}_j^{(i)} - 1(z^{(i)}=j)) x^{(i)}$$

$$\hat{z}_j^{(i)} = ? \quad z^{(i)} = ?$$

4.

$$(4) \text{ MSE } E = \frac{1}{2} \sum_{i=1}^k (y^{(i)} - \hat{y}^{(i)})^2 \quad \hat{y} = z \quad \text{hidden layer: } \frac{\partial E}{\partial w_j} = \frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_j}$$

$$z = \sigma(w_j^T x)$$

hidden

$$- \sum_{i=1}^k (y^{(i)} - \hat{y}^{(i)}) \cdot 1 \cdot \sigma(w_j^T x) (1 - \sigma(w_j^T x)) \cdot x^{(i)}$$

output layer:

$$\frac{\partial E}{\partial v} = \frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial v} = \left[ - \sum_{i=1}^k (y^{(i)} - \hat{y}^{(i)}) \cdot 1 \right]$$

5.

(5) outputs  $\hat{y}_j = v_j^T z = v_{j0} + v_{j1} z_1 + v_{j2} z_2 + \dots + v_{jn} z_n$

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

$$\text{loss, } E = \frac{1}{2} \sum_{i=1}^m \sum_{e=1}^k (\hat{y}_e^{(i)} - y_e^{(i)})^2$$

$$\frac{\partial E}{\partial v_j} = \frac{\partial E}{\partial \hat{y}_j} \cdot \frac{\partial \hat{y}_j}{\partial v_j} \quad \frac{\partial E}{\partial \hat{y}_j} = 2 \cdot \frac{1}{2} \sum_{i=1}^m (\hat{y}_j^{(i)} - y_j^{(i)})^2$$

$$= \sum_{i=1}^m (\hat{y}_j^{(i)} - y_j^{(i)})^2 \cdot z^{(i)}$$

$$\frac{\partial \hat{y}_j}{\partial v_j} = z^{(i)}$$

Hidden layer

$$\frac{\partial E}{\partial w_j} = \sum_{q=1}^k \frac{\partial E}{\partial \hat{y}_q} \cdot \frac{\partial \hat{y}_q}{\partial z_j} \cdot \frac{\partial z_j}{\partial w_j} \quad \frac{\partial \hat{y}_q}{\partial z_j} = v_{jq}$$

$$\frac{\partial z_j}{\partial w_j} = z_j^{(i)} (1 - z_j^{(i)}) x^{(i)}$$

$$\therefore \sum_{q=1}^k \sum_{i=1}^m (\hat{y}_q^{(i)} - y_q^{(i)})^2 \cdot v_{jq} \cdot z_j^{(i)} (1 - z_j^{(i)}) x^{(i)}$$

6.

(b) loss: binary cross entropy  $\rightarrow E = \sum_{j=1}^k y_j^{(i)} \log(\hat{y}_j^{(i)})$   
 hidden layer: sigmoid  $\rightarrow z = \sigma(w^T x)$   $\hat{y} = \sigma(v^T z)$   
 output: sigmoid.  $\hat{y} = \hat{y}(1 - \hat{y})$

hidden:  $\frac{\partial E}{\partial w_j} = \left( \frac{\partial E}{\partial \hat{y}} \right) \cdot \frac{\partial \hat{y}}{\partial z_j} \cdot \frac{\partial z_j}{\partial w_j}$

$$\begin{aligned} \frac{\partial E}{\partial \hat{y}} &= f'_{\text{get}} \hat{y} = \hat{y}_j(1 - \hat{y}_j) \log(\hat{y}_j) + y_j^{(i)} \cdot \frac{1}{\hat{y}_j} \cdot \hat{y}_j(1 - \hat{y}_j) \\ &= \hat{y}_j(1 - \hat{y}_j) \log(\hat{y}_j) + y_j^{(i)}(1 - \hat{y}_j) \end{aligned}$$

$$\frac{\partial \hat{y}}{\partial z_j} = \sigma(v^T z)(1 - \sigma(v^T z)) \cdot v_j = \hat{y}_j(1 - \hat{y}_j) \cdot v_j$$

$$\frac{\partial z_j}{\partial w_j} = \hat{y}_j(1 - \hat{y}_j) \cdot x^{(i)}$$

hidden:

$$\sum_{i=1}^m \hat{y}_j(1 - \hat{y}_j) \log(\hat{y}_j) + y_j^{(i)}(1 - \hat{y}_j) \cdot \hat{y}_j(1 - \hat{y}_j) \cdot v_j \cdot z_j^{(i)}(1 - z_j^{(i)}) \cdot x^{(i)}$$

output layer:

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

7.

7) multi-class  $L = \text{cat. cross entr} \rightarrow \text{true label}$

hidden: sigmoid  $\rightarrow \sum_i \hat{y}_i \log(\hat{y}_i) \rightarrow \text{probability}$   
 output: softmax  $\rightarrow \text{its final softmax output}$

$$\frac{\partial E}{\partial w_j} = \sum_{i=1}^k \frac{\partial E}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial z_j} \cdot \frac{\partial z_j}{\partial w_j} = \sum_{q=1}^m \sum_{i=1}^k (\hat{y}_i^{(q)} - y_i^{(q)}) v_i q \frac{\partial \hat{y}_i^{(q)}}{\partial z_j} (1 - z_j^{(q)}) x^{(q)}$$

output layer  $= \frac{\partial E}{\partial v_j} = \frac{\partial E}{\partial \hat{y}_j} \cdot \frac{\partial \hat{y}_j}{\partial v_j} = \sum_{i=1}^m (\hat{y}_i^{(i)} - y_i^{(i)}) z_i^{(i)}$

8. The weights are initialized randomly but they must not be large values, they need to be closer to 0. The aim of the initialization is to provide a good starting point for the optimization algorithm to take over and ensure that the weights either don't explode or vanish within the first few iterations. If the weights are initialized all as 0 then the neurons will learn the same features during training, in-fact any initialization scheme that is initializing the weights to be constant will all result in the same weights being trained. You need the weights to be closer to 0 so they don't vanish and/or explode and different so that the weights all learn and pick-up different features and values.

---

## Computation Graphs

1. The advantage of using computation graphs is that it allows for easier computation of the gradients in order to update the weights and optimize them for the network. Computation graphs provide a way to optimize the computation and so that it is also more intuitive for someone to look at and quickly and easily be able to figure out the gradients as needed for back-prop.

In the forward pass the weights and features are passed through the network and the activation functions are applied to the weights and features as designated to each node. Whereas in the backward pass the gradients are calculated with respect to the nodes and each input.

Each node in the graph needs to be able to compute a value (an activation) based on the features / data it's being fed in and this value is configured by the user (so the user can choose softmax, sigmoid, linear, relu, whatever activation they want). In addition, in the backward pass the each node must be able to compute the gradients, which is the derivatives for a specific path (the derivative with respect to a feature / input).

Each node needs to store the outputs and the values from the forward pass so that it can compute the gradients based off that and adjust the values as calculated from the gradient to optimize the weights for the network.

2) hidden out:  $z = f_1(w_1, x)$        $\hat{y} = f_2(w_2, f_1(w_1, x))$

output out:  $\hat{y} = f_2(w_2, z)$        $E = \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$

$$E = \frac{1}{2} \sum_{i=1}^m (y^{(i)} - f_2(w_2, f_1(w_1, x)))^2$$

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

$$\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2} \quad \frac{\partial E}{\partial w_1} = \cancel{\frac{\partial E}{\partial \hat{y}}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_1} = \frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial f_2}{\partial z} \cdot \frac{\partial f_1}{\partial w_1}$$

$$\frac{\partial E}{\partial \hat{y}} = 2 \cdot \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) (-1) = - \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})$$

$$\frac{\partial \hat{y}}{\partial z} = \frac{\partial f_2}{\partial z} \quad \frac{\partial z}{\partial w_1} = \frac{\partial f_1}{\partial w_1}$$

Since  $z = f_1(\dots)$ ,  $\frac{\partial z}{\partial w_1} = \frac{\partial f_1}{\partial w_1}$  ... same logic applies for  $\frac{\partial \hat{y}}{\partial z}$ .

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

$$\frac{\partial E}{\partial w_2} = \cancel{\frac{\partial E}{\partial \hat{y}}} \cdot \frac{\partial \hat{y}}{\partial w_2} = - \left[ \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) \cdot \frac{\partial f_2}{\partial w_2} \right] = - \sum_{i=1}^m (y^{(i)} - f_2(w_2, z)) \cdot \frac{\partial f_2}{\partial w_2}$$

2. cross entropy  $L = - \sum_i y_i \ln(\hat{y}_i)$        $\hat{y} = f_2(w_2, f_1(w_1, x))$

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

hidden out =  $f_1(w_1, x) = z$ , out out =  $\hat{y}$ .

$$\frac{\partial L}{\partial w_1} = \cancel{\frac{\partial L}{\partial \hat{y}}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_1} = \left( \frac{\partial L}{\partial f_2} \right) \frac{\partial f_2}{\partial f_1} \cdot \frac{\partial f_1}{\partial w_1}$$

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

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

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

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

~~$\sigma + \eta$~~

3. L<sub>2</sub> loss:  $\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^k (y_j^{(i)} - f_2(w_i, f_i(w_i, x^{(i)})))^2 = L$

$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial f_2} \cdot \frac{\partial f_2}{\partial f_i} \cdot \frac{\partial f_i}{\partial w_i}$

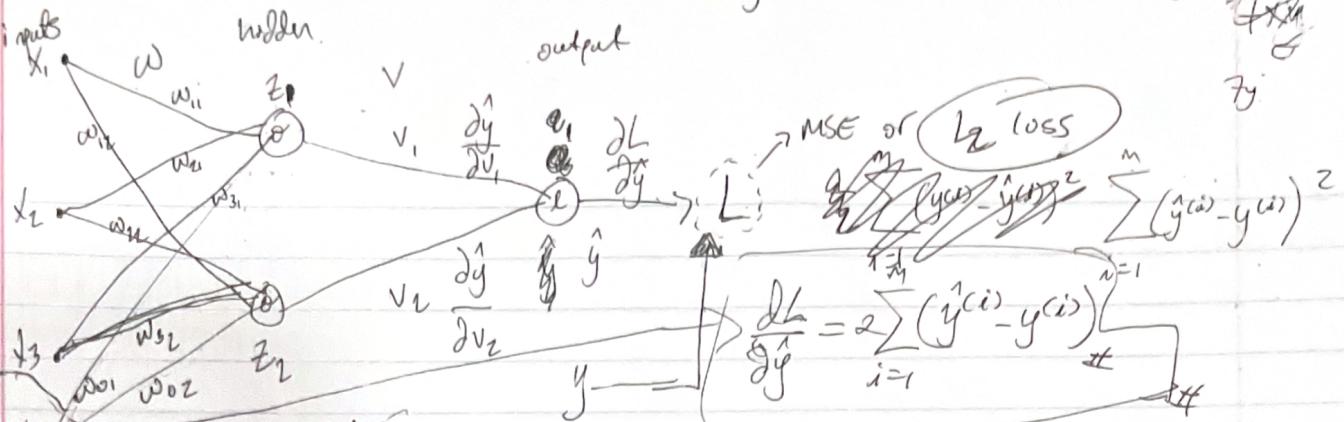
$\Rightarrow \sum_{i=1}^m \sum_{j=1}^k (y_j^{(i)} - f_2(w_i, f_i(w_i, x^{(i)}))) \cdot \frac{\partial f_2}{\partial f_i} \cdot \frac{\partial f_i}{\partial w_i} \#$

cross entropy  $L = \sum_{i=1}^m \sum_{j=1}^k y_j^{(i)} \cdot \log(f_2(w_i, f_i(w_i, x^{(i)}))) \approx$

$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial f_2} \cdot \frac{\partial f_2}{\partial f_i} \cdot \frac{\partial f_i}{\partial w_i} = \sum_{i=1}^m \sum_{j=1}^k \frac{y_j^{(i)}}{f_2(w_i, f_i(w_i, x^{(i)}))} \cdot \frac{\partial f_2}{\partial f_i} \cdot \frac{\partial f_i}{\partial w_i} \#$

$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial f_2} \cdot \frac{\partial f_2}{\partial w_2} = \sum_{i=1}^m \sum_{j=1}^k \frac{y_j^{(i)}}{f_2(w_i, f_i(w_i, x^{(i)}))} \frac{\partial f_2}{\partial w_2} \#$

# regression network



Bias

$$\text{justification: } \hat{y} = V^T z \quad z_j = \sigma(w_j^T x)$$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial (V_0 + V_1 z_1 + \dots + V_n z_n)} = z_1 \star$$

gradients

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

$$\frac{\partial \hat{y}}{\partial V_1} = \frac{\partial L}{\partial V_1} =$$

$$= 2 \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \cdot z_1^{(i)} \star$$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial y}$$

$$= 2 \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \cdot z_1^{(i)} \star$$

$$= \frac{\partial L}{\partial \hat{y}} \cdot z_1^{(i)} \star$$

$$w = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_1} = \frac{\partial L}{\partial \hat{y}} \cdot$$

$$\frac{\partial z_1}{\partial w_{11}} = \frac{\partial (\sigma(w_{11}^T x))}{\partial w_{11}}$$

$$= \sigma'(w_{11} + w_{11} x_1 + w_{21} x_2 + \dots + w_{31} x_3) \star$$

$$= \sigma'(w_{11}^T x) (1 - \sigma(w_{11}^T x)) \cdot x_1^{(i)} \star$$

$$\frac{\partial L}{\partial w_{11}} = \left( \frac{\partial L}{\partial \hat{y}} \right) \cdot V_1 \cdot z_1^{(i)} (1 - z_1^{(i)}) \cdot x_1^{(i)} \star$$

$$\frac{\partial L}{\partial w_{12}} = \left( \frac{\partial L}{\partial \hat{y}} \right) \cdot \frac{\partial \hat{y}}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_{12}} = \frac{\partial L}{\partial \hat{y}} \cdot V_2 \cdot z_2^{(i)} (1 - z_2^{(i)}) \cdot x_1^{(i)} \star$$

$$\frac{\partial L}{\partial w_{21}} = \left( \frac{\partial L}{\partial \hat{y}} \right) \cdot V_1 \cdot z_1^{(i)} (1 - z_1^{(i)}) \cdot x_2^{(i)} \star$$

$$\frac{\partial L}{\partial w_{22}} = \left( \frac{\partial L}{\partial \hat{y}} \right) \cdot V_2 \cdot z_2^{(i)} (1 - z_2^{(i)}) \cdot x_2^{(i)} \star$$

$$\frac{\partial L}{\partial w_{31}} = \frac{\partial L}{\partial \hat{y}} \cdot V_1 \cdot z_1^{(i)} (1 - z_1^{(i)}) \cdot x_3^{(i)} \star$$

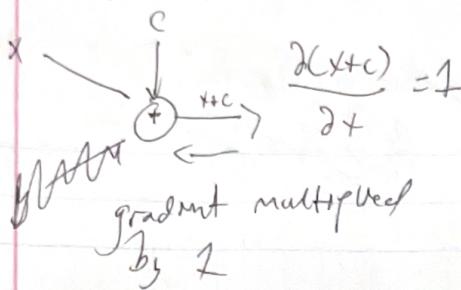
$$\frac{\partial L}{\partial w_{32}} = \frac{\partial L}{\partial \hat{y}} \cdot V_2 \cdot z_2^{(i)} (1 - z_2^{(i)}) \cdot x_3^{(i)} \star$$

Bias

$$\frac{\partial L}{\partial w_{01}} = \frac{\partial L}{\partial \hat{y}} \cdot V_1 \cdot z_1^{(i)} (1 - z_1^{(i)}) \cdot 1 \star$$

$$\frac{\partial L}{\partial w_{02}} = \frac{\partial L}{\partial \hat{y}} \cdot V_2 \cdot z_2^{(i)} (1 - z_2^{(i)}) \cdot 1 \star$$

5. addition of a constant



addition of two inputs

$$\begin{aligned} a &\xrightarrow{\text{arb}} \text{sum} \\ b &\xrightarrow{\text{arb}} \text{sum} \end{aligned}$$

$$\frac{\partial(\text{sum})}{\partial a} = 1$$

$$\frac{\partial(\text{sum})}{\partial b} = 1$$

In both cases, gradient is multiplied by 1.

multiplication of a constant

$$a \xrightarrow{\text{mult}} ac$$

$$\frac{\partial(ac)}{\partial a} = c$$

gradient multiplied by constant 'c'

Mult. of two inputs

$$\begin{aligned} a &\xrightarrow{\text{mult}} ab \\ b &\xrightarrow{\text{mult}} ab \end{aligned}$$

$$\frac{\partial(ab)}{\partial a} = b$$

$$\frac{\partial(ab)}{\partial b} = a$$

gradient multiplied by the other input,

max of two inputs

$$a \xrightarrow{\text{max}} \max(a, b)$$

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

$$\frac{\partial(\max(a, b))}{\partial a} = \begin{cases} 1 & \text{if } a \geq b \\ 0 & \text{else} \end{cases}$$

$$\frac{\partial(\max(a, b))}{\partial b} = \begin{cases} 0 & \text{if } a \geq b \\ 1 & \text{else} \end{cases}$$

the gradient is multiplied by 1 for the maximum input only,

min of two inputs

$$\min(a, b) = \begin{cases} a & \text{if } a \leq b \\ b & \text{else} \end{cases}$$

for min of two inputs,

the gradient will also be multiplied by 1 for the minimum input only

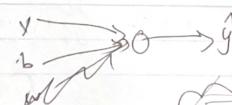
6. If the input node is a vector (rank 1) and the output is a scalar (rank 0) then the derivative that will flow back will be a tensor of rank  $(1) + (0) = 1$ .

7. If the input and output nodes are both vectors (both rank 1) then the rank of the gradient tensor will be  $(1) + (1) = \text{rank 2 tensor}$  flowing back.

8. If the input node is a matrix (rank 2) and the output is a scalar (rank 0) then the rank of the gradient tensor will be  $2+0 = 2$ .

$$6. \begin{bmatrix} \frac{\partial y}{\partial w_1} \\ \vdots \\ \frac{\partial y}{\partial w_n} \end{bmatrix}$$

$$7. \begin{bmatrix} \frac{\partial y}{\partial w_1} & \dots & \frac{\partial y}{\partial w_m} \\ \frac{\partial y}{\partial w_1} & \dots & \frac{\partial y}{\partial w_m} \end{bmatrix}$$



$$9. (F \circ G) = F(g(x)) \quad \text{and} \quad \underline{\nabla(F \circ G)} = \underline{\frac{\partial F}{\partial g}} \cdot \underline{\nabla g}$$

let  $f_{\text{def}} = F(x_0, x_1, \dots, x_n) = \begin{bmatrix} f_1(x_0, x_1, \dots, x_n) \\ \vdots \\ f_m(x_0, x_1, \dots, x_n) \end{bmatrix}$

let  $g = G(x_0, x_1, \dots, x_n) = \begin{bmatrix} g_1(x_0, x_1, \dots, x_n) \\ \vdots \\ g_m(x_0, x_1, \dots, x_n) \end{bmatrix}$

and  $\therefore (F \circ G) = F(g) = \begin{bmatrix} f_1(g_1(x_0), \dots, g_1(x_n)) \\ \vdots \\ f_m(g_m(x_0), \dots, g_m(x_n)) \end{bmatrix}$

$\nabla F \circ G|_x = \text{Jacobian}$

$$= \begin{bmatrix} \nabla f_1(g_1(x_0), \dots, g_1(x_n)) \\ \vdots \\ \nabla f_m(g_m(x_0), \dots, g_m(x_n)) \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1(g_1(x_0))}{x_0} \dots \frac{\partial f_1(g_1(x_n))}{x_n} \\ \vdots \\ \frac{\partial f_m(g_m(x_0))}{x_0} \dots \frac{\partial f_m(g_m(x_n))}{x_n} \end{bmatrix}$$

10.

$$F(x, y, z) = [3xy \ y - z]^T \quad (F \circ G)(x, y) = F(g(x, y))$$

$$G(x, y) = [x - 5y \ xy \ xy - x^2]^T \quad \text{Jacobian: } \begin{bmatrix} 6xy - 15y^2 & 3x^2 - 30xy \\ y - 1 & x + 1 \end{bmatrix}$$

$$\begin{bmatrix} 3(x - 5y)(xy) & xy - (x^2 - y) \end{bmatrix}$$

$$(F \circ G)(x, y) = \begin{bmatrix} 3x^2y - 15xy^2 & xy - x^2 + y \end{bmatrix} \xrightarrow{*} \nabla F \circ G(x, y) = \begin{bmatrix} \frac{\partial f_1}{\partial x} = 6xy - 15y^2 & \frac{\partial f_1}{\partial y} = 3x^2 - 30xy \\ \frac{\partial f_2}{\partial x} = y - 1 & \frac{\partial f_2}{\partial y} = x + 1 \end{bmatrix}$$

10. chain rule  $f(x, y, z) = \begin{bmatrix} 3xy \\ y - z \end{bmatrix} \quad f(x, y) = \begin{bmatrix} x - 5y \\ xy \end{bmatrix} \xrightarrow{*} 3x$

$$J_{f \circ g}(x, y) = J_f(g(x, y)) \cdot J_g(x, y) \xrightarrow{*} \text{chain rule}$$

first lets compute  $J_g(x, y) = \begin{bmatrix} 1 & -5 \\ y & x \\ +1 & -1 \end{bmatrix}$

now compute  $J_f(g(x, y)) = \underbrace{J_f(x, y)}_{= J_f(x, y) \circ g(x, y)} \xrightarrow{*} = \begin{bmatrix} 3y & 3x & 0 \\ 0 & 1 & -1 \end{bmatrix} \xrightarrow{dM=23}$

$$J_f(x, y, z) \circ g(x, y) = \begin{bmatrix} 3xy & 3x - 15y & 0 \\ 0 & 1 & -1 \end{bmatrix} \xrightarrow{dM=23} = J_f(g(x, y))$$

Finally:  $J_{f \circ g}(x, y) = J_f(g(x, y, z)) \cdot J_g(x, y) = \begin{bmatrix} 3xy & 3x - 15y & 0 \\ 0 & 1 & -1 \end{bmatrix} \cdot \begin{bmatrix} 1 & -5 \\ y & x \\ +1 & -1 \end{bmatrix} \xrightarrow{*}$

$$\rightarrow \begin{bmatrix} 3xy & 3x^2 - 15y & 0 \\ 0 & 1 & -1 \end{bmatrix} \cdot \begin{bmatrix} 1 & -5 \\ y & x \\ 1 & -1 \end{bmatrix} =$$

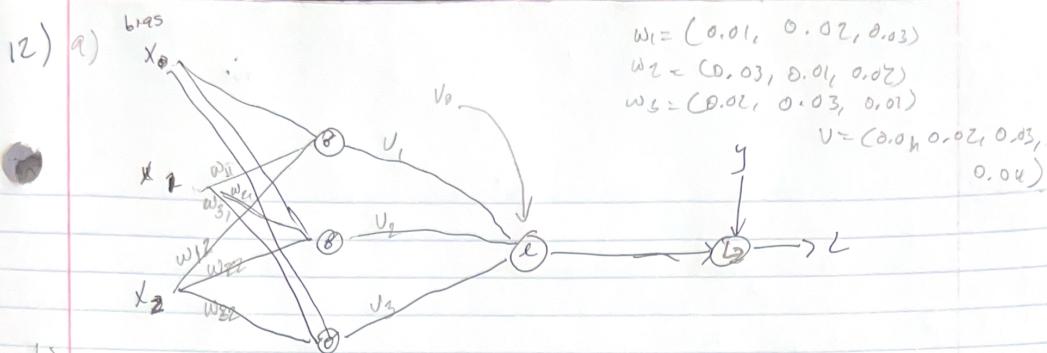
10 continued.

$$\begin{bmatrix} 3xy + 3xy \cdot (5y^2 + 0) & -15xy + 3x^2 - 15xy + 0 \\ 0 + y - 1 & 0 + x + 1 \end{bmatrix} = \begin{bmatrix} 6xy - 15y^2 & 3x^2 - 30xy \\ y - 1 & x + 1 \end{bmatrix}$$

Same as when calculated  
Via direct method ✓

11. It is important to order the nodes in order to properly take the necessary derivatives to compute the gradient for the update rule so that we can optimally converge the optimization algorithm.

12.



$$b_1, b_2, b_3$$

$$x_1, x_2$$

$$w_1, w_2, w_3$$

$$u_1, u_2$$

$$l$$

$$y$$

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

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

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

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

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

$$z_1^{(1)} = \sigma(w_1^T x_1) = \sigma(1(0.01) + 1(0.02) + 2(0.03)) = 0.52248$$

$$z_2^{(1)} = \sigma(w_2^T x_1) = \sigma(1(0.03) + 1(0.01) + 2(0.02)) = 0.519989$$

$$z_3^{(1)} = \sigma(w_3^T x_1) = \sigma(1(0.02) + 1(0.03) + 2(0.01)) = 0.517493$$

$$y = \sqrt{z^{(1)}} = 0.01 + (0.02)(0.52248) + (0.03)(0.519989) + (0.04)(0.517493) = 0.05675$$

$$L_2 \text{ loss} = (8 - 0.05675)^2 = 63.095$$

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

$$z_1^{(2)} = \sigma(w_1^T x_2) = \sigma(0.03 + 0.01(1) + 0.02(3)) = \sigma(0) = 0.52498$$

$$z_2^{(2)} = \sigma(w_2^T x_2) = \sigma(0.01 + 0.02(1) + 0.03(3)) = \sigma(1) = 0.52996$$

$$z_3^{(2)} = \sigma(w_3^T x_2) = \sigma(0.02 + 0.03(1) + 0.01(3)) = \sigma(0.08) = 0.519989$$

$$y = \sqrt{z^{(2)}} = 0.01 + (0.02)(0.52498) + (0.03)(0.52996) + (0.04)(0.519989) = 0.05715$$

$$L_2 \text{ loss} = (11 - 0.05715)^2 = 119.74559$$

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

$$y = \sqrt{z^{(3)}}$$

$$z_1^{(3)} = \sigma(w_1^T x_3) = \sigma(1) = 0.52747 = 0.5722303$$

$$z_2^{(3)} = \sigma(w_2^T x_3) = \sigma(0.9) = 0.522184 L_2 \text{ loss} = (10 - 0.5722303)^2$$

$$z_3^{(3)} = \sigma(w_3^T x_3) = \sigma(1) = 0.524999 = 0.519943$$

$$12c) \frac{\partial L_2}{\partial w_{11}w_{12}w_{13}} = \frac{\partial L_2}{\partial y} \frac{\partial y}{\partial v} \frac{\partial v}{\partial (w_{11} \dots w_3)} \quad \text{find } \frac{\partial L_2}{\partial w_1}, \frac{\partial L_2}{\partial w_2}, \frac{\partial L_2}{\partial w_3}$$

$$\frac{\partial L_2}{\partial v} = \frac{\partial L_2}{\partial y^{(i)}} \cdot \frac{\partial y^{(i)}}{\partial v}$$

12 continued.

$$\frac{\partial L_2}{\partial v}$$

$$\frac{\partial L_2}{\partial v} = \sum_{i=1}^m \left[ (\sigma(w_1 \times v^{(i)}), \sigma(w_2 \times v^{(i)}), \sigma(w_3 \times v^{(i)}) \cdot v) - y^{(i)} \cdot (\sigma(w_1 \times v^{(i)}), \sigma(w_2 \times v^{(i)}), \sigma(w_3 \times v^{(i)})) \right]$$

$$\frac{\partial L_2}{\partial (w_1 w_2 w_3)} = \sum_{i=1}^m \left[ (\sigma(w_1 \times v^{(i)}), \sigma(w_2 \times v^{(i)}), \sigma(w_3 \times v^{(i)})) \cdot V \cdot \frac{\partial Z^T}{\partial x} \right]$$

$$= \begin{bmatrix} \sigma(w_1 \times v^{(i)}) (1 - \sigma(w_1 \times v^{(i)}) \cdot x^{(i)}) \\ \sigma(w_2 \times v^{(i)}) (1 - \sigma(w_2 \times v^{(i)}) \cdot x^{(i)}) \\ \sigma(w_3 \times v^{(i)}) (1 - \sigma(w_3 \times v^{(i)}) \cdot x^{(i)}) \end{bmatrix}$$

$$12c) \frac{\partial L_2}{\partial w_{11}w_{12}w_{13}} = \frac{\partial L_2}{\partial y} \frac{\partial y}{\partial v} \frac{\partial v}{\partial (w_{11} \dots w_3)} \quad \text{find } \frac{\partial L_2}{\partial w_1}, \frac{\partial L_2}{\partial w_2}, \frac{\partial L_2}{\partial w_3}$$

$$\frac{\partial L_2}{\partial v} = \frac{\partial L_2}{\partial y^{(i)}} \cdot \frac{\partial y^{(i)}}{\partial v}$$

$$\frac{\partial L_2}{\partial v} = \sum_{i=1}^m \left[ (\sigma(w_1 \times v^{(i)}), \sigma(w_2 \times v^{(i)}), \sigma(w_3 \times v^{(i)}) \cdot v) - y^{(i)} \cdot (\sigma(w_1 \times v^{(i)}), \sigma(w_2 \times v^{(i)}), \sigma(w_3 \times v^{(i)})) \right]$$

$$\frac{\partial L_2}{\partial (w_1 w_2 w_3)} = \sum_{i=1}^m \left[ (\sigma(w_1 \times v^{(i)}), \sigma(w_2 \times v^{(i)}), \sigma(w_3 \times v^{(i)})) \cdot V \cdot \frac{\partial Z^T}{\partial x} \right]$$

$$= \begin{bmatrix} \sigma(w_1 \times v^{(i)}) (1 - \sigma(w_1 \times v^{(i)}) \cdot x^{(i)}) \\ \sigma(w_2 \times v^{(i)}) (1 - \sigma(w_2 \times v^{(i)}) \cdot x^{(i)}) \\ \sigma(w_3 \times v^{(i)}) (1 - \sigma(w_3 \times v^{(i)}) \cdot x^{(i)}) \end{bmatrix}$$

$$\frac{\partial L_2}{\partial v} = (\sigma(.70575) \cdot (.7^{(i)})) = 7.94325 (\sigma(.7 \cdot X_0)) =$$

$$\frac{\partial L_2}{\partial v} = \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) Z_1^{(i)} = 7.94325 (.52248) + (1 - .05715) (.52996) + (10 - .057273) (.52747)$$

$$\cong 15.019$$

$$\frac{\partial L_2}{\partial v} = \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) Z_2^{(i)} = 15.066$$

$$\frac{\partial L_2}{\partial v} = (15.019, 15.066, 15.017)$$

$$\frac{\partial L_2}{\partial v} = | 1 | = 15.017$$

$$\frac{\partial L_2}{\partial w_1} = \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \cdot v_i \cdot Z_1^{(i)} (1 - Z_1^{(i)}) \cdot x^{(i)}$$

$$, 0545$$

$$\frac{\partial L_2}{\partial w_1} = (8 - .05675) \cdot 0.02 \cdot (.52248)(1 - .52248) \cdot | + (10 - .05715) \cdot 0.02 \cdot (.52996) \cdot | + (10 - .057273) \cdot 0.02 \cdot (.52747)$$

$$+ 0.3964$$

$$+ (10 - .057273) \cdot 0.02 \cdot (.52747)$$

$$\frac{\partial L_2}{\partial w_2} = .3472$$

$$\frac{\partial L_2}{\partial w_10} = .1439$$

$$\frac{\partial L_2}{\partial w_2} = .1439$$

$$\frac{\partial L_2}{\partial w_2} = (.216, .29, .5133)$$

$$\frac{\partial L_2}{\partial w_3} = .1788$$

$$\frac{\partial L_2}{\partial w_3} = .1387$$

$$\frac{\partial L_2}{\partial w_3} = .1685$$

$$13a) f(x, y) = (2x+3y)^2 \quad \nabla f(x, y)$$

$$= [2(2x+3y) \cdot 2, 2(2x+3y) \cdot 3] \\ = [8x+12y, 12x+18y]$$

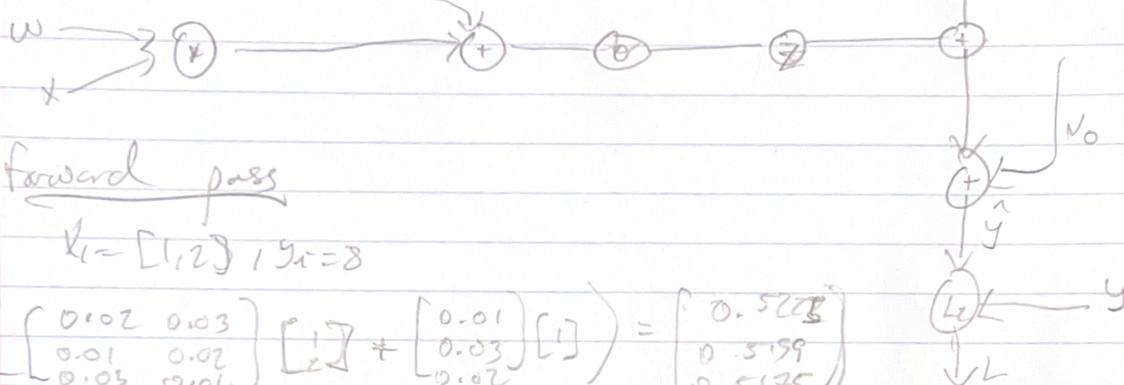
$$13b) F(x, y) = \begin{bmatrix} x^2+2y \\ 3x+4y^2 \end{bmatrix} \quad J_F \leftarrow \begin{bmatrix} 2x & 2 \\ 3 & 8y \end{bmatrix} \Big|_{(1,2)} = \begin{bmatrix} 2 & 2 \\ 3 & 16 \end{bmatrix}$$

13c) ~~FE~~

$$G(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix} \quad F \circ G = F(g(x)) \quad \begin{bmatrix} (x)^2 + 2(x^2) \\ 3(x) + 4(x^4) \end{bmatrix} = \begin{bmatrix} x^2 + 2x^2 \\ 3x + 4x^4 \end{bmatrix} = \alpha$$

$$J_G(x) = \frac{2x+4x}{3+16x^3} = \begin{bmatrix} 6x \\ 16x^3+3 \end{bmatrix} \Big|_{x=2} = \begin{bmatrix} 12 \\ 131 \end{bmatrix}$$

13d) ~~No~~



forward pass

$$x_i = [1, 2], y_i = 8$$

$$\sigma \left( \begin{bmatrix} 0.02 & 0.03 \\ 0.01 & 0.02 \\ 0.03 & 0.01 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} + \begin{bmatrix} 0.01 \\ 0.03 \\ 0.02 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} \right) = \begin{bmatrix} 0.5223 \\ 0.5199 \\ 0.5175 \end{bmatrix}$$

$$\sigma \left[ (0.02 \ 0.03 \ 0.04) \cdot 2 \right] = 0.01$$

$$y \approx 0.5223, y = 8 \quad L_2(y, y) = 63.1$$

backward pass

$$\frac{\partial L_2}{\partial y} = 2(8 - 0.5223) = 15.8, \quad \frac{\partial y}{\partial v} (v_z + v_b) = 2, \quad \frac{\partial y}{\partial v_b} = 1$$

$$\frac{\partial L_2}{\partial y} \frac{\partial y}{\partial v} = 15.8 \quad \begin{bmatrix} 0.5223 \\ 0.5199 \\ 0.5175 \end{bmatrix}$$

$$\frac{\partial L_2}{\partial y} \frac{\partial y}{\partial v_b} = 15.8 \quad \frac{\partial y}{\partial v} (v_z + v_b) = \begin{bmatrix} 0.02 & 0.03 \\ 0.01 & 0.02 \\ 0.03 & 0.01 \end{bmatrix} = V$$

$$\frac{\partial z}{\partial w} (w_a + w_b) = \begin{bmatrix} 1 \\ z \end{bmatrix} \times \frac{\partial h_z}{\partial w} = \frac{\partial L_z}{\partial g} \cdot \left( \frac{\partial g}{\partial z} \right) \cdot \left( \frac{\partial z}{\partial w} \right)$$

$$= 15.8 \cdot \begin{bmatrix} 0.02 & 0.03 \\ 0.01 & 0.02 \\ 0.03 & 0.01 \end{bmatrix} \begin{bmatrix} 1 \\ z \end{bmatrix}$$

vector

14.

