

Assignment 2

Artificial neurons

$$\begin{aligned} \Rightarrow g(x_1, x_2) &= \Theta^T x + \Theta_0 \\ &= (0.2x_1) + (0.3x_2) + 1 \\ &= 0.6 \end{aligned}$$

2) On one side of the boundary, the value should be greater than 0. On the other side, less than 0. On the boundary, 0.

3) Θ_0 is the bias, or the negative distance from the origin.

Θ_1, Θ_2 are the weights and make up for the normal vector

$$\Rightarrow 4) \text{normal} = [2, 3]$$

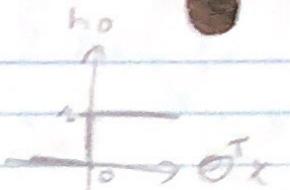
$$\text{normal normalized} = \frac{1}{\sqrt{13}} [2, 3]$$

$$\text{Distance} = \frac{1}{\sqrt{13}}$$

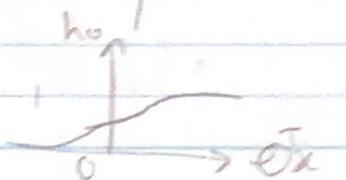
5) Augment the x vector by one and add 1 to the beginning of the vector.

$$\Theta^T x + \Theta_0 = \tilde{\Theta}^T \tilde{x} \quad \text{where } \begin{aligned} \tilde{\Theta} &= [\Theta_0, \dots, \Theta_n] \\ \tilde{x} &= [1, x_1, \dots, x_n] \end{aligned}$$

6) Step function: $h_0 = \begin{cases} 1 & \text{if } \theta^T x > 0 \\ 0 & \text{otherwise} \end{cases}$



Sigmoid function: $h_0 = \frac{1}{1 + e^{-\theta^T x}}$



Sigmoid function is a continuous function so it outputs values between 0 and 1 whereas step function only output 0 or 1 when $\theta^T x$ is small, the graph resembles more like a straight line increasing from 0 to 1. As $\theta^T x$ increases, it resembles more and more as a step function.

7) $P(y=1|x) > 1 \rightarrow c_1$ let c_0 and c_1
 $P(y=0|x) < 1 \rightarrow c_0$ be 2 labels

$$\log \left(\frac{P(y=1|x)}{P(y=0|x)} \right) > 0 \rightarrow c_1$$

$$< 0 \rightarrow c_0$$

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

↑ model
as a linear function

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

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

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

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

↳ Sigmoid

8)

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

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

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

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

9)

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

After initially guessing the parameters (θ), take the gradient at the solution and subtract it from the parameter. The size of the update is controlled by the learning rate.

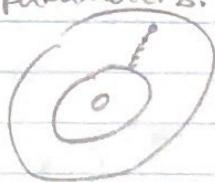
10)

Stop : $(J(\theta^{(i+1)}) - J(\theta^{(i)})) < \gamma$ where γ is the threshold. It should use the loss change since the goal is to minimize loss, regardless of the change in parameters.

11)



Too large

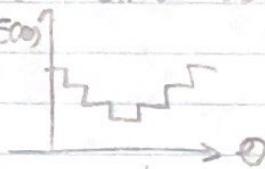


Too small

If the learning rate is too large, the function might not converge or if it does, it might skip the global minimum and settle for a local maximum.

If the learning rate is too small, the steps taken by the function will be small and time consuming. The model might settle at a local minima and the value of J could be reached prematurely.

- incorrect labels
- 12) Empirical error loss is computed by summing the number of errors at a given parameter. Since we are using count as a feature, the graph would be discrete and would resemble something as such :



Therefore, the gradient at a certain point could be zero and therefore, the formula of gradient descent would be unusable ($\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta \nabla J(\theta^{(i)})$)

$$\begin{aligned}
 13) \quad l(\theta) &= \log \left[\prod_{\substack{x^{(i)} \in C_1 \\ x^{(i)} \in C_0}} p(y=1|x^{(i)}) \right] \\
 &\quad \times \log \left[\prod_{\substack{x^{(i)} \in C_0 \\ x^{(i)} \in C_1}} p(y=0|x^{(i)}) \right] \\
 &= \log \prod_{i=1}^m p(y^{(i)}=1|x^{(i)})^{y^{(i)}} p(y^{(i)}=0|x^{(i)})^{1-y^{(i)}} \\
 &= \sum_{i=1}^m y^{(i)} \log \left(\frac{p(y=1|x^{(i)})}{p(y=0|x^{(i)})} \right) + (1-y^{(i)}) \log \left(1 - \frac{p(y=1|x^{(i)})}{p(y=0|x^{(i)})} \right) \\
 l(\theta) &= \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)}) + (1-y^{(i)}) \log(1 - h_\theta(x^{(i)})) \\
 &\quad \uparrow \text{Binary Cross-entropy.}
 \end{aligned}$$

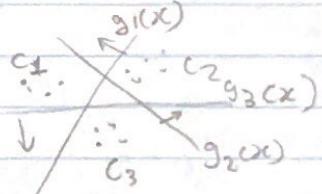
$$\begin{aligned}
 14) \quad \frac{\partial}{\partial \theta} (-l(\theta)) &= \frac{\partial}{\partial \theta} \sum_{i=1}^m y^{(i)} \log(h_\theta(x)) + (1-y^{(i)}) \log(1-h_\theta(x)) \\
 &= \sum_{i=1}^m y^{(i)} (1-h_\theta(x^{(i)})) - \sum_{i=1}^m (1-y^{(i)}) (-h_\theta(x^{(i)})) x^{(i)} \\
 &= \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x^{(i)}
 \end{aligned}$$

minimize -ve log likelihood:

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

The new parameters are equal to $\theta^{(i)}$ minus the gradient of negative log likelihood.

15) One against all others :

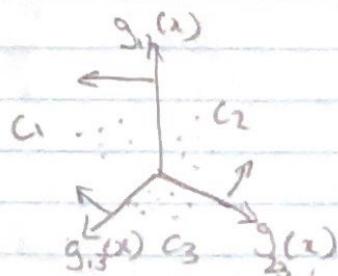


It has K discriminant functions and the label \hat{y} , is calculated by $\hat{y} = \operatorname{argmax}_j g_j(x)$

One against each others:

It has $\frac{1}{2}K(K+1)$ discriminant

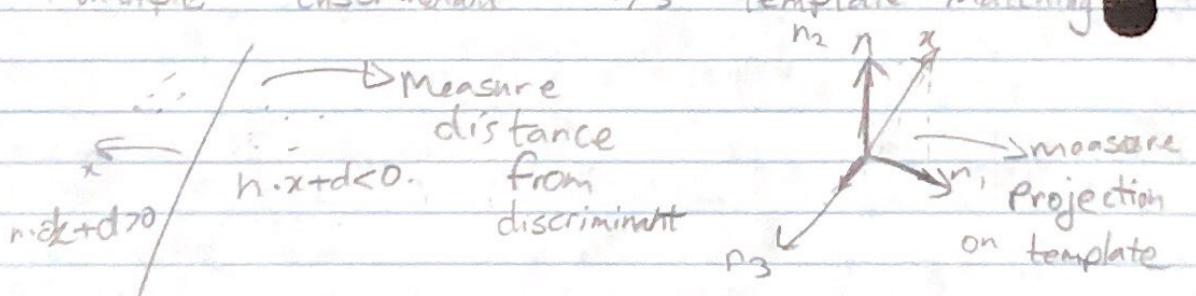
functions and the label, \hat{y} , is calculated by $\hat{y} = \operatorname{argmax}_i \sum_{d \neq i} g_{id}(x)$



One against each other is easier to compare but it is easier to discriminate the data is one against all others.

(6) The rows of Θ^T are templates. With K rows of Θ^T , we have K templates. $\Theta^T x$ measures how well x matches each of the K template. Each class has a template and x is labelled with the template with whom it has the highest similarities. In multiple linear discriminant, Θ^T are parameters of K linear discriminant functions and each function separates one class from all others. Each class has a linear discriminant function. An input x is passed through every function and x is labelled with the function outputting the greatest value (largest distance from discriminant).

multiple discriminant v/s template matching



(7) Softmax converts the output to the probability that a certain input belongs to a certain class. The sum of all output after softmax will be equal to 1. The input will be classified with the class which has highest probability for this input.

$$18) \text{ softmax} = \frac{\exp(\theta_j^T x)}{\sum_{i=1}^k \exp(\theta_i^T x)}$$

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

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

$$19) l(\theta) = -\log \left[\prod_{i=1}^m \prod_{j=1}^k p(y^{(i)} = j | x^{(i)}) \right] \stackrel{-\text{ve log likelihood}}{\rightarrow}$$

$$= -\sum_{i=1}^m \sum_{j=1}^k I(y^{(i)} = j) \log p(y^{(i)} = j | x^{(i)}) \stackrel{y_j^{(i)}}{\rightarrow} \stackrel{g_j^{(i)}}{\rightarrow}$$

$$= -\sum_{i=1}^m \sum_{j=1}^k I(y^{(i)} = j) \log h\theta_j(x^{(i)}) \stackrel{\text{categorical cross-entropy}}{\rightarrow}$$

$$20) \frac{\partial}{\partial \theta_j} (-l(\theta)) = \sum_{i=1}^m (h\theta_j(x^{(i)}) - I(y^{(i)} = j)) x^{(i)}$$

$$\theta_j \leftarrow \theta_j + n \sum_{i=1}^m (h\theta_j(x^{(i)}) - I(y^{(i)} = j)) x^{(i)}$$

↓ new parameter ↑ initial parameter ↑ learning rate ↓ gradient of
 -ve log likelihood

After training, we have $\theta_1, \dots, \theta_k$

$$\text{and } \hat{y} = \underset{j}{\operatorname{argmax}} \theta_j^T x.$$

Neural networks

- 2) Dimensionality increase can be useful to allow us to find some relationships that are only available in higher dimension space. It could also allow the use of linear discriminants.
- 3) Dimensionality decrease is used to condense the data and extract only the useful features.
- 4) The hidden layer may not have the same number of inputs as the output layer.
Since w_0 and y_p are initially unknown we have to guess the parameters of w_j and b_j until we get a satisfactory $E^{(1)}$. Therefore, the same equation cannot be used for both layers

$$\text{loss} \quad E = \frac{1}{n} \sum_{i=1}^n (y^{(1)} - \hat{y}^{(1)})^2$$

$$\text{Grad. rule: } \frac{\partial E}{\partial v} = \frac{\partial E}{\partial g} \frac{\partial g}{\partial v}$$

$$\frac{\partial E}{\partial w_j} = \frac{\partial E}{\partial g} \cdot \frac{\partial g}{\partial z_j} \frac{\partial z_j}{\partial w_j}$$
$$= \sum_{i=1}^n (y^{(1)} - \hat{y}^{(1)}) V_i z_j^{(1)} (1 + z_j^{(1)}) z_j^{(1)}$$

Let K be number of outputs.

$$5) \text{ Loss: } E(\sum_{i=1}^m \sum_{j=1}^K (y_i^{(j)} - \hat{y}_i^{(j)})^2)$$

$$\text{Chain rule: } \frac{\partial E}{\partial v_j} = \frac{\partial E}{\partial \hat{y}_j} \frac{\partial \hat{y}_j}{\partial v_j}$$

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

$$\frac{\partial E}{\partial w_j} = \sum_{i=1}^m \frac{\partial E}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_i} \frac{\partial z_i}{\partial w_j}$$

$$\sum_{i=1}^m \sum_{l=1}^K (\hat{y}_l^{(i)} - y_l^{(i)}) v_{lj} z_j^{(i)} (1 - z_j^{(i)}) z_l^{(i)}$$

$$6) \text{ Loss: } L(\theta) = \sum_{i=1}^m p(y=1|x^{(i)})^{y^{(i)}} \frac{(1-p(y=1|x^{(i)}))^{1-y^{(i)}}}{p(y=0|x^{(i)})}$$

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

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

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

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

7) Let K be number of classes:

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

The gradient is the same as for regression but with different activation. (output activation is softmax for multiclass classification network)

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

$$\frac{\partial E}{\partial w_j} = \sum_{i=1}^m \frac{\partial E}{\partial \hat{y}_j} \frac{\partial \hat{y}_j}{\partial z_j} \frac{\partial z_j}{\partial w_j}$$

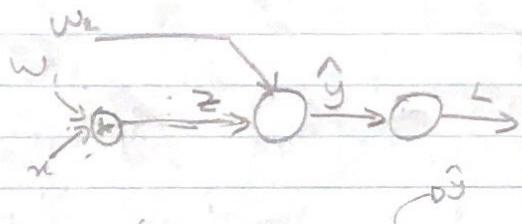
$$\rightarrow \sum_{i=1}^m \sum_{l=1}^k (\hat{y}_j^{(i)} - y_l^{(i)}) v_{lj} z_j^{(i)} (1 - z_j^{(i)}) z^{(i)}$$

- 8) Since the initial output of the hidden layer is unknown as well as $\sum w_j z_j$ and $\sum v_j z_j$, the approach is to guess $\sum w_j z_j$ and $\sum v_j z_j$ to compute $\hat{z}^{(i)}$. Then modify the parameters until a satisfactory $\hat{z}^{(i)}$ is computed.

Computation graphs

2) The advantage of computational graph is that each node is simple and has a simple explicit function for derivatives. During the forward pass, the functions at all intermediate nodes are computed using the input. During the backward pass, the gradient of each node with respect to the input is calculated, starting from the end node. Each node should be able to compute the gradient to update the parameters and it should store the value computed during the forward pass to be used when calculating the gradient.

2)



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

$$\hat{y} = f_2(w_2, z)$$

$$z = f_1(w_1, x)$$

$$i) L = L_2(f_2(w_2, f_1(w_1, x)), y) = (f_2(w_2, f_1(w_1, x)) - y)^2$$

$$\frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y) \quad \frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_2}, \quad \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_1}$$

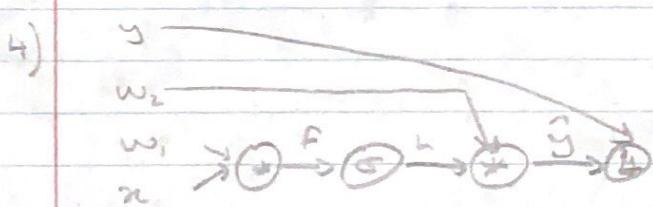
$$ii) L = \text{Cross}(f_2(w_2, f_1(w_1, x)), y) = \prod_{j=1}^m \prod_{k=1}^K \left(p(y=j|z^{(k)}) \right)$$

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

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z},$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_1}$$

3)



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

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

$$h = \sigma(n w_1)$$

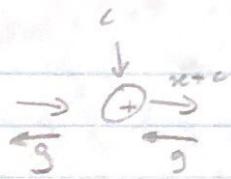
$$\frac{\partial L}{\partial \hat{y}} = 2(y - \hat{y}) \quad \frac{\partial \hat{y}}{\partial h} = w_2 \Rightarrow \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h} = \frac{\partial L}{\partial \hat{y}} w_2$$

$$\frac{\partial h}{\partial w_2} = 1 \Rightarrow \frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h} \frac{\partial h}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot h.$$

$$\frac{\partial h}{\partial w_1} = h(1-h)n \Rightarrow \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h} \frac{\partial h}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} w_2 h(1-h)n$$

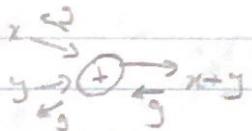
5)

addition of a constant: $x \rightarrow \textcircled{+} \rightarrow x+c$



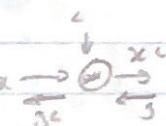
$$\frac{\partial}{\partial x}(x+c) = 1$$

addition of 2 inputs: $x \rightarrow \textcircled{+} \rightarrow x+y$



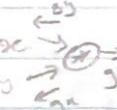
$$\frac{\partial}{\partial x}(x+y) = \frac{\partial}{\partial y}(x+y) = 1$$

multiplication of a constant: $x \rightarrow \textcircled{\times c} \rightarrow xc$



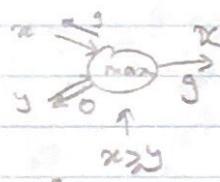
$$\frac{\partial}{\partial x}(xc) = c$$

multiplication of 2 inputs: $x \rightarrow \textcircled{\times y} \rightarrow xy$



$$\frac{\partial}{\partial x}(xy) = y \quad \frac{\partial}{\partial y}(xy) = x$$

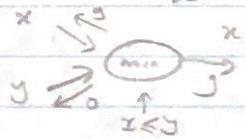
max of 2 inputs:



$$\frac{\partial}{\partial x}(\max(x,y)) = \begin{cases} 1 & x > y \\ 0 & x \leq y \end{cases}$$

$$\frac{\partial}{\partial y}(\max(x,y)) = \begin{cases} 0 & x > y \\ 1 & x \leq y \end{cases}$$

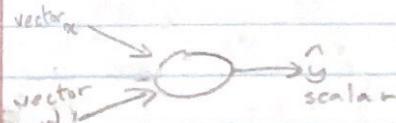
min of 2 inputs:



$$\frac{\partial}{\partial x}(\min(x,y)) = \begin{cases} 1 & x < y \\ 0 & x \geq y \end{cases}$$

$$\frac{\partial}{\partial y}(\min(x,y)) = \begin{cases} 0 & x < y \\ 1 & x \geq y \end{cases}$$

6)



$\frac{\partial g}{\partial w} \rightarrow \text{gradient (vector)}$

$\frac{\partial g}{\partial x} \rightarrow \text{gradient (vector)}$

The derivative will be a vector

7)

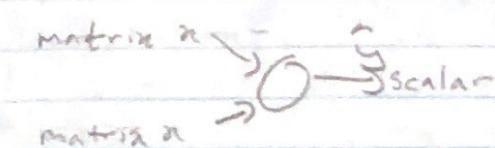


$\frac{\partial g}{\partial w} \rightarrow \text{Rank 2 tensor}$

$\frac{\partial g}{\partial x} \rightarrow \text{Rank 2 tensor}$

The derivative will be a Rank 2 tensor

8)



$\frac{\partial g}{\partial w} \rightarrow \text{Rank 2 tensor}$

$\frac{\partial g}{\partial x} \rightarrow \text{Rank 2 tensor}$

The derivative will be a Rank 2 tensor

$$9) f(x) = f(x_1, \dots, x_n) = \begin{bmatrix} f_1(x_1, \dots, x_n) \\ \vdots \\ f_m(x_1, \dots, x_n) \end{bmatrix}$$

Jacobian matrix

$$\nabla f(x) = \begin{bmatrix} \nabla f_1 \\ \vdots \\ \nabla f_m \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$$(F \circ g)(x) = f(g(x))$$

$$((F \circ g)(x))' = f'(g(x))g'(x)$$

$$\begin{array}{c} F \rightarrow G \\ \frac{\partial F}{\partial x} \rightarrow \frac{\partial G}{\partial x} \\ \frac{\partial F}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x} \end{array}$$

$$10) f(x, y, z) = \begin{bmatrix} 3xy \\ y-z \end{bmatrix} \quad g(x, y, z) = \begin{bmatrix} x-5y \\ xy \\ x+y \end{bmatrix}$$

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

$$\nabla f(x, y, z) = \begin{bmatrix} 3y & 3x & 0 \\ 0 & 1 & -1 \end{bmatrix} \quad \nabla g(x, y) = \begin{bmatrix} 1 & -5 \\ y & x \\ 1 & 1 \end{bmatrix}$$

$$\begin{aligned} \nabla(F \circ g)(x, y) &= \begin{bmatrix} \frac{\partial}{\partial x}(3x^2y - 15xy^2) & \frac{\partial}{\partial y}(3x^2y - 15xy^2) \\ \frac{\partial}{\partial y}(xy - x+y) & \frac{\partial}{\partial y}(xy - x+y) \end{bmatrix} \\ &= \begin{bmatrix} 6xy - 15y^2 & 3x^2 - 30xy \\ y-1 & x-1 \end{bmatrix} \leftarrow \end{aligned}$$

using chain rule:

$$\text{Plug } g(x, y) \text{ in } \nabla f(g(x, y)) = \begin{bmatrix} 3xy & 3(x-5y) & 0 \\ 0 & 1 & -1 \end{bmatrix}$$

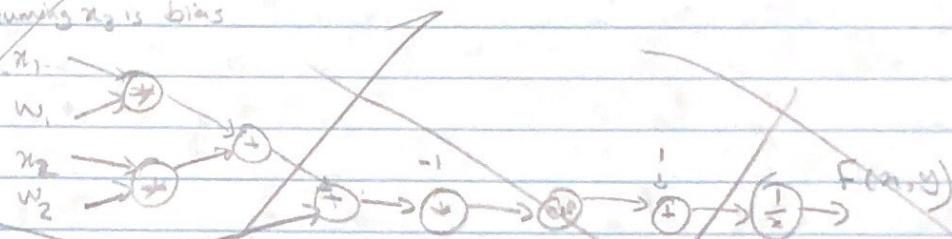
$$\begin{aligned} \text{Chain rule: } \nabla f(g(x, y)) \nabla g(x, y) &= \begin{bmatrix} 3xy & 3(x-5y) & 0 \\ 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & -5 \\ y & x \\ 1 & 1 \end{bmatrix} \\ &= \begin{bmatrix} 6xy - 15y^2 & 3x^2 - 30xy \\ y-1 & x-1 \end{bmatrix} \leftarrow \end{aligned}$$

11) In the forward pass, the nodes compute function values and these values can be the input to another function. For example, if node a, is an input to node b, we have to ensure that node a is computed first to be able to compute node b, hence the need of ordering the nodes.

In the backward pass, the inputs become the outputs and vice versa. Instead of computing the function we compute derivatives and pass back the output to previous nodes. The nodes need to be in reverse order of the forward pass to ensure that all derivatives can be computed.

Assuming w_0 is bias

12) a)



b) f

Assuming $x_0 = 1$ for x_3, y_3

$$z_1 = [w_0]^T [x_0 \dots x_2] = [0.01, 0.02, 0.03] \begin{bmatrix} 1 \\ x_2 \end{bmatrix} = 0.11$$

$$z_2 = [w_1]^T [x_0 \dots x_2] = [0.03, 0.01, 0.02] \begin{bmatrix} 1 \\ x_2 \end{bmatrix} = 0.09$$

$$z_3 = [w_2]^T [x_0 \dots x_2] = [0.02, 0.03, 0.01] \begin{bmatrix} 1 \\ x_2 \end{bmatrix} = 0.08$$

$$A_1 = \text{sigmoid}(z_1) = 0.5274$$

$$A_2 = \text{sigmoid}(z_2) = 0.5224$$

$$A_3 = \text{sigmoid}(z_3) = 0.5199$$

$$\hat{y}_3 = v^T A = [0.01, 0.02, 0.03, 0.04] \begin{bmatrix} 1 \\ A_1 \\ A_2 \\ A_3 \end{bmatrix} = 0.0561$$

$$L_3 = \frac{1}{2} (\hat{y}_3 - y_3)^2 = 49.432 \leftarrow$$

c) $\frac{\partial L}{\partial w_1} = (\hat{y} - y) \cdot v_1 \cdot (A)(1-A)x$
 $\frac{\partial L}{\partial v_1} = (\hat{y} - y) \cdot z$

	x_1, y_1	x_2, y_2	x_3, y_3
$\frac{\partial L}{\partial w_1}$	$-0.0397, -0.0397, -0.0793$	$-0.055, -0.055, -0.164$	$-0.0496, -0.0496$
$\frac{\partial L}{\partial w_2}$	$-0.600, 0.600, -0.119$	$-0.4082, -0.082, -0.243$	$-0.0715, -0.193, -0.16$
$\frac{\partial L}{\partial w_3}$	$-0.080, -0.070, -0.159$	$-0.09, -0.09, -0.26$	$-0.99, -0.99, -0.98$
$\frac{\partial L}{\partial v}$	$-7.94, -4.19, -4.13, -4.12$	$-10.9, 5.80, -5.75, -5.76$	$-9.95, 5.25, -5.25, -5.25$

d) 1' d im. mit der vidi

$$13) a) \nabla f(x,y) = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} = \begin{bmatrix} \frac{4(2x+3y)}{6(2x+3y)} \\ 8x+12y \\ 12x+18y \end{bmatrix}$$

$$b) f(x,y) = \begin{bmatrix} g_1(x,y) \\ g_2(x,y) \end{bmatrix} \quad g_1(x,y) = x^2 + 2y \\ g_2(x,y) = 3x + 4y^2$$

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

$$\Delta F(1,2) = \begin{bmatrix} 2 \\ \frac{2}{3} \\ 16 \end{bmatrix}$$

$$c) f'(x,y) = \begin{bmatrix} 2x \\ \frac{2x}{3} \\ 8y \end{bmatrix} \quad g'(x) = \begin{bmatrix} 1 \\ 2x \end{bmatrix} \rightarrow II$$

$$f'(g(x)) = \begin{bmatrix} 2x \\ 3 \\ 8x^2 \end{bmatrix} - I$$

$$II \times I = \begin{bmatrix} 2x + 4x \\ 3 + 16x^3 \end{bmatrix} = \begin{bmatrix} 6x \\ 3 + 16x^3 \end{bmatrix} \Rightarrow D(f \circ g_1)(x)$$

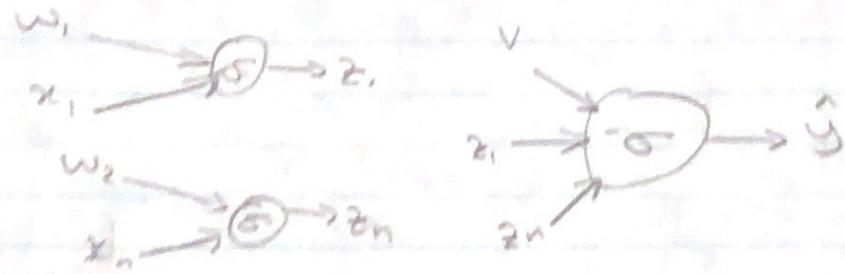
$$\text{for } x=2 \quad \nabla(f \circ g_1)(2) = \begin{bmatrix} 12 \\ 131 \end{bmatrix}$$

$$d) F \cdot g = \begin{bmatrix} x^2 + 2x^2 \\ 3x + 4(x^2)^2 \end{bmatrix} = \begin{bmatrix} 3x^2 \\ 3x + 4x^4 \end{bmatrix}$$

$$\nabla(F \cdot g) = \begin{bmatrix} 6x \\ 3 + 16x^3 \end{bmatrix}$$

$$\nabla(F \cdot g)(2) = \begin{bmatrix} 12 \\ 3 + 16(8) \end{bmatrix} = \begin{bmatrix} 12 \\ 131 \end{bmatrix}$$

p1)



$$z_i = h(w_i^T x)$$

$$\text{then } \hat{y}_i = \sigma(y_i + \alpha) = h(w_i^T \cdot z)$$

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

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

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

$$\frac{\partial \hat{y}}{\partial v} = h(v^T z)(1 - h(v^T z)) \cdot z$$

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

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

$$v_i = v_i - \eta \left(\frac{\partial L}{\partial v} \right)$$

$$v_i = v_i - \eta \left(- \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) \cdot h(v_i^T \cdot z^{(i)}) \cdot (1 - h(v_i^T \cdot z^{(i)})) \right)$$

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial v} \cdot \frac{\partial v}{\partial w_i} \quad \frac{\partial \hat{y}}{\partial v} = \hat{y}_i \cdot (1 - \hat{y}_i) \cdot v$$

$$z_i = h(w_i^T \cdot x) \quad \frac{\partial v}{\partial w_i} = h(w_i^T \cdot x) \cdot (1 - h(w_i^T \cdot x)) \cdot x_i$$

$$\frac{\partial L}{\partial w_i} = \sum_{i=1}^m \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{(1 - \hat{y}_i)}{\partial v} \cdot \frac{\partial v}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_i} = \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2 \cdot \hat{y}_i \cdot (1 - \hat{y}_i) \cdot v_i \cdot h'(w_i^T \cdot x^{(i)}) \cdot x_i$$

$$w_i = w_{i-1} + \eta \left(- \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2 \cdot \hat{y}_i \cdot (1 - \hat{y}_i) \cdot v_i \cdot h'(w_i^T \cdot x^{(i)}) \cdot x_i \right)$$