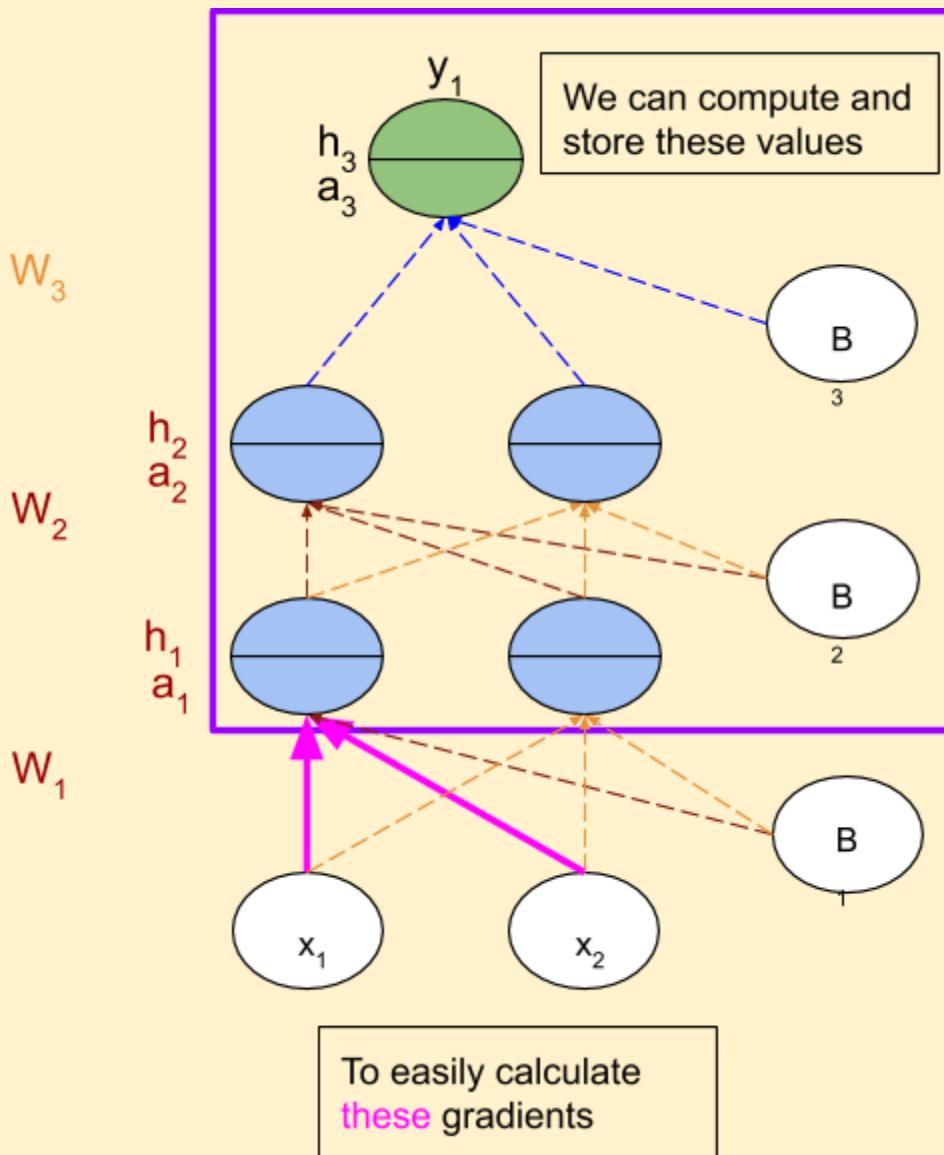


Backpropagation (Math-heavy/Vectorized)

Setting the context

How does this differ from the previous section

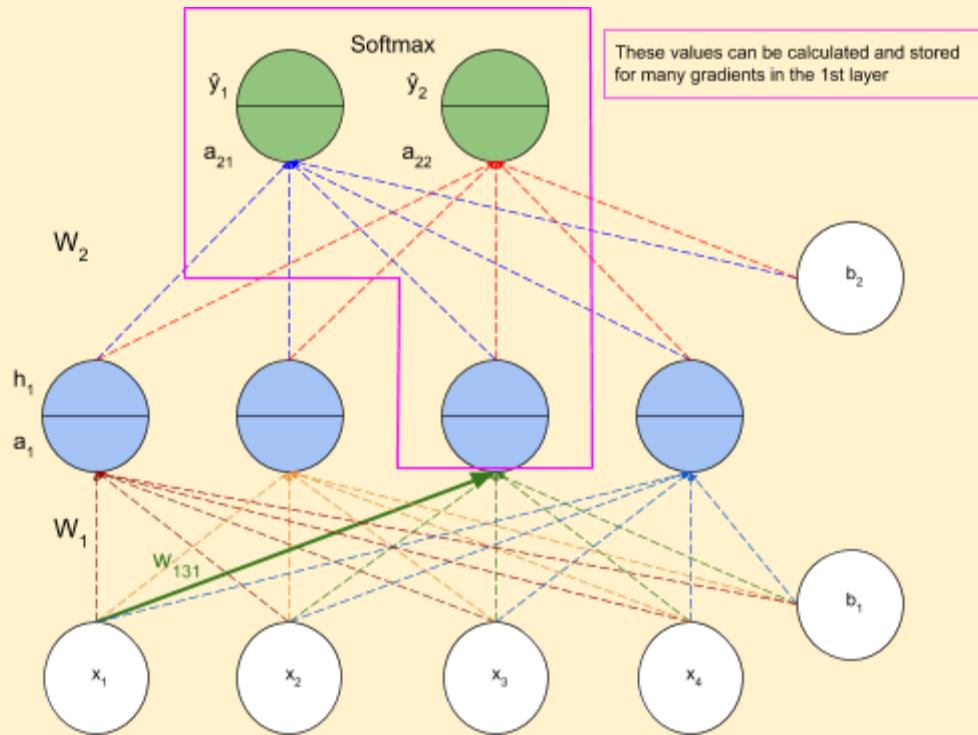
1. We've looked at two different levels of complexity to the backpropagation algorithm so far
 - a. No-math: A simple forward pass with no gradient calculation
 - b. Light-math: Gradient calculation for each weight using chain rule
2. Now, with the Heavy-math version, our objective is to identify common calculations between different weights and re-use them to make our work simpler



PadhAI: Backpropagation - the full version

One Fourth Labs

3. Let us consider the example from the light-math backpropagation chapter



4. Consider $d\mathbf{w}_{131}$ or $\frac{\partial L}{\partial w_{131}}$

- Here, we are certain about using the highlighted values for gradient computation. So we can pre-calculate and store them
- In the outermost layer

$$\frac{\partial L}{\partial a_2} = \begin{bmatrix} \frac{\partial L}{\partial a_{21}} \\ \frac{\partial L}{\partial a_{22}} \end{bmatrix}$$

5. Similarly, storing the values of $\frac{\partial h_1}{\partial a_2}$ will prove useful down the line

One Fourth Labs

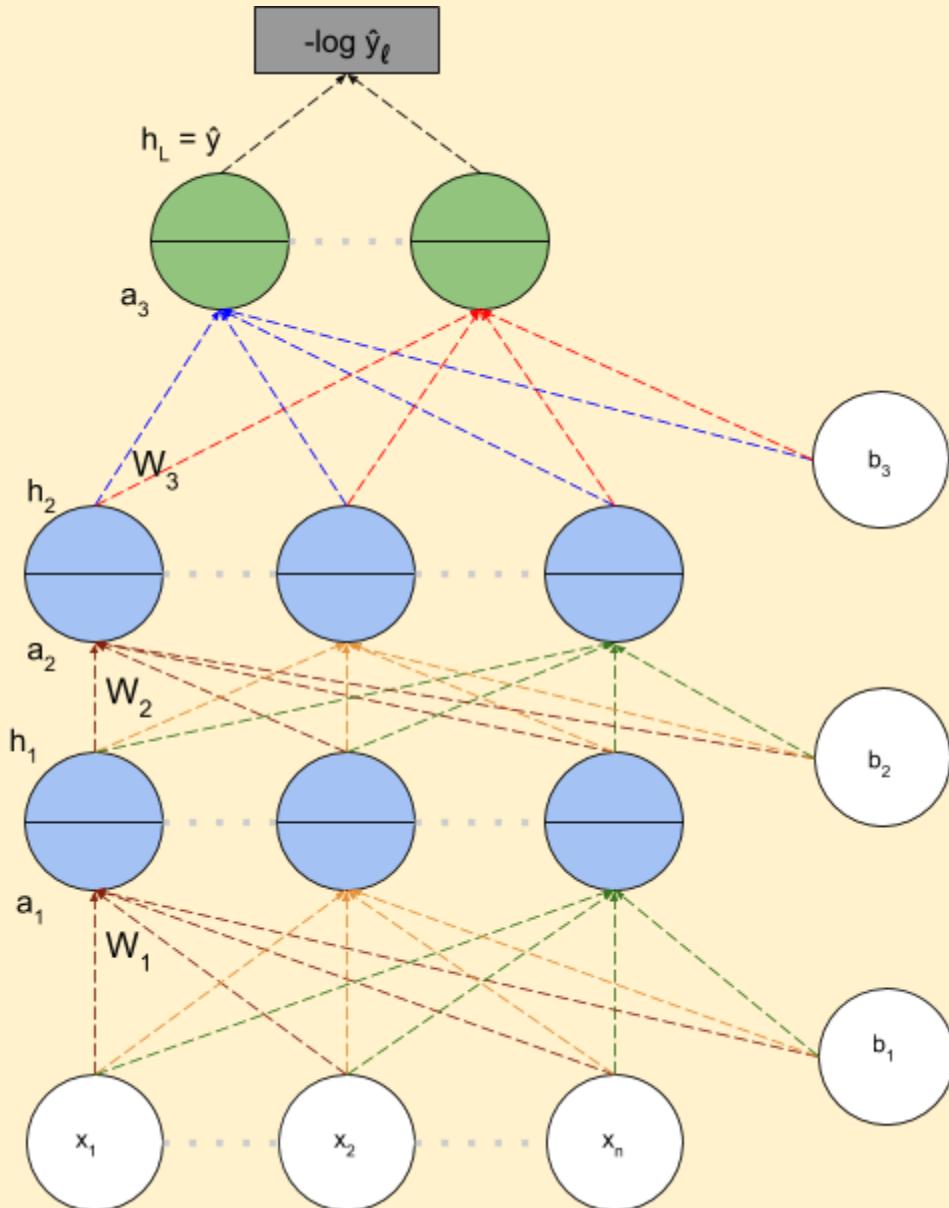
Intuition behind backpropagation

Let us look at an intuitive explanation of backpropagation

1. Let's have a quick calculus recap

- a. $f(x) = x^2$, *Derivative* : $\frac{df(x)}{dx}$
- b. $f(x) = x^2 + y^2$, *Partial Derivatives* : $\frac{\partial f(x)}{\partial x}$ and $\frac{\partial f(x)}{\partial y}$ a
- c. $f(\theta) = [x, y]$, *Gradient* $\nabla_{\theta} = [\frac{\partial f(x)}{\partial x}, \frac{\partial f(x)}{\partial y}]$

2. Consider the following sample network

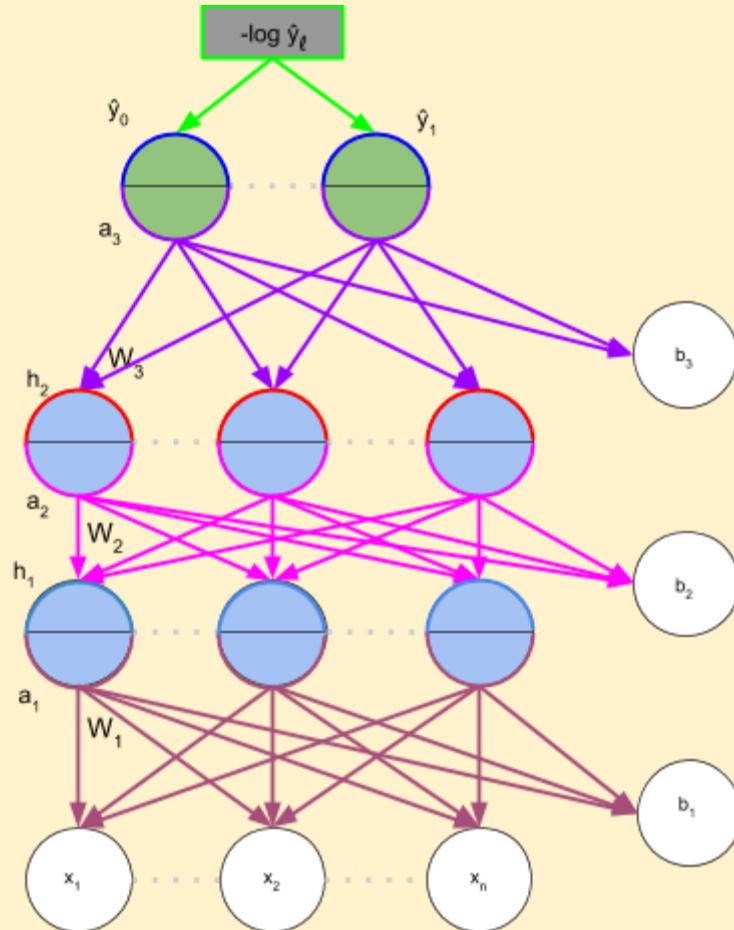


3. Now, let us assume our network has undergone some training, so we have a Loss value in hand.
4. $W_1, W_2, W_3, b_1, b_2, b_3$ have been updated and reflect accordingly in the Loss function

PadhAI: Backpropagation - the full version

One Fourth Labs

5. Now, we want to scrutinise how each parameter is responsible for the loss. So we move backwards from the Loss in a step-by-step manner



6. Stepwise calculation (while personifying the neurons and layers 😊)
- Step 1: The **loss** talks to the output layer, saying “You better take responsibility for the poor output!”
 - Step 2: The **output activation layer** says, “Hey, I’m simply applying the softmax function to the input given to me by the **output preactivation layer**”
 - Step 3: The **output preactivation layer** says, “I take responsibility for my part, but I am only as good as the hidden layer and the weights below me.” After all $f(x) = \hat{y} = O(W_L h_{L-1} + b_L)$
 - Step 4: The parameters W_3 and b_3 say “It is our mistake, **please update our values**”
 - Step 5: However, the **second hidden activation layer** says “Hey, I’m simply applying the sigmoid function to the input given to me by the **second hidden preactivation layer**”
 - Step 6: The **second hidden preactivation layer** says, “I am only as good as the hidden layer and weights below me”
 - Step 7: The parameters W_2 and b_2 say “It is our mistake, **please update our values**”
 - Step 8: However, the **first hidden activation layer** says “Hey, I’m simply applying the sigmoid function to the input given to me by the **first hidden preactivation layer**”
 - Step 9: The **first hidden preactivation layer** says, “I am only as good as the weights below me, we cannot blame the input layer.”
 - Step 10: The last set of parameters W_1 and b_1 say “It is our mistake, please update us”

PadhAI: Backpropagation - the full version

One Fourth Labs

7. Thus, to arrive at the derivative of the Loss function w.r.t any of the weights, we must proceed downwards from the top. We cannot simply calculate it without knowing the preceding values.
8. Our roadmap for the rest of the module
 - a. To calculate the desired gradient, we need to compute
 - b. Gradient w.r.t output units
 - c. Gradient w.r.t hidden units
 - d. Gradient w.r.t weights and biases
 - e.

$$\frac{\partial L(\theta)}{\partial W_{111}} = \frac{\partial L(\theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_3} \frac{\partial a_3}{\partial h_2} \frac{\partial h_2}{\partial a_2} \frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1} \frac{\partial a_1}{\partial W_{111}}$$

Talk to the weight directly	Talk to the output layer	Talk to the previous hidden layer	Talk to the previous hidden layer	Talk to the weights
		works for any number of output layers		

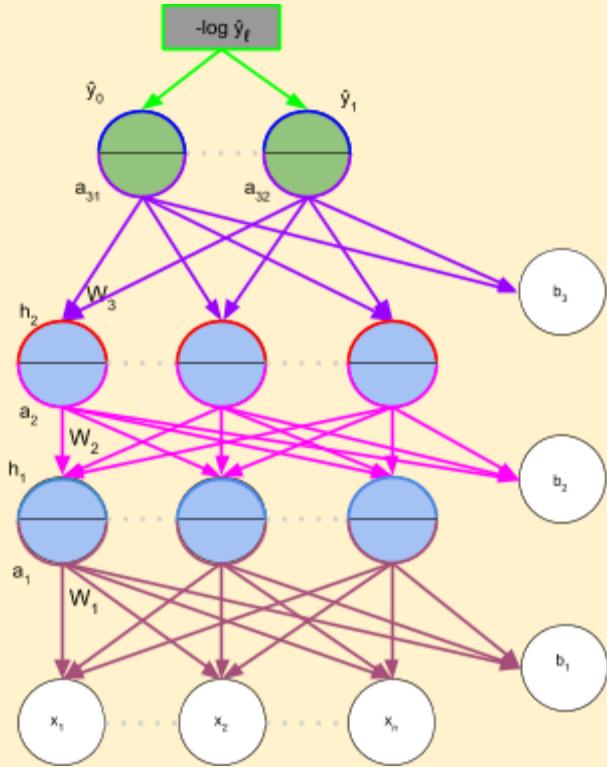
- f. Our aim is to do these calculations for any of the possible weights using notation i, j, k instead of numbers
- g. For the rest of this exercise, our focus is on Cross Entropy loss and Softmax output.
9. To reduce the tediousness of applying the chain rule each time to get the desired gradient, we will look to re-use common values and pathways to more efficiently calculate any gradient.

One Fourth Labs

Understanding the dimensions of gradients

What are we interested in?

1. Consider the backpropagation illustration from the previous section



2. What we are interested in is $\nabla_{a_3} L(\theta) = \frac{\partial L(\theta)}{\partial a_{Li}}$ (where true output $y = 1$, $L = \text{Layer number}$, i is the index of the correct class-label for the given input, and i is the neuron number)
3. We know that $\frac{\partial L(\theta)}{\partial a_{Li}}$ is dependent on a_{31} and a_{32}
4. Therefore, the derivative at the output layer

$$\nabla_{a_3} L(\theta) = \begin{bmatrix} \frac{\partial L}{\partial a_{31}} \\ \frac{\partial L}{\partial a_{32}} \end{bmatrix}$$

5. In the above gradient, $L = 3$ and $i \in \{1, 2\}$
6. Henceforth, we can use these notations in place of numbers to simplify gradient calculation for all possible gradients.

One Fourth Labs

Computing derivatives w.r.t Output Layer

Part 1

The first derivative in the chain

1. What we are actually interested in is:

- a. $\frac{\partial L(\theta)}{\partial a_{Li}} = \frac{\partial(-\log \hat{y}_l)}{\partial a_{Li}}$
 - b. Where L = layer number, i = neuron (from 1 to k), l = index of correct output
 - c. Here, we use the cross entropy loss function
 - d. In the output layer L, assume we have neurons $a_{L1}, a_{L2} \dots a_{Lk}$
 - e. The output layer L involves applying the softmax function to all the neurons
 - f. $\hat{y}_l = \frac{e^{a_{Ll}}}{\sum_i e^{a_{Li}}}$ again, (l refers to the index of the correct output neuron)
 - g. Thus, \hat{y}_l depends on all the neurons' outputs as they all appear in the denominator, thereby making the derivative non-zero for all the output neurons
2. $\frac{\partial L(\theta)}{\partial a_{Li}} = \frac{\partial(-\log \hat{y}_l)}{\partial a_{Li}} = \frac{\partial(-\log \hat{y}_l)}{\partial \hat{y}_l} \frac{\partial \hat{y}_l}{\partial a_{Li}}$
 3. From the previous points, we know that \hat{y}_l depends on a_{Li}
 4. The first part of the derivative is fairly straightforward (of the form $\frac{\partial \log x}{\partial x}$)
 5. $\frac{\partial(-\log \hat{y}_l)}{\partial \hat{y}_l} = \frac{-1}{\hat{y}_l}$

PadhAI: Backpropagation - the full version

One Fourth Labs

Computing derivatives w.r.t Output Layer

Part 2

1. Continuing from where we left off $\frac{\partial L(\theta)}{\partial a_{Li}} = \frac{-1}{\hat{y}_l} \frac{\partial \hat{y}_l}{\partial a_{Li}}$
2. Here, we know that $\hat{y}_l = \frac{e^{a_{Ll}}}{\sum_i e^{a_{Li}}}$ (taking the l-th entry of the softmax function applied to vector a_L)
3. $\frac{\partial L(\theta)}{\partial a_{Li}} = \frac{-1}{\hat{y}_l} \frac{\partial \text{softmax}(a_L)_l}{\partial a_{Li}}$ where $a_L = [a_{L1}, a_{L2} \dots a_{Lk}]$
 - a. Where $\text{softmax}(a_L) = [\frac{e^{a_{L1}}}{\sum_i e^{a_{Li}}}, \frac{e^{a_{L2}}}{\sum_i e^{a_{Li}}} \dots \frac{e^{a_{Lk}}}{\sum_i e^{a_{Li}}}]$
 - b. Selecting the l-th entry would give us the value $\text{softmax}(a_L)_l = \frac{\exp(a_L)_l}{\sum_i \exp(a_L)_i}$
4. $\frac{\partial L(\theta)}{\partial a_{Li}} = \frac{-1}{\hat{y}_l} \frac{\partial}{\partial a_{Li}} \frac{\exp(a_L)_l}{\sum_i \exp(a_L)_i}$
 - a. This is of the form $\frac{g(x)}{h(x)}$ which gives derivatives $\frac{\partial \frac{g(x)}{h(x)}}{\partial x} = \frac{\partial g(x)}{\partial x} \frac{1}{h(x)} - \frac{g(x)}{h(x)^2} \frac{\partial h(x)}{\partial x}$
 - b. Here $g(x) = \exp(a_L)_l$ and $h(x) = \sum_i \exp(a_L)_i$
 - c. Substitute the values and expand the formula
5. $\frac{\partial L(\theta)}{\partial a_{Li}} = \frac{-1}{\hat{y}_l} \left(\frac{\frac{\partial}{\partial a_{Li}} \exp(a_L)_l}{\sum_i \exp(a_L)_i} - \frac{\exp(a_L)_l (\frac{\partial}{\partial a_{Li}} \sum_i \exp(a_L)_i)}{(\sum_i \exp(a_L)_i)^2} \right)$
 - a. Here, consider $\frac{\partial}{\partial a_{Li}} \exp(a_L)_l$, this value is 0 for all values of i : 0 to k except for when i = l
 - b. Thus, we use an indicator variable $\mathbf{1}_{(l=i)} \exp(a_L)_l$ to denote that all other values except i=l resolve to 0
 - c. Now consider $\frac{\partial}{\partial a_{Li}} \sum_i \exp(a_L)_i$, here i' ranges from 1 to k. When taking the derivative, only the index i=i' remains, which is simply a derivative of an exponent.
6. $\frac{\partial L(\theta)}{\partial a_{Li}} = \frac{-1}{\hat{y}_l} \left(\frac{\mathbf{1}_{(l=i)} \exp(a_L)_l}{\sum_i \exp(a_L)_i} - \frac{\exp(a_L)_l}{\sum_i \exp(a_L)_i} \frac{\exp(a_L)_i}{\sum_i \exp(a_L)_i} \right)$
 - a. This is can be rewritten in terms of the softmax function for the different variables
7. $\frac{\partial L(\theta)}{\partial a_{Li}} = \frac{-1}{\hat{y}_l} (\mathbf{1}_{(l=i)} \text{softmax}(a_L)_l - \text{softmax}(a_L)_l \text{softmax}(a_L)_i)$
 - a. We know that the Softmax function is \hat{y} , so we rewrite it.
8. $\frac{\partial L(\theta)}{\partial a_{Li}} = \frac{-1}{\hat{y}_l} (\mathbf{1}_{(l=i)} \hat{y}_l - \hat{y}_l \hat{y}_i)$
9. After cancellation $\frac{\partial L(\theta)}{\partial a_{Li}} = -(\mathbf{1}_{(l=i)} - \hat{y}_i)$

PadhAI: Backpropagation - the full version

One Fourth Labs

Computing derivatives w.r.t Output Layer

Part 3

1. So far, we have derived the partial derivative with respect to the i -th element of layer a_L
- a. $\frac{\partial L(\theta)}{\partial a_{Li}} = -(1_{(l=i)} - \hat{y}_i)$
2. We can now write the gradient w.r.t the vector a_L
3. As we saw earlier, $a_L = [a_{L1}, a_{L2} \dots a_{Lk}]$
4. Going by the indicator variable in step 1, it resolves to 0 for all values of i except for $i = l$
5. Let us assume a scenario where $k = 4$, and $l = 2$
- a. $\frac{\partial L(\theta)}{\partial a_{L1}} = -(0 - \hat{y}_i)$
- b. $\frac{\partial L(\theta)}{\partial a_{L2}} = -(1 - \hat{y}_i)$
- c. $\frac{\partial L(\theta)}{\partial a_{L3}} = -(0 - \hat{y}_i)$
- d. $\frac{\partial L(\theta)}{\partial a_{L4}} = -(0 - \hat{y}_i)$
- e. Here, the indicator variable values as a vector would be
6. The gradient w.r.t a_L is $\nabla_{a_L} =$

$$\nabla_{a_L} = \begin{bmatrix} \frac{\partial L(\theta)}{\partial a_{L1}} \\ \vdots \\ \frac{\partial L(\theta)}{\partial a_{Lk}} \end{bmatrix}$$

$$\nabla_{a_L} = \begin{bmatrix} -(1_{(l=1)} - \hat{y}_i) \\ \vdots \\ -(1_{(l=k)} - \hat{y}_i) \end{bmatrix}$$

7. The above can be seen as a difference of two vectors, $[0, 1, 0, \dots 0_k]$ and \hat{y}
8. The first vector is essentially the one hot representation of the true output $e(l)$: $-(e(l) - \hat{y}_i)$
9. In reality, this is simply the difference between the true distribution y and the predicted distribution \hat{y}
10. $\nabla_{a_L} L(\theta) = -(y - \hat{y}_i)$

PadhAI: Backpropagation - the full version

One Fourth Labs

Quick recap of the story so far

What have we covered up till this point?

1. Our roadmap for this module
 - a. To calculate the desired gradient, we need to compute
 - b. Gradient w.r.t output units
 - c. Gradient w.r.t hidden units
 - d. Gradient w.r.t weights and biases
 - e.

$\frac{\partial L(\theta)}{\partial W_{111}}$	$=$	$\frac{\partial L(\theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_3}$	$\frac{\partial a_3}{\partial h_2} \frac{\partial h_2}{\partial a_2}$	$\frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1}$	$\frac{\partial a_3}{\partial W_{111}}$	
Talk to the weight directly	Talk to the output layer	Talk to the previous hidden layer	Talk to the previous hidden layer	Talk to the weights	works for any number of output layers	

- f. For the rest of this exercise, our focus is on Cross Entropy loss and Softmax output.
2. Here, what the sections highlighted in green are what we have covered so far, i.e. the derivative with respect to the last layer a_L
3. The gradient was calculated to be $\nabla_{a_L} L(\theta) = \frac{\partial L(\theta)}{\partial a_{L_i}} = -(y - \hat{y}_i)$

PadhAI: Backpropagation - the full version

One Fourth Labs

Computing derivatives w.r.t Hidden Layers

Part 1

The derivatives corresponding to the hidden layers

- What we are interested in is

- a. $\frac{\partial L(\theta)}{\partial h_{ij}} = \sum_{m=1}^k \frac{\partial L(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}}$
- b. This formula is the summation of all the paths that lead from the concerned neuron to the loss function
- c. Here, i = layer number, m = neuron number for a , j = neuron number for h
- d. From the previous section, we already know how to compute $\frac{\partial L(\theta)}{\partial a_{i+1,m}}$ so we need to only focus on $\frac{\partial a_{i+1,m}}{\partial h_{ij}}$
- e. However, when we compute the derivative of the neuron $a_{i+1,m}$ w.r.t h_{ij} we are left with the weight component $W_{i+1,m,j}$
- f. This refers to the weight component between the output neuron($a_{i+1,m}$) and input neuron (h_{ij})

- Thus we have $\frac{\partial L(\theta)}{\partial h_{ij}} = \sum_{m=1}^k \frac{\partial L(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}$

- Now consider these two vectors

a.

$$\nabla_{a_{i+1}} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial a_{i+1,1}} \\ \vdots \\ \frac{\partial L(\theta)}{\partial a_{i+1,k}} \end{bmatrix}$$

$$W_{i+1,\cdot,j} = \begin{bmatrix} W_{i+1,1,j} \\ \vdots \\ W_{i+1,k,j} \end{bmatrix}$$

- b. Here, $\nabla_{a_{i+1}} L(\theta)$ refers to the gradient vector of the loss function w.r.t to all output neurons from $a_{i+1,1}$ to $a_{i+1,k}$
 - c. And $W_{i+1,\cdot,j}$ refers to all rows of the j -th column of the W_{i+1} matrix, ie a vector.
- The dot product of these two vectors is $(W_{i+1,\cdot,j})^T \nabla_{a_{i+1}} L(\theta) = \sum_{m=1}^k \frac{\partial L(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}$

PadhAI: Backpropagation - the full version

One Fourth Labs

5. Here, the RHS is the same as the value from step 2. Therefore, the derivative of the loss function with respect to the hidden layers is the dot-product between the gradient of loss w.r.t output layer and the corresponding weights.

PadhAI: Backpropagation - the full version

One Fourth Labs

Computing derivatives w.r.t Hidden Layers

Part 2

1. We have $\frac{\partial L(\theta)}{\partial h_{ij}} = (W_{i+1, \cdot, j})^T \nabla_{a_{i+1}} L(\theta)$
 - a. This is with respect to one neuron
 - b. We would like to speed up this computation by solving all the derivatives in one go
2. We can now write the gradient w.r.t h_i
 - a.

$$\nabla_{h_i} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial a_{h_i 1}} \\ \vdots \\ \frac{\partial L(\theta)}{\partial a_{h_i n}} \end{bmatrix}$$

$$\nabla_{h_i} L(\theta) = \begin{bmatrix} (W_{i+1, \cdot, 1})^T \nabla_{a_{i+1}} L(\theta) \\ \vdots \\ (W_{i+1, \cdot, n})^T \nabla_{a_{i+1}} L(\theta) \end{bmatrix}$$

- b. Can be written more compactly as $(W_{i+1})^T \nabla_{a_{i+1}} L(\theta)$
3. Thus, the formula for gradient of loss function for the last hidden layer before the output layer is given by $\nabla_{h_i} L(\theta) = (W_{i+1})^T \nabla_{a_{i+1}} L(\theta)$
4. This calculates the gradient w.r.t all neurons of layer i . It uses simple matrix-vector multiplication to achieve this.
5. Now, we have seen a special case applied to the last hidden layer. We must figure out how to make this formula applicable for any generic hidden layer.

One Fourth Labs

Computing derivatives w.r.t Hidden Layers

Part 3

1. Consider the next layer a_i

a. $\frac{\partial L(\theta)}{\partial a_{ij}} = \frac{\partial L(\theta)}{\partial h_{ij}} \frac{\partial h_{ij}}{\partial a_{ij}}$

- b. The first derivative is what we computed in part 2

c. We need to compute the second derivative $\frac{\partial h_{ij}}{\partial a_{ij}}$

- d. We know that h_{ij} is simply the application of an activation function (sigmoid, tanh etc) to a_{ij}

e. So it can be rewritten as $\frac{\partial h_{ij}}{\partial a_{ij}} = g'(a_{ij})$ where $h_{ij} = g(a_{ij})$ and $g'(a_{ij})$ is its derivative

2. $\frac{\partial L(\theta)}{\partial a_{ij}} = \frac{\partial L(\theta)}{\partial h_{ij}} g'(a_{ij})$

3. The full gradient can be written as

a.

$$\nabla_{a_i} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial h_{i1}} g'(a_{i1}) \\ \vdots \\ \frac{\partial L(\theta)}{\partial h_{in}} g'(a_{in}) \end{bmatrix}$$

- b. This vector is the element-wise product of two vectors $\nabla_{h_i} L(\theta)$ and $[..., g'(a_{ik}), ...]$ (which is a vector of derivations of the activation function w.r.t the pre-activation layer. They are both vectors of n-terms)

4. Thus $\nabla_{a_i} L(\theta) = \nabla_{h_i} L(\theta) \odot [..., g'(a_{ik}), ...]$ (\odot refers to element-wise multiplication)

5. This formula can be applied to any of the hidden layers

PadhAI: Backpropagation - the full version

One Fourth Labs

Computing derivatives w.r.t one weight in any layer

1. Our roadmap for this module
 - a. To calculate the desired gradient, we need to compute
 - b. Gradient w.r.t output units
 - c. Gradient w.r.t hidden units
 - d. Gradient w.r.t weights and biases
 - e.

$\frac{\partial L(\theta)}{\partial W_{111}}$	$=$	$\frac{\partial L(\theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_3}$	$\frac{\partial a_3}{\partial h_2} \frac{\partial h_2}{\partial a_2}$	$\frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1}$	$\frac{\partial a_3}{\partial W_{111}}$
Talk to the weight directly	Talk to the output layer	Talk to the previous hidden layer	Talk to the previous hidden layer	Talk to the weights	
works for any number of output layers					

- f. For the rest of this exercise, our focus is on Cross Entropy loss and Softmax output.
2. Here, what the sections highlighted in green are what we have covered so far, i.e. the derivative with respect to the last hidden layer and all other subsequent hidden layers
3. The gradients were calculated to be
 - a. $\nabla_{h_i} L(\theta) = (W_{i+1})^T \nabla_{a_{i+1}} L(\theta)$
 - b. $\nabla_{a_i} L(\theta) = \nabla_{h_i} L(\theta) \odot [..., g'(a_{ik}), ...]$
4. Now, we will be computing the derivative of the loss function w.r.t weights and biases.
5. Recall that $a_k = b_k + W_k h_{k-1}$
6. $\frac{\partial a_{ki}}{\partial W_{kij}} = h_{k-1,j}$
 - a. k = layer number
 - b. i = current layer neuron number
 - c. j = previous/input layer neuron number
7. Now $\frac{\partial L(\theta)}{\partial W_{kij}} = \frac{\partial L(\theta)}{\partial a_{ki}} h_{k-1,j}$
8. We can use this to update the Weight by Gradient Descent
9. $W_{kij} = W_{kij} - \eta \frac{\partial L(\theta)}{\partial W_{kij}}$
10. In the next step, we will look at updating all the weights in a layer simultaneously.

PadhAI: Backpropagation - the full version

One Fourth Labs

Computing derivatives w.r.t all weights in any layer

- Let's take a simple example of a $W_k \in \mathbb{R}^{3 \times 3}$ and see what each entry looks like

a.

$$\nabla_{W_k} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial W_{k11}} & \frac{\partial L(\theta)}{\partial W_{k12}} & \frac{\partial L(\theta)}{\partial W_{k13}} \\ \frac{\partial L(\theta)}{\partial W_{k21}} & \frac{\partial L(\theta)}{\partial W_{k22}} & \frac{\partial L(\theta)}{\partial W_{k23}} \\ \frac{\partial L(\theta)}{\partial W_{k31}} & \frac{\partial L(\theta)}{\partial W_{k32}} & \frac{\partial L(\theta)}{\partial W_{k33}} \end{bmatrix} \quad \nabla_{W_{kj}} L(\theta) = \frac{\partial L(\theta)}{\partial a_{ki}}$$

$$\nabla_{W_k} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial a_{k1}} h_{k-1,1} & \frac{\partial L(\theta)}{\partial a_{k1}} h_{k-1,2} & \frac{\partial L(\theta)}{\partial a_{k1}} h_{k-1,3} \\ \frac{\partial L(\theta)}{\partial a_{k2}} h_{k-1,1} & \frac{\partial L(\theta)}{\partial a_{k2}} h_{k-1,2} & \frac{\partial L(\theta)}{\partial a_{k2}} h_{k-1,3} \\ \frac{\partial L(\theta)}{\partial a_{k3}} h_{k-1,1} & \frac{\partial L(\theta)}{\partial a_{k3}} h_{k-1,2} & \frac{\partial L(\theta)}{\partial a_{k3}} h_{k-1,3} \end{bmatrix} = \nabla_{a_k} L(\theta) \cdot h_{k-1}^T$$

- Thus we can update all the weights in one go using $W_k = W_k - \eta(\nabla_{a_k} L(\theta) \cdot h_{k-1}^T)$
- $\nabla_{W_k} L(\theta) = \nabla_{a_k} L(\theta) \cdot h_{k-1}^T$
- Finally coming to the biases, $a_{ki} = b_{ki} \sum_j W_{kij} h_{k-1,j}$
- $\frac{\partial L(\theta)}{\partial b_{ki}} = \frac{\partial L(\theta)}{\partial a_{ki}} \frac{\partial a_{ki}}{\partial b_{ki}} = \frac{\partial L(\theta)}{\partial a_{ki}}$
- We can now write the gradient w.r.t the vector b_k

a.

$$\nabla_{b_k} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial a_{k1}} \\ \vdots \\ \frac{\partial L(\theta)}{\partial a_{kn}} \end{bmatrix} = \nabla_{b_a} L(\theta)$$

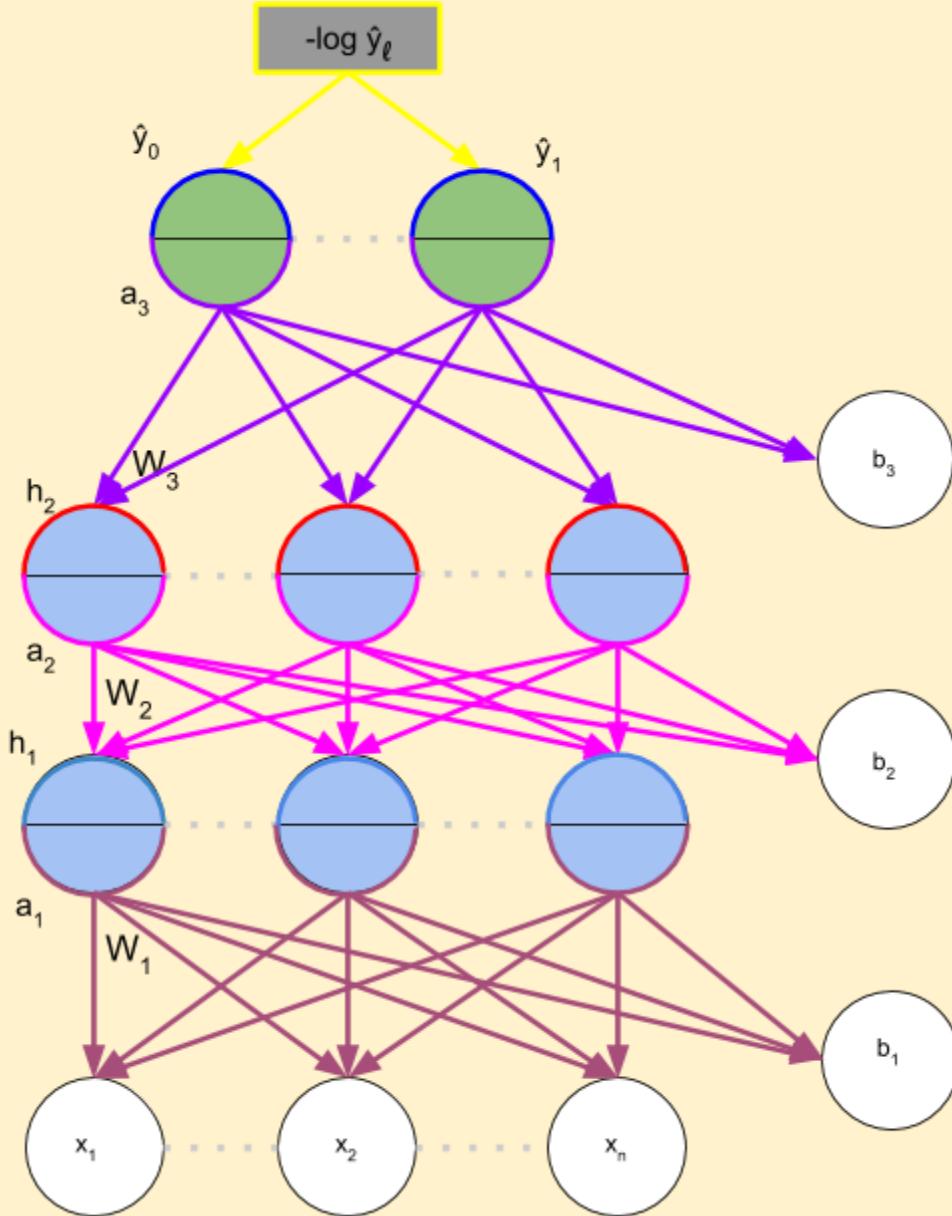
- Thus, we can update all biases using $b_k = b_k - \eta(\nabla_{b_k} L(\theta))$
- $\nabla_{b_k} L(\theta) = \nabla_{a_k} L(\theta)$

One Fourth Labs

A running example of backpropagation

The full story so far

1. Consider the following sample Neural Network



2. Steps to implement Backpropagation

- a. Randomly initialise W and b
- b. Forward propagation
 - i. For $k = 1$ to $L-1$ do

$$a_k = b_k + W_k h_{k-1}$$

$$h_k = g(a_k)$$
 - ii. end
 - iii. $a_L = b_L + W_L h_{L-1};$
 - iv. $\hat{y} = O(a_L)$

PadhAI: Backpropagation - the full version

One Fourth Labs

- c. Backpropagation
 - i. //Compute output gradient
 - ii. $\nabla_{a_L} L(\theta) = -(y - \hat{y}_i)$
 - iii. **For k = L to 1 do**
 //Compute gradients w.r.t parameters
$$\nabla_{W_k} L(\theta) = \nabla_{a_k} L(\theta) \cdot h_{k-1}^T$$
$$\nabla_{b_k} L(\theta) = \nabla_{a_k} L(\theta)$$

 //Compute gradients w.r.t layer below
$$\nabla_{h_{k-1}} L(\theta) = (W_k)^T (\nabla_{a_k} L(\theta))$$

 /// Compute gradients w.r.t layer below (pre-activation)
$$\nabla_{a_{k-1}} L(\theta) = \nabla_{h_{k-1}} L(\theta) \odot [..., g'(a_{k-1,j}), ...]$$
 - iv. **end**
3. Sample calculations to be added at a later date. $\Delta y_n / \Delta x_n$

PadhAI: Backpropagation - the full version

One Fourth Labs

Summary

What are the new things we've learned in this module

1. Here are some of the takeaways from this chapter
 - a. Data: Real inputs $x_i \in \mathbb{R}$
 - b. Task:
 - i. Binary classification
 - ii. Multi-class classification
 - iii. **Regression**
 - c. Model: Deep Neural Network to deal with complex decision boundaries
 - d. Loss:
 - i. Cross entropy loss: $L(\Theta) = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d y_{ij} \log(\hat{y}_{ij})$
 - ii. **Square Error Loss:** $L(\Theta) = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d (y_{ij} - \hat{y}_{ij})^2$
 - e. Learning: Gradient Descent with backpropagation
 - f. Evaluation:
 - i. Accuracy = $\frac{\text{Number of correct predictions}}{\text{Total Number of predictions}}$
 - ii. Per-class Accuracy = $\frac{\text{Number of correct predictions of a class}}{\text{Total Number of true values of that class}}$
2. Topics highlighted in red are to be covered in future segments