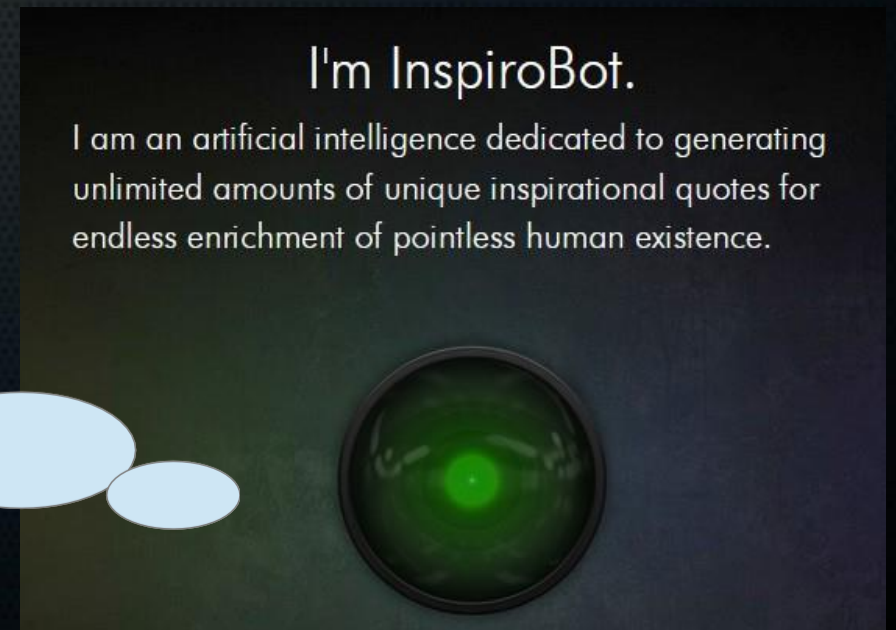
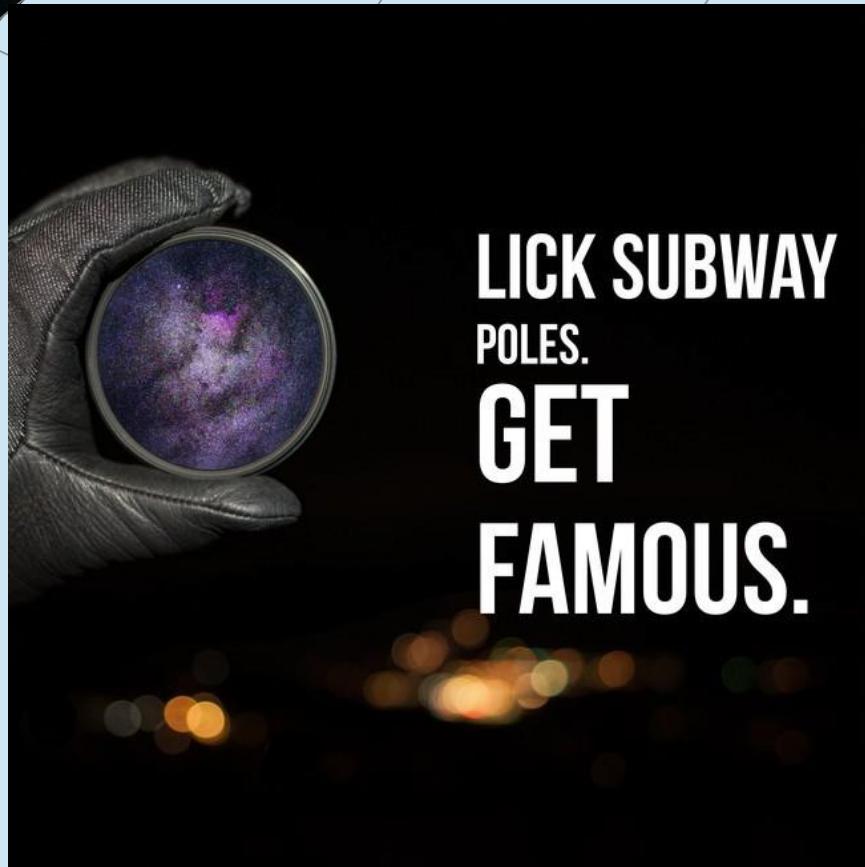


Daily Inspiration



Today

- Recap yesterday
- Neural network backpropagation
- Convolutional neural networks

Yesterday

- Logistic regression:

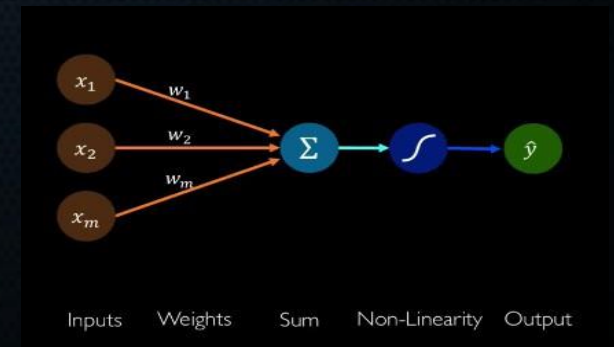
- use the sigmoid function to turn the tools of regression into classification
- Logistic regression cost function:

$$\text{Cost}(x) = -y \cdot \log(h_{\theta}(x)) - (1 - y) \cdot \log(1 - h_{\theta}(x))$$

- Multiclass (i classes): train i binary classifiers

- Neural networks:

- Hierarchically ordered units that calculate ever more complex functions using simpler calculated features
- Universal approximation
- Forward propagation
- Multiclass: turn training samples into i -dimensional vectors



Exam

- Value of this course is in the practicals, in (trying to) do it yourself.
- Exam will ask mostly basic understanding:
 - I have these predictions and these Y , what is the MSE?
 - Showing that you understood how to get linear regression predictions with linear algebra
 - Calculating a few Euclidean distances
 - Showing that you know how nested cross-validation works and why we use it

Exam

- Since we worked with them extensively, you should be able to write down the partial derivatives of a simple linear regression
- For a simple neural network image: should be able to write down in linear algebra how a certain layer calculates its activations.
- Should be able to write down how backprop works for a NN
- Bit of *pseudocode*: I will probably give you part of the steps for a certain function you implemented. You complete it by writing down the logical steps you would need to do in words (doesn't need to be working code, but the right steps).

Before we start

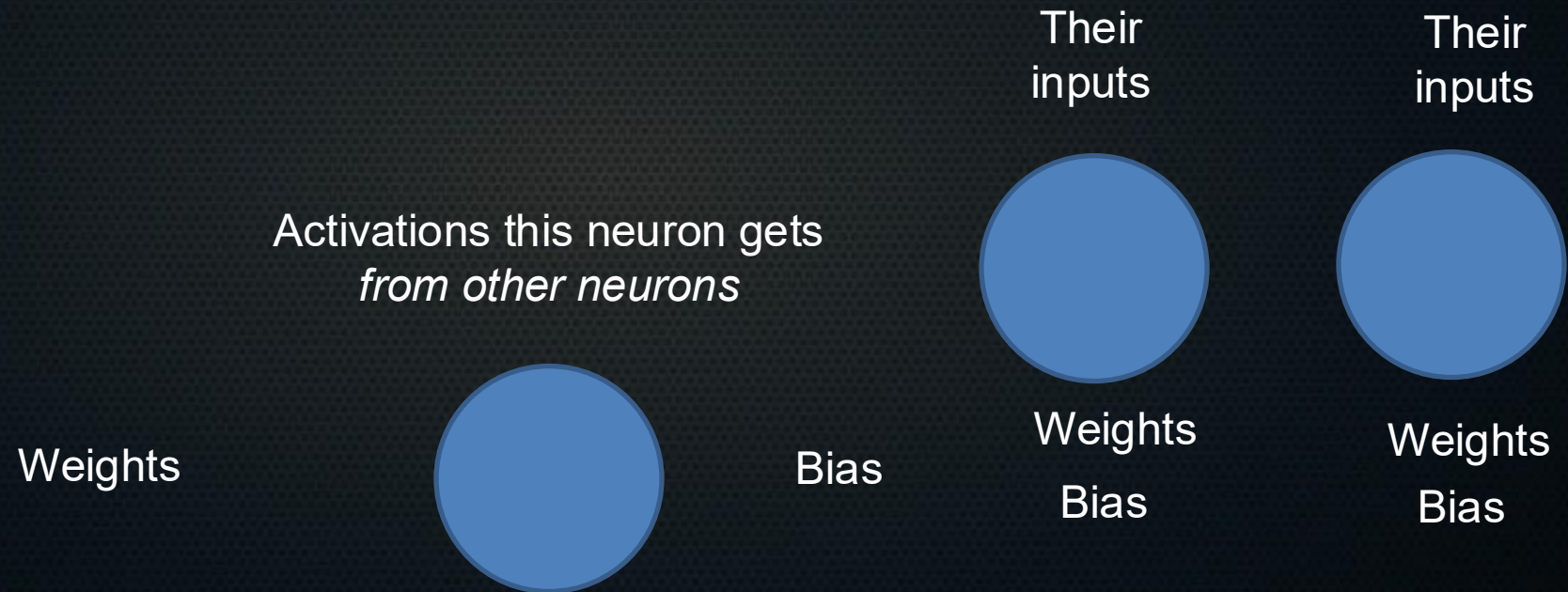
$J(\text{neuron}_1, \text{neuron}_2, \text{neuron}_3 \dots \text{neuron}_n) = \text{network predictions} - \text{real values}$ (actually use binary cross-entropy)



What knobs can we
turn in our function?

Before we start

$J(\text{neuron}_1, \text{neuron}_2, \text{neuron}_3 \dots \text{neuron}_n) = \text{network predictions} - \text{real values}$ (actually use binary cross-entropy)



What knobs can we turn in our function?

Before we start

$J(\text{neuron}_1, \text{neuron}_2, \text{neuron}_3 \dots \text{neuron}_n) = \text{network predictions} - \text{real values}$ (actually use binary cross-entropy)

$a(b(c(d(e(f(g(h(i(j(k(x))))))))))$

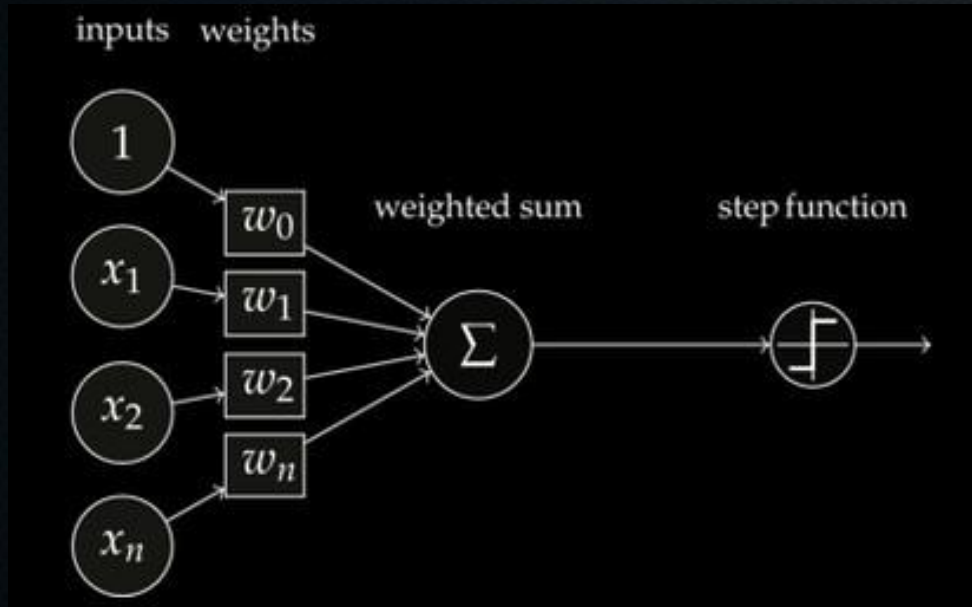
Lots of chain rules for the partial
derivatives!

Before we start

- This is going to get quite mathy.
- Implementing this yourself is a clear *stretch goal* . I want you to understand how backpropagation works, and the basic partial derivatives you use and chain for it. Actually making it yourself with linear algebra is the dream, but probably not the reality.

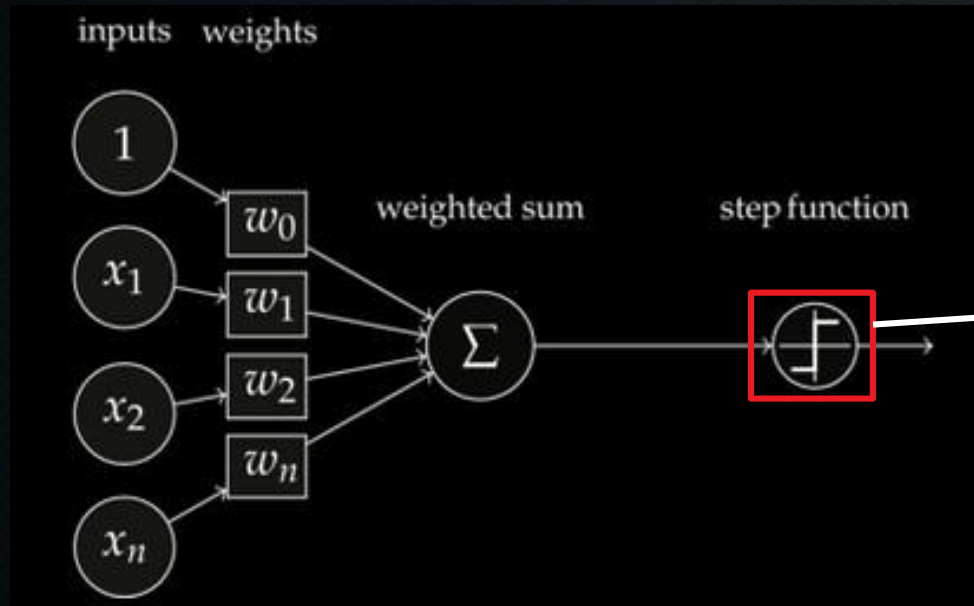
Backpropagation

- Some history: neural networks started out as perceptrons.



Backpropagation

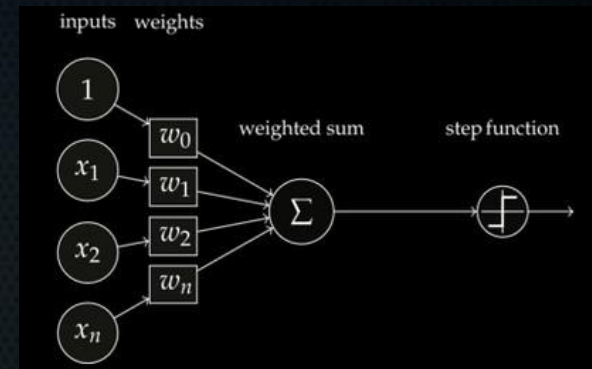
- Some history: neural networks started out as perceptrons.



Output is either 1 or 0. No in-between like the sigmoid!

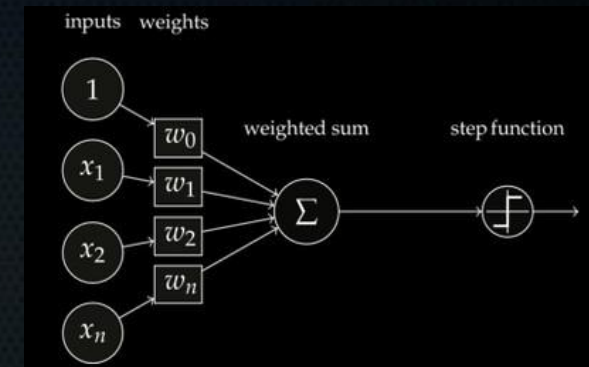
Backpropagation

- Some history: neural networks started out as perceptrons.
- In 1969, a paper was published by Marvin Minsky and Seymour Papert that showed that a single perceptron could not learn XOR



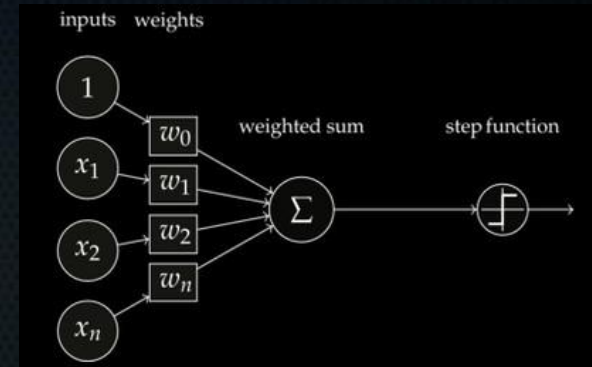
Backpropagation

- Some history: neural networks started out as perceptrons.
- In 1969, a paper was published by Marvin Minsky and Seymour Papert that showed that a single perceptron could not learn XOR
- While it was known that multi-layer perceptrons could, the perceptron learning rule was not good at learning multi-layer networks



Backpropagation

- Some history: neural networks started out as perceptrons.
- In 1969, a paper was published by Marvin Minsky and Seymour Papert that showed that a single perceptron could not learn XOR
- While it was known that multi-layer perceptrons could, the perceptron learning rule was not good at learning multi-layer networks
- Unfortunately, this mostly killed neural network research for a decade (!)



Backpropagation

- 1986 to the rescue!*

Learning representations by back-propagating errors

David E. Rumelhart*, Geoffrey E. Hinton†
& Ronald J. Williams*

* Institute for Cognitive Science, C-015, University of California,
San Diego, La Jolla, California 92093, USA

† Department of Computer Science, Carnegie-Mellon University,
Pittsburgh, Philadelphia 15213, USA

We describe a new learning procedure, back-propagation, for networks of neurone-like units. The procedure repeatedly adjusts the weights of the connections in the network so as to minimize a measure of the difference between the actual output vector of the net and the desired output vector. As a result of the weight adjustments, internal 'hidden' units which are not part of the input or output come to represent important features of the task domain, and the regularities in the task are captured by the interactions of these units. The ability to create useful new features distinguishes back-propagation from earlier, simpler methods such as the perceptron-convergence procedure¹.

*Not completely true, approach stemmed from the 70s, and someone else even published on it in control systems in 1960, but this paper really showed the *power* of the approach and renewed interest!

Backpropagation

- 1986 to the rescue!
- Idea: rather than the step-function we use a smooth, *differentiable* function (sigmoid or other).
- We know the error in the last layer (we know true classes and we know the vector that our NN outputs)
- Due to this, we can take the error, change the parameters of a layer, then go back a layer and change the parameters there, etc.

Learning representations by back-propagating errors

David E. Rumelhart*, Geoffrey E. Hinton†
& Ronald J. Williams*

* Institute for Cognitive Science, C-015, University of California,
San Diego, La Jolla, California 92093, USA

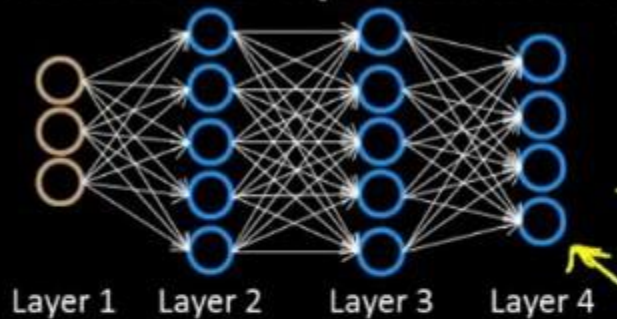
† Department of Computer Science, Carnegie-Mellon University,
Pittsburgh, Philadelphia 15213, USA

We describe a new learning procedure, back-propagation, for networks of neurone-like units. The procedure repeatedly adjusts the weights of the connections in the network so as to minimize a measure of the difference between the actual output vector of the net and the desired output vector. As a result of the weight adjustments, internal 'hidden' units which are not part of the input or output come to represent important features of the task domain, and the regularities in the task are captured by the interactions of these units. The ability to create useful new features distinguishes back-propagation from earlier, simpler methods such as the perceptron-convergence procedure¹.

Backpropagation - terminology

- How does that work?

Neural Network (Classification)



→ $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

→ $L =$ total no. of layers in network $\underline{L = 4}$

→ $s_l =$ no. of units (not counting bias unit) in layer l $s_1 = 3, s_2 = 4, s_3 = 4, s_4 = 4$

Binary classification

$y = 0$ or 1 ←

1 output unit ←

$$h_{\Theta}(x) \in \mathbb{R}$$

$$s_L = 1, \quad \underline{K = 1}$$



Multi-class classification (K classes)

$y \in \mathbb{R}^K$ E.g. $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$ ←
pedestrian car motorcycle truck

K output units

$$h_{\Theta}(x) \in \mathbb{R}^K$$

$$s_L = K \quad (K \geq 3)$$

Backpropagation

- First, we need a cost function. For logistic regression we used:
$$\text{Cost}(x) = -y \cdot \log(h_{\theta}(x)) - (1 - y) \cdot \log(1 - h_{\theta}(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)

Backpropagation

- First, we need a cost function. For logistic regression we used:
$$\text{Cost}(x) = -y \cdot \log(h_\theta(x)) - (1 - y) \cdot \log(1 - h_\theta(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)

$$\text{Cost}(x) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_\theta(x^{(i)}))_k) \right]$$

Backpropagation

- First, we need a cost function. For logistic regression we used:
$$\text{Cost}(x) = -y \cdot \log(h_\theta(x)) - (1-y) \cdot \log(1-h_\theta(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)

$$\text{Cost}(x) = -\frac{1}{m} \left[\underbrace{\sum_{i=1}^m}_{\text{over samples}} \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_\theta(x^{(i)}))_k) \right]$$

We have m training samples and calculate the cost over each sample i

Backpropagation

- First, we need a cost function. For logistic regression we used:
$$\text{Cost}(x) = -y \cdot \log(h_{\theta}(x)) - (1 - y) \cdot \log(1 - h_{\theta}(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)

$$\text{Cost}(x) = -\frac{1}{m} \left[\underbrace{\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\theta}(x^{(i)}))_k)}_{\text{We have m training samples and calculate the cost over each sample } i} \right]$$

We have m training samples and calculate the cost over each sample i

We have K classes and sum the cost of each entry k in the class vector

Backpropagation

- First, we need a cost function. For logistic regression we used:

$$\text{Cost}(x) = -y \cdot \log(h_{\theta}(x)) - (1 - y) \cdot \log(1 - h_{\theta}(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)

$$\text{Cost}(x) = -\frac{1}{m} \left[\underbrace{\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\theta}(x^{(i)}))_k)}_i \right]$$

We have m training samples and calculate the cost over each sample i

We have K classes and sum the cost of each entry k in the class vector



$$y^{(3)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$h_{\theta}(x^{(3)}) = \begin{bmatrix} 0.234 \\ 0.678 \\ 0.102 \\ 0.020 \\ 0.800 \end{bmatrix}$$



Backpropagation

- First, we need a cost function. For logistic regression we used:

$$\text{Cost}(x) = -y \cdot \log(h_{\theta}(x)) - (1 - y) \cdot \log(1 - h_{\theta}(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)

$$\text{Cost}(x) = -\frac{1}{m} \left[\underbrace{\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\theta}(x^{(i)}))_k)}_{\text{We have m training samples and calculate the cost over each sample } i} \right]$$

We have m training samples and calculate the cost over each sample i

We have K classes and sum the cost of each entry k in the class vector



$k = 1$
 $y^{(3)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
 $h_{\theta}(x^{(3)}) = \begin{bmatrix} 0.234 \\ 0.678 \\ 0.102 \\ 0.020 \\ 0.800 \end{bmatrix}$

The first row of the matrix $h_{\theta}(x^{(3)})$ is highlighted with a red box. To the right of the matrix are five empty circles.

Backpropagation

- First, we need a cost function. For logistic regression we used:

$$\text{Cost}(x) = -y \cdot \log(h_{\theta}(x)) - (1 - y) \cdot \log(1 - h_{\theta}(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)

$$\text{Cost}(x) = -\frac{1}{m} \left[\underbrace{\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\theta}(x^{(i)}))_k)}_{\text{We have m training samples and calculate the cost over each sample } i} \right]$$

We have m training samples and calculate the cost over each sample i

We have K classes and sum the cost of each entry k in the class vector



$$y^{(3)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad h_{\theta}(x^{(3)}) = \begin{bmatrix} 0.234 \\ 0.678 \\ 0.102 \\ 0.020 \\ 0.800 \end{bmatrix}$$

○
○
○
○
○

Backpropagation

- First, we need a cost function. For logistic regression we used:

$$\text{Cost}(x) = -y \cdot \log(h_{\theta}(x)) - (1 - y) \cdot \log(1 - h_{\theta}(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)

$$\text{Cost}(x) = -\frac{1}{m} \left[\underbrace{\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\theta}(x^{(i)}))_k)}_{\text{We have m training samples and calculate the cost over each sample } i} \right]$$

We have m training samples and calculate the cost over each sample i

We have K classes and sum the cost of each entry k in the class vector

$y^{(3)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
 $h_{\theta}(x^{(3)}) = \begin{bmatrix} 0.234 \\ 0.678 \\ 0.102 \\ 0.020 \\ 0.800 \end{bmatrix}$

$k = K = 5$

Backpropagation

- First, we need a cost function. For logistic regression we used:

$$\text{Cost}(x) = -y \cdot \log(h_\theta(x)) - (1-y) \cdot \log(1-h_\theta(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)

$$\text{Cost}(x) = -\frac{1}{m} \left[\underbrace{\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_\theta(x^{(i)}))_k)}_{\text{We have m training samples and calculate the cost over each sample } i} \right]$$

We have m training samples and calculate the cost over each sample i

We have K classes and sum the cost of each entry k in the class vector

$$y^{(3)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad h_\theta(x^{(3)}) = \begin{bmatrix} 0.234 \\ 0.678 \\ 0.102 \\ 0.020 \\ 0.800 \end{bmatrix}$$

○
○
○
○
○

$k = K = 5$

$$\text{Cost}(x^{(3)}; \text{5th neuron}) = -[\underbrace{y_5^{(3)}}_{\text{green box}} \log(h_\theta(x^{(3)}))_5 + (1 - y_5^{(3)}) \log(1 - (h_\theta(x^{(3)}))_5)]$$

Backpropagation

- First, we need a cost function. For logistic regression we used:

$$\text{Cost}(x) = -y \cdot \log(h_\theta(x)) - (1 - y) \cdot \log(1 - h_\theta(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)

$$\text{Cost}(x) = -\frac{1}{m} \left[\underbrace{\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_\theta(x^{(i)}))_k)}_{\text{We have m training samples and calculate the cost over each sample } i}$$

We have m training samples and calculate the cost over each sample i

We have K classes and sum the cost of each entry k in the class vector

$$y^{(3)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad h_\theta(x^{(3)}) = \begin{bmatrix} 0.234 \\ 0.678 \\ 0.102 \\ 0.020 \\ 0.800 \end{bmatrix}$$

$k = K = 5$

○
○
○
○
○

$$\text{Cost}(x^{(3)}; 5\text{th neuron}) = - \left[\boxed{0} + (1 - y_5^{(3)}) \log(1 - (h_\theta(x^{(3)}))_5) \right]$$

Backpropagation

- First, we need a cost function. For logistic regression we used:

$$\text{Cost}(x) = -y \cdot \log(h_\theta(x)) - (1-y) \cdot \log(1-h_\theta(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)

$$\text{Cost}(x) = -\frac{1}{m} \left[\underbrace{\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_\theta(x^{(i)}))_k)} \right]$$

We have m training samples and calculate the cost over each sample i

We have K classes and sum the cost of each entry k in the class vector

$$y^{(3)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad h_\theta(x^{(3)}) = \begin{bmatrix} 0.234 \\ 0.678 \\ 0.102 \\ 0.020 \\ 0.800 \end{bmatrix}$$

$k = K = 5$

$$\text{Cost}(x^{(3)}; 5\text{th neuron}) = - \left[\boxed{0} + (1 - \boxed{y_5^{(3)}}) \log(1 - \boxed{(h_\theta(x^{(3)}))_5}) \right]$$

Backpropagation

- First, we need a cost function. For logistic regression we used:

$$\text{Cost}(x) = -y \cdot \log(h_\theta(x)) - (1-y) \cdot \log(1-h_\theta(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)

$$\text{Cost}(x) = -\frac{1}{m} \left[\underbrace{\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_\theta(x^{(i)}))_k)} \right]$$

We have m training samples and calculate the cost over each sample i

We have K classes and sum the cost of each entry k in the class vector

$$y^{(3)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad h_\theta(x^{(3)}) = \begin{bmatrix} 0.234 \\ 0.678 \\ 0.102 \\ 0.020 \\ 0.800 \end{bmatrix}$$

○
○
○
○
○

$$\text{Cost}(x^{(3)}; 5\text{th neuron}) = -[\boxed{0} + (1 - \boxed{0}) \log(1 - \boxed{0.8})]$$

Backpropagation

- First, we need a cost function. For logistic regression we used:

$$\text{Cost}(x) = -y \cdot \log(h_{\theta}(x)) - (1 - y) \cdot \log(1 - h_{\theta}(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)

$$\text{Cost}(x) = -\frac{1}{m} \left[\underbrace{\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\theta}(x^{(i)}))_k)}_{\text{We have m training samples and calculate the cost over each sample } i} \right]$$

We have m training samples and calculate the cost over each sample i

We have K classes and sum the cost of each entry k in the class vector

$$y^{(3)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad k = K = 5$$

$$h_{\theta}(x^{(3)}) = \begin{bmatrix} 0.234 \\ 0.678 \\ 0.102 \\ 0.020 \\ 0.800 \end{bmatrix} \quad \begin{matrix} \bigcirc \\ \bigcirc \\ \bigcirc \\ \bigcirc \\ \bigcirc \end{matrix}$$

$$\text{Cost}(x^{(3)}; 5\text{th neuron}) = -[\log(0.2)] \approx 0.69897$$

Backpropagation

- First, we need a cost function. For logistic regression we used:

$$\text{Cost}(x) = -y \cdot \log(h_\theta(x)) - (1-y) \cdot \log(1-h_\theta(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)

$$\text{Cost}(x) = -\frac{1}{m} \left[\underbrace{\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)}))_k}_{\text{Error if } y = 1 \text{ for this position in the class vector}} + \underbrace{\sum_{i=1}^m \sum_{k=1}^K (1 - y_k^{(i)}) \log(1 - (h_\theta(x^{(i)}))_k)}_{\text{Error if } y = 0 \text{ for this position in the class vector}} \right]$$

We have m training samples and calculate the cost over each sample i

We have K classes and sum the cost of each entry k in the class vector

Error if $y = 0$ for this position in the class vector

Backpropagation

- First, we need a cost function. For logistic regression we used:

$$\text{Cost}(x) = -y \cdot \log(h_\theta(x)) - (1 - y) \cdot \log(1 - h_\theta(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)

$$\text{Cost}(x) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K \underbrace{y_k^{(i)} \log(h_\theta(x^{(i)}))_k}_{\text{Error if } y = 1 \text{ for this position in the class vector}} + \underbrace{(1 - y_k^{(i)}) \log(1 - (h_\theta(x^{(i)}))_k)}_{\text{Error if } y = 0 \text{ for this position in the class vector}} \right]$$

We have m training samples and calculate the cost over each sample i

We have K classes and sum the cost of each entry k in the class vector

Average over all training samples m

Error if $y = 0$ for this position in the class vector

Backpropagation

- **We (may) want to regularise all weights and biases for each layer in the network**

$$\text{Cost}(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\theta}(x^{(i)}))_k) \right]$$

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\theta_{ij}^{(l)})^2$$

l = which layer's parameters are we looking at?

i = which unit in that layer (row in theta matrix) are we looking at?

j = what parameter of that unit (column) are we looking at?

Backpropagation

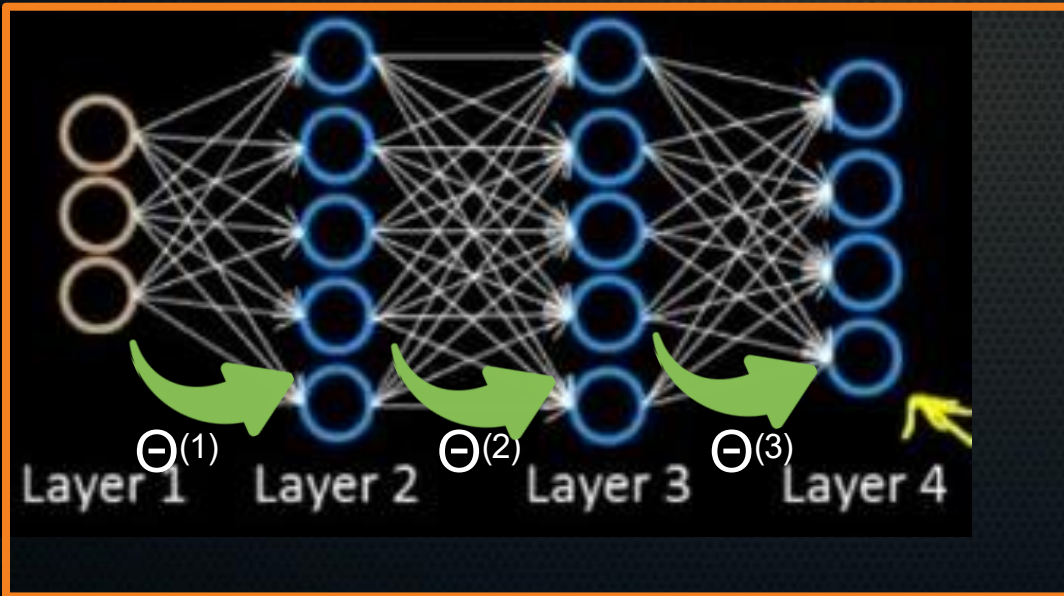
$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\theta_{ij}^{(l)})^2$$

l = which layer's parameters are we looking at?

i = which unit in that layer (row in theta matrix) are we looking at?

j = what parameter of that unit (column) are we looking at?

$\Theta^{(j)}$ = matrix of weights controlling function mapping from layer j to layer $j + 1$



Have 3 matrices of parameters.

$\Theta^{(1)}$ = parameters layer 2 uses to calculate values with layer 1 as input.

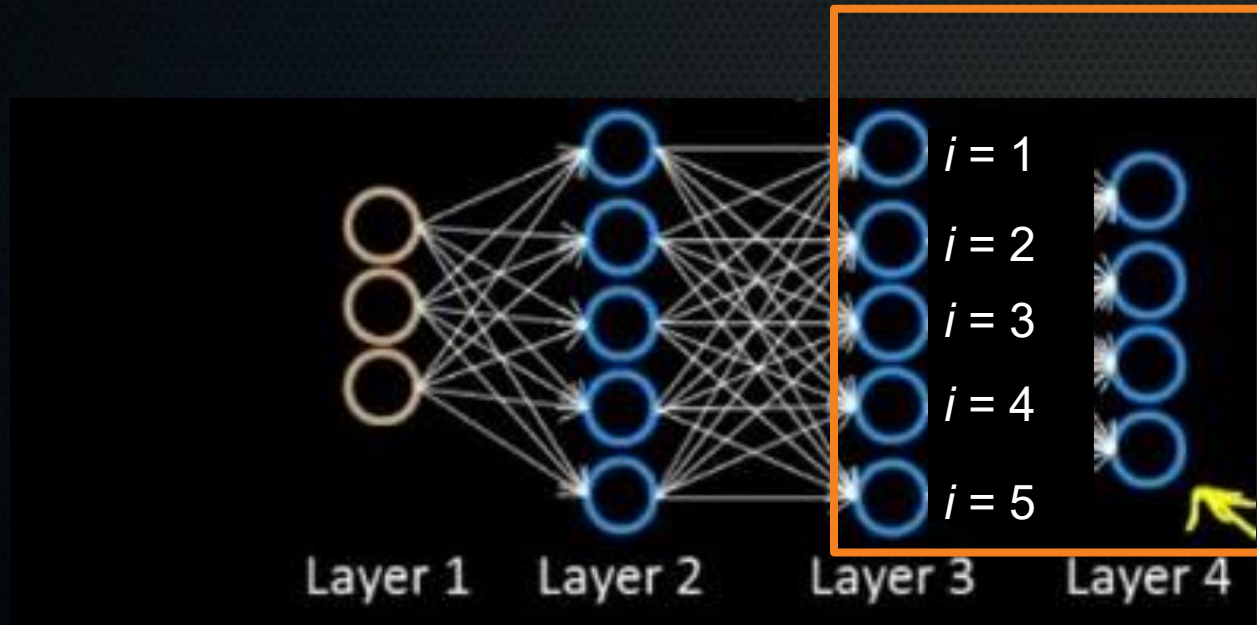
Backpropagation

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\theta_{ij}^{(l)})^2$$

l = which layer's parameters are we looking at?

i = which unit in that layer (row in theta matrix) are we looking at?

j = what parameter of that unit (column) are we looking at?



Each unit's parameters are in a certain row of that layer's Θ matrix.

Backpropagation

- First, we need a cost function. For logistic regression we used:
$$\text{Cost}(x) = -y \cdot \log(h_{\theta}(x)) - (1 - y) \cdot \log(1 - h_{\theta}(x))$$
- We use the same thing, only generalised for the fact that:
 - Our output is a vector (not a 1 or 0 as for logistic regression)
 - We (may) want to regularise all weights and biases for each layer in the network

$$\text{Cost}(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\theta}(x^{(i)}))_k) \right]$$

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\theta_{ij}^{(l)})^2$$

l = which layer's parameters are we looking at?

i = which unit in that layer (row in theta matrix) are we looking at?

j = what parameter of that unit (column) are we looking at?

How do we use this cost to update parameters?

Given one training example (x, y) :

Forward propagation:

$$a^{(1)} = x$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

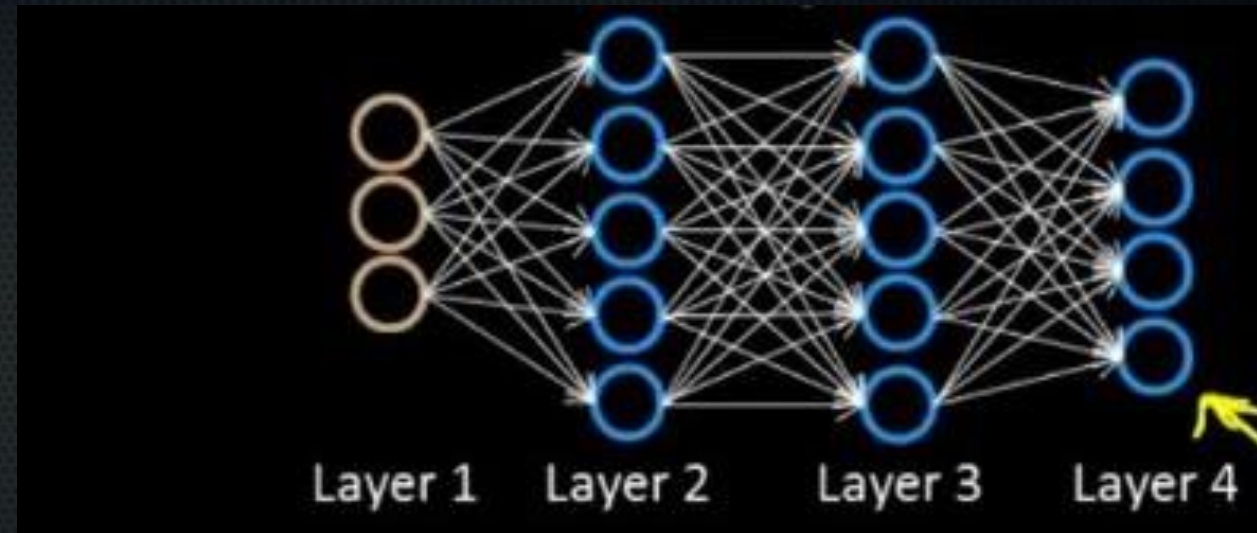
$$a^{(2)} = g(z^{(2)}) \quad (\text{add } a_0^{(2)})$$

$$z^{(3)} = \Theta^{(2)} a^{(2)}$$

$$a^{(3)} = g(z^{(3)}) \quad (\text{add } a_0^{(3)})$$

$$z^{(4)} = \Theta^{(3)} a^{(3)}$$

$$a^{(4)} = h_{\Theta}(x) = g(z^{(4)})$$

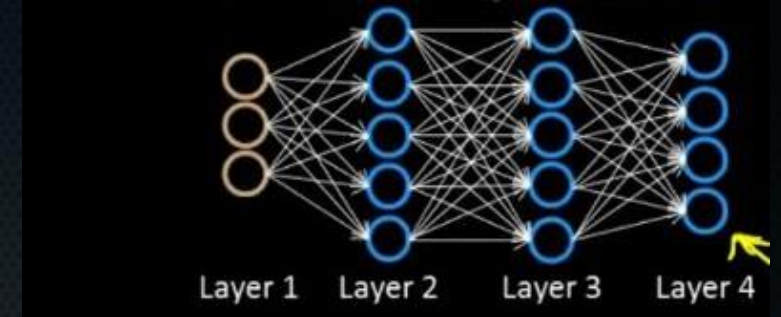
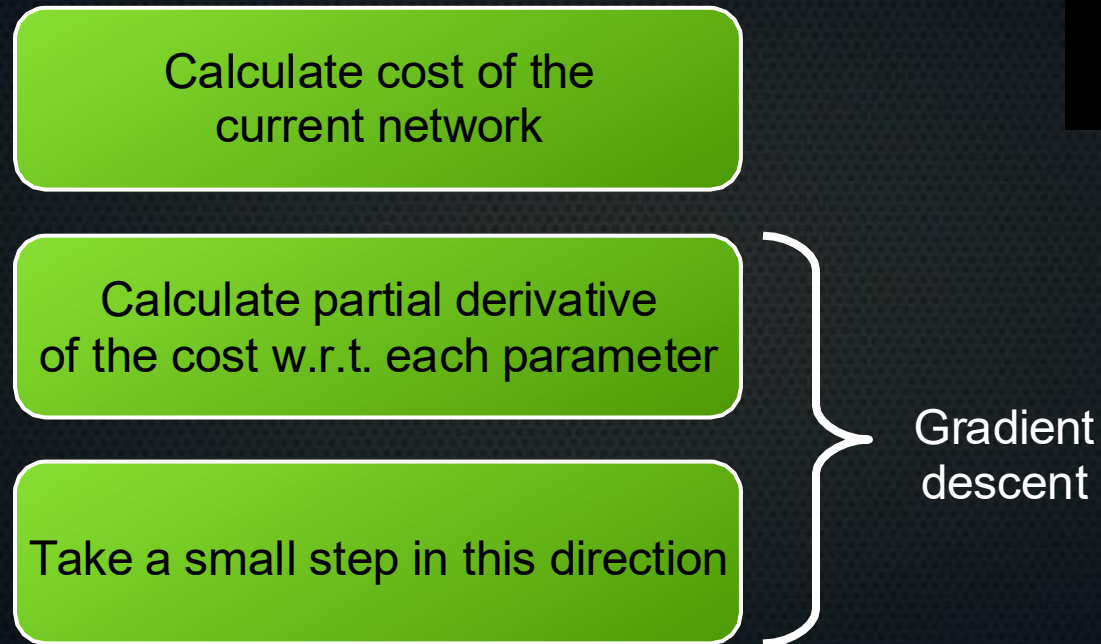


$$z = \left(\sum_{i=1}^{n_{\text{input neurons}}} w_i \cdot x_i \right) + b$$

$$g = \text{sigmoid} = \frac{1}{1 + e^{-z}}$$

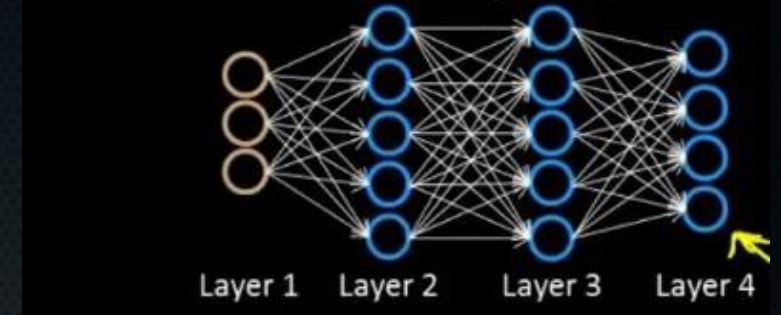
How do we use this cost to update parameters?

- Scheme as follows:



How do we use this cost to update parameters?

- With all these neurons in all these layers, we have a cost function that is dependent on *many* parameters.
- To tackle this complexity: let's first see how one *neuron* influences the cost function.

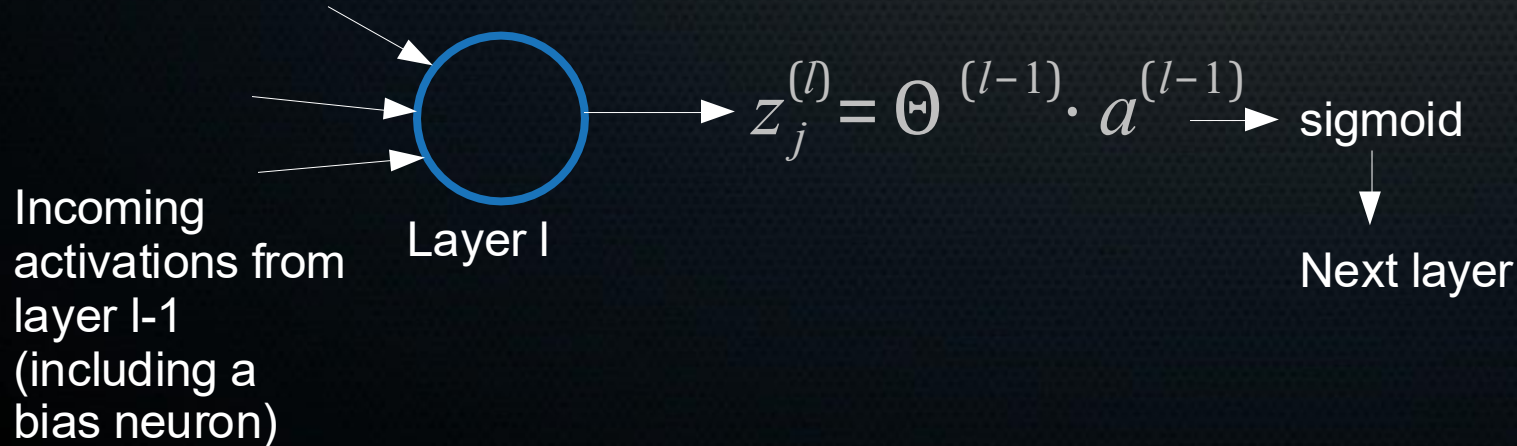


Gradient descent

Calculate cost of the current network

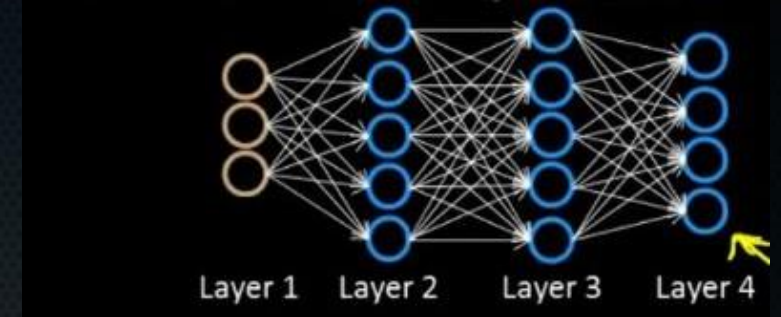
Calculate partial derivative of the cost w.r.t. each parameter

Take a small step in this direction



How do we use this cost to update parameters?

- Partial derivative w.r.t. this neuron's $z = \frac{\partial C}{\partial z_j^{(l)}}$
- If we somehow change this neuron's z with small $\Delta z_j^{(l)}$, then total cost will change with $\frac{\partial C}{\partial z_j^{(l)}} \cdot \Delta z_j^{(l)}$

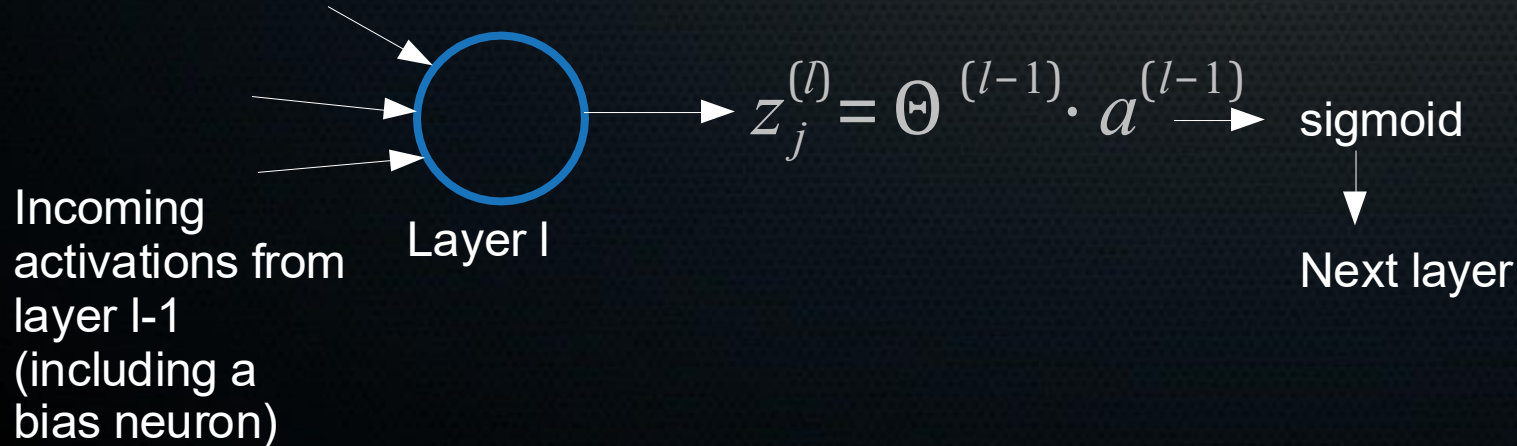


Gradient descent

Calculate cost of the current network

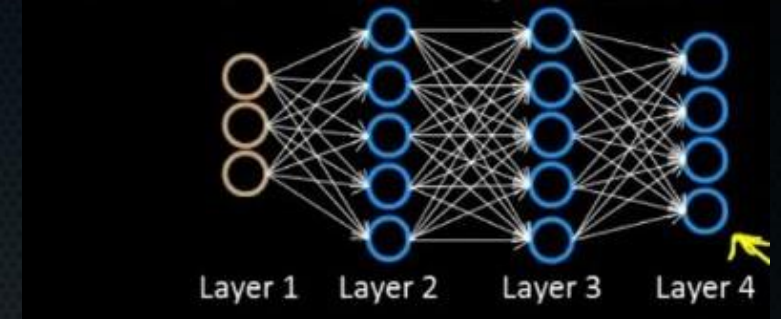
Calculate partial derivative of the cost w.r.t. each parameter

Take a small step in this direction



How do we use this cost to update parameters?

- Partial derivative w.r.t. this neuron's $z = \frac{\partial C}{\partial z_j^{(l)}}$
- If we somehow change this neuron's z with small $\Delta z_j^{(l)}$, then total cost will change with $\frac{\partial C}{\partial z_j^{(l)}} \cdot \Delta z_j^{(l)}$
- If $\frac{\partial C}{\partial z_j^{(l)}}$ is 20, then a small nudge of -0.1 gives a change in total cost of -2.

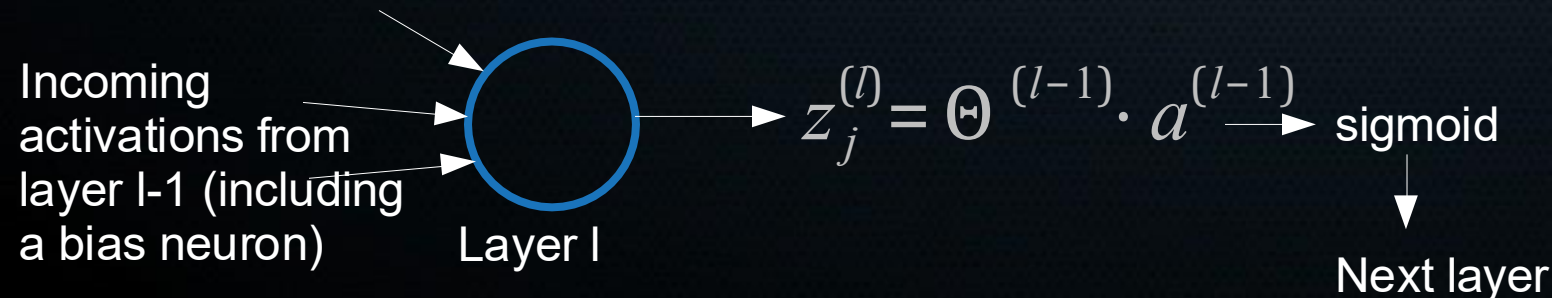


Gradient descent

Calculate cost of the current network

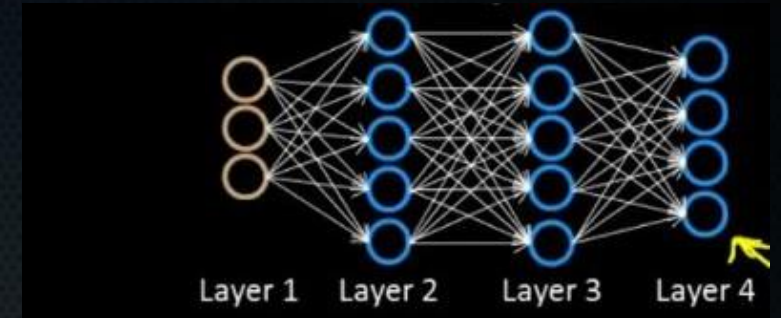
Calculate partial derivative of the cost w.r.t. each parameter

Take a small step in this direction



How do we use this cost to update parameters?

- Partial derivative w.r.t. this neuron's $z = \frac{\partial C}{\partial z_j^{(l)}}$
- If we somehow change this neuron's z with small $\Delta z_j^{(l)}$, then total cost will change with $\frac{\partial C}{\partial z_j^{(l)}} \cdot \Delta z_j^{(l)}$
- If $\frac{\partial C}{\partial z_j^{(l)}}$ is 0.01, then a small nudge of -0.1 gives a change in cost of -0.001.

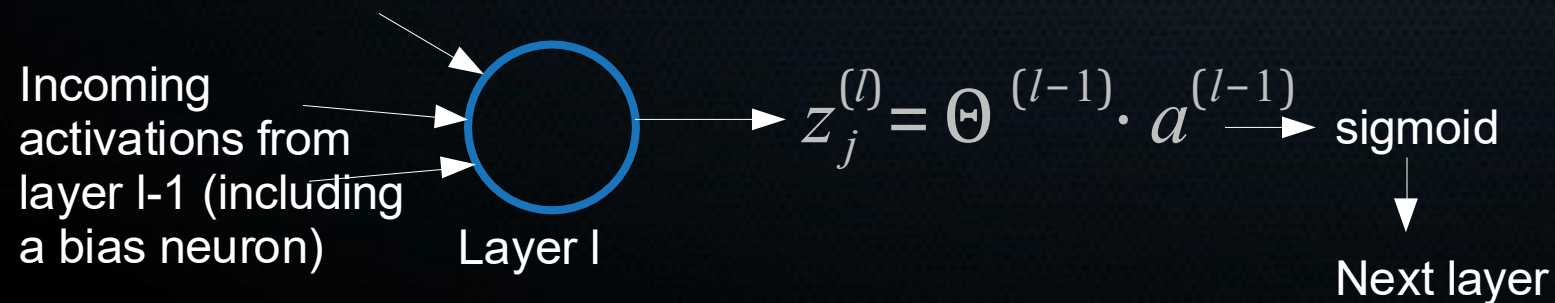


Gradient descent

Calculate cost of the current network

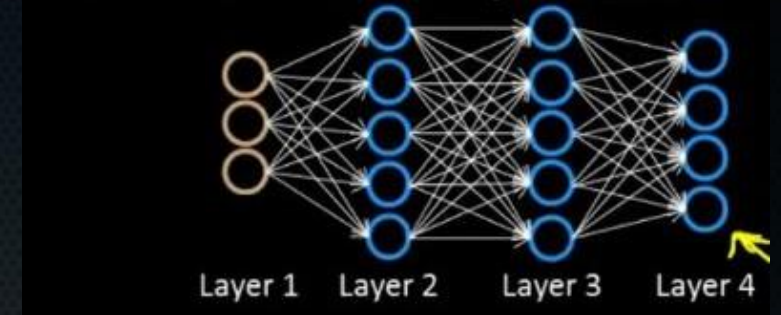
Calculate partial derivative of the cost w.r.t. each parameter

Take a small step in this direction



How do we use this cost to update parameters?

- Partial derivative w.r.t. this neuron's $z = \frac{\partial C}{\partial z_j^{(l)}}$
- If we somehow change this neuron's z with small $\Delta z_j^{(l)}$, then total cost will change with $\frac{\partial C}{\partial z_j^{(l)}} \cdot \Delta z_j^{(l)}$
- Hence, we see that **this quantity** shows heuristically how wrong a neuron is.

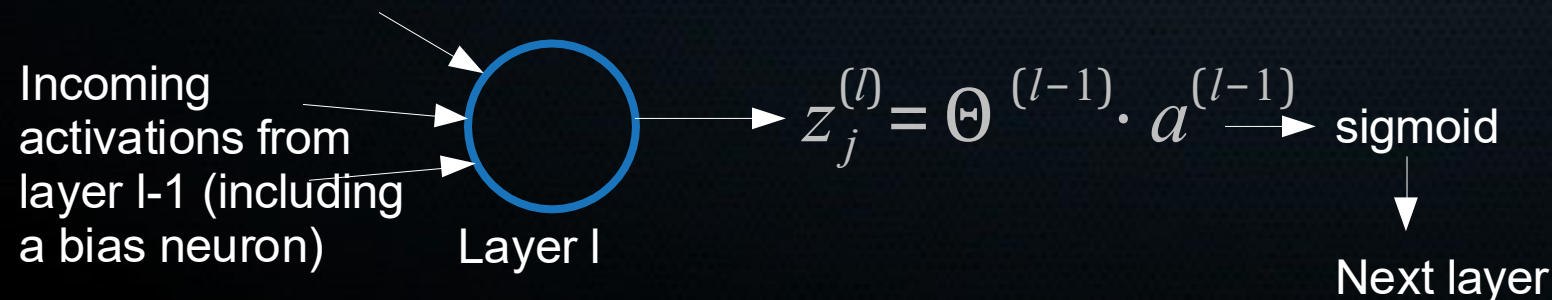


Gradient descent

Calculate cost of the current network

Calculate partial derivative of the cost w.r.t. each parameter

Take a small step in this direction

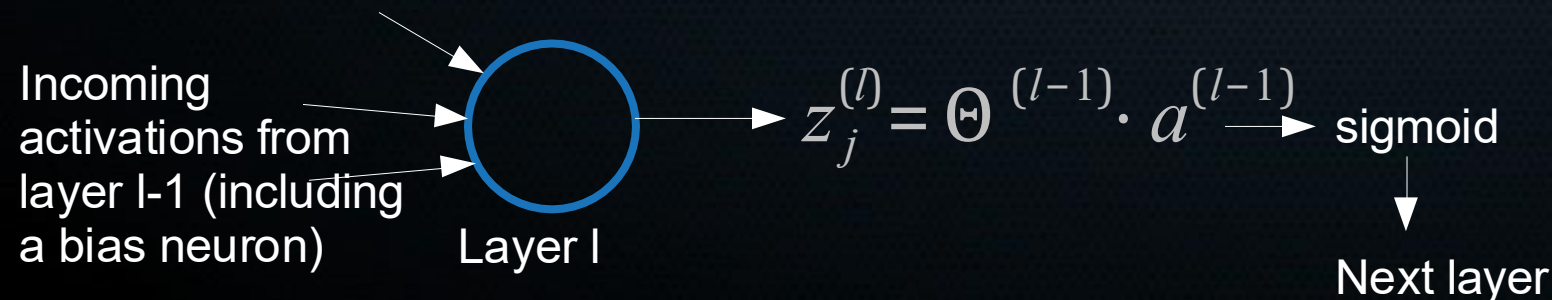


How do we use this cost to update parameters?

- So, the error for a neuron:

$$\delta_j^{(l)} = \text{error of node } j \text{ in layer } l = \frac{\partial C}{\partial z_j^{(l)}}$$

- Remember, we want to change the weights and biases, that is, calculate all $\frac{\partial C}{\partial w_{jk}^{(l)}}$ $\frac{\partial C}{\partial b_j^{(l)}}$
- Then we can take a small step and update.



How do we use this cost to update parameters?

- So, the error for a neuron:
 $\delta_j^{(l)} = \text{error of node } j \text{ in layer } l = \frac{\partial C}{\partial z_j^{(l)}}$
- If we can :
 - Find this error for every neuron
 - Relate this quantity to how to change the weights and biases

We are set.

How do we use this cost to update parameters?

- So, the error for a neuron:

$$\delta_j^{(l)} = \text{error of node } j \text{ in layer } l = \frac{\partial C}{\partial z_j^{(l)}}$$

Abusing the chain rule:

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \frac{\delta_j^{(l)} \cdot \sum_{k=1}^{n_{\text{weights}}} w_{jk}^{(l)} \cdot a^{(l-1)} + b_j^{(l)}}{\partial b_j^{(l)}}$$

How do we use this cost to update parameters?

- So, the error for a neuron:

$$\delta_j^{(l)} = \text{error of node } j \text{ in layer } l = \frac{\partial C}{\partial z_j^{(l)}}$$

Using the chain rule:

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \frac{\delta_j^{(l)} \cdot \sum_{k=1}^{n_{\text{weights}}} w_{jk}^{(l)} \cdot a^{(l-1)} + b_j^{(l)}}{\partial b_j^{(l)}}$$

These cancel out, giving what we want: how to change the bias of a neuron

How do we use this cost to update parameters?

- So, the error for a neuron:

$$\delta_j^{(l)} = \text{error of node } j \text{ in layer } l = \frac{\partial C}{\partial z_j^{(l)}}$$

Using the chain rule:

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \cdot \frac{\sum_{k=1}^{n_{\text{weights}}} w_{jk}^{(l)} \cdot a^{(l-1)} + b_j^{(l)}}{\partial b_j^{(l)}}$$

We just defined this delta term as shorthand for the error of each neuron's product of the weights with its inputs plus its bias

How do we use this cost to update parameters?

- So, the error for a neuron:

$$\delta_j^{(l)} = \text{error of node } j \text{ in layer } l = \frac{\partial C}{\partial z_j^{(l)}}$$

Using the chain rule:

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \cdot \frac{\partial \left(\sum_{k=1}^{n_{\text{weights}}} w_{jk}^{(l)} \cdot a^{(l-1)} + b_j^{(l)} \right)}{\partial b_j^{(l)}}$$

This is just how a neuron calculates the weighted sum of the inputs (+ bias):
weights * activations previous layer + bias

How do we use this cost to update parameters?

- So, the error for a neuron:

$$\delta_j^{(l)} = \text{error of node } j \text{ in layer } l = \frac{\partial C}{\partial z_j^{(l)}}$$

Using the chain rule:

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \cdot \frac{\partial \left(\sum_{k=1}^{n_{\text{weights}}} w_{jk}^{(l)} \cdot a^{(l-1)} + b_j^{(l)} \right)}{\partial b_j^{(l)}}$$

The whole thing is the partial derivative of z w.r.t. the bias.

How should we nudge the bias such that the cost C decreases?

How does the cost change if we keep everything the same but increase or decrease this bias?

How do we use this cost to update parameters?

- So, the error for a neuron:

$$\delta_j^{(l)} = \text{error of node } j \text{ in layer } l = \frac{\partial C}{\partial z_j^{(l)}}$$

Using the chain rule:

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \cdot \frac{\partial \left(\sum_{k=1}^{n_{\text{weights}}} w_{jk}^{(l)} \cdot a^{(l-1)} + b_j^{(l)} \right)}{\partial b_j^{(l)}}$$

Who can tell me what the partial derivative is?

Hint: $f(x, k) = 12x - 3k$

$$\frac{\partial f(x, k)}{\partial k} = -3$$

How do we use this cost to update parameters?

- So, the error for a neuron:

$$\delta_j^{(l)} = \text{error of node } j \text{ in layer } l = \frac{\partial C}{\partial z_j^{(l)}}$$

Using the chain rule:

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \cdot \frac{\partial \left(\sum_{k=1}^{n_{\text{weights}}} \boxed{w_{jk}^{(l)} \cdot a^{(l-1)}} + b_j^{(l)} \right)}{\partial b_j^{(l)}}$$

Contains no $b_j^{(l)}$, so part. derivative w.r.t.
 $b_j^{(l)} = 0$.

How do we use this cost to update parameters?

- So, the error for a neuron:

$$\delta_j^{(l)} = \text{error of node } j \text{ in layer } l = \frac{\partial C}{\partial z_j^{(l)}}$$

Using the chain rule:

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \cdot \frac{\partial \left(\sum_{k=1}^{n_{\text{weights}}} w_{jk}^{(l)} \cdot a^{(l-1)} + \boxed{b_j^{(l)}} \right)}{\partial b_j^{(l)}}$$

Just one, like the derivative of $f(x) = x$ equals 1.

How do we use this cost to update parameters?

- So, the error for a neuron:

$$\delta_j^{(l)} = \text{error of node } j \text{ in layer } l = \frac{\partial C}{\partial z_j^{(l)}}$$

Using the chain rule:

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \cdot 1 = \delta_j^{(l)}$$

How do we use this cost to update parameters?

- So, the error for a neuron:

$$\delta_j^{(l)} = \text{error of node } j \text{ in layer } l = \frac{\partial C}{\partial z_j^{(l)}}$$

Using the chain rule:

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \cdot 1 = \delta_j^{(l)}$$

Now for the weights:

$$\frac{\partial C}{\partial w_{jk}^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot \frac{\partial \left(\sum_{k=1}^{n_{\text{weights}}} w_{jk}^{(l)} \cdot a^{(l-1)} + b_j^{(l)} \right)}{\partial w_{jk}^{(l)}}$$

How do we use this cost to update parameters?

- So, the error for a neuron:

$$\delta_j^{(l)} = \text{error of node } j \text{ in layer } l = \frac{\partial C}{\partial z_j^{(l)}}$$

Using the chain rule:

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \cdot 1 = \delta_j^{(l)}$$

Now for the weights:

$$\frac{\partial C}{\partial w_{jk}^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot \frac{\partial \left(\sum_{k=1}^{n_{\text{weights}}} w_{jk}^{(l)} \cdot a^{(l-1)} + b_j^{(l)} \right)}{\partial w_{jk}^{(l)}}$$

What is this partial derivative?

How do we use this cost to update parameters?

- So, the error for a neuron:

$$\delta_j^{(l)} = \text{error of node } j \text{ in layer } l = \frac{\partial C}{\partial z_j^{(l)}}$$

Using the chain rule:

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \cdot 1 = \delta_j^{(l)}$$

Now for the weights:

$$\frac{\partial C}{\partial w_{jk}^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot \frac{\partial \left(\sum_{k=1}^{n_{\text{weights}}} \overbrace{w_{jk}^{(l)} \cdot a^{(l-1)}}^{\mathbf{a}^{(l-1)}} + \underbrace{b_j^{(l)}}_0 \right)}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)}$$

How do we use this cost to update parameters?

- This is a lot of math. Why are we doing this again?

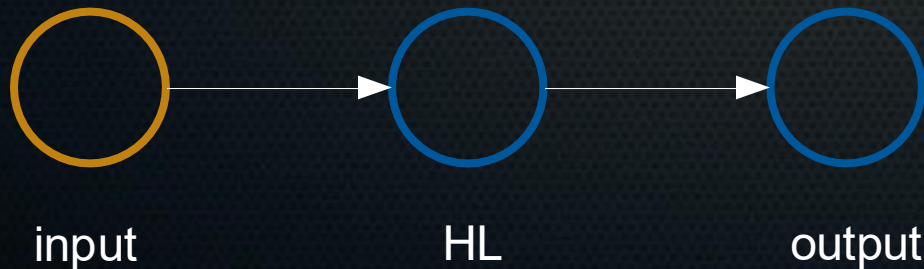
$$\delta_j^{(l)} = -\frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \left| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)} \right.$$

Gradient
descent

Calculate cost of the
current network

Calculate partial derivative
of the cost w.r.t. each parameter

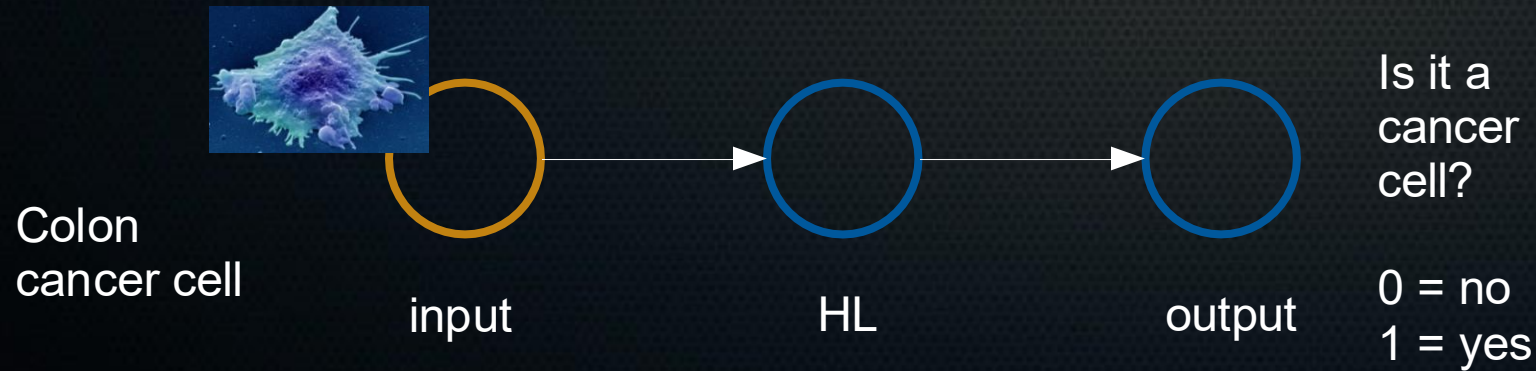
Take a small step in this direction



How do we use this cost to update parameters?

- This is a lot of math. Why are we doing this again?

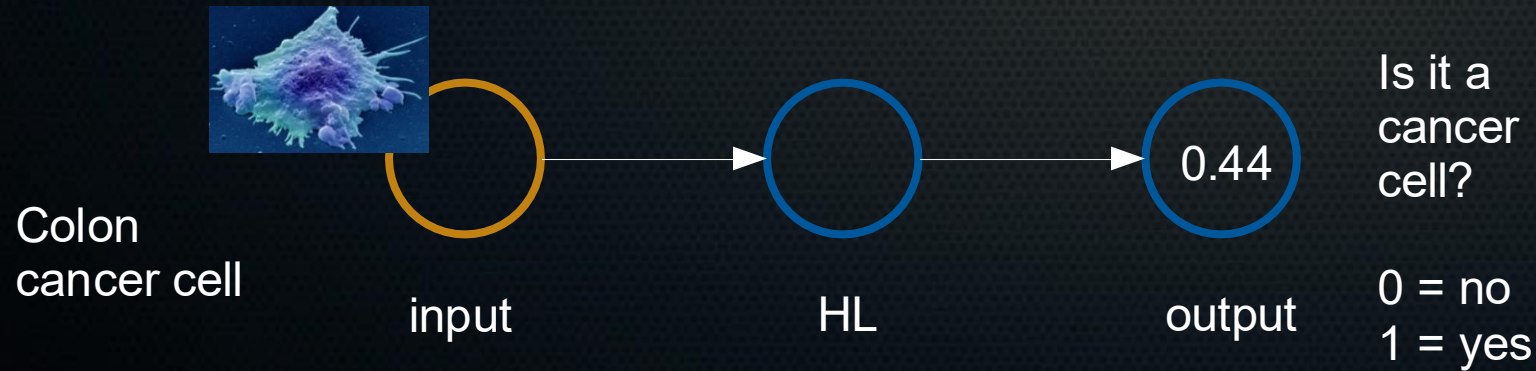
$$\delta_j^{(l)} = -\frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \left| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)} \right.$$



How do we use this cost to update parameters?

- This is a lot of math. Why are we doing this again?

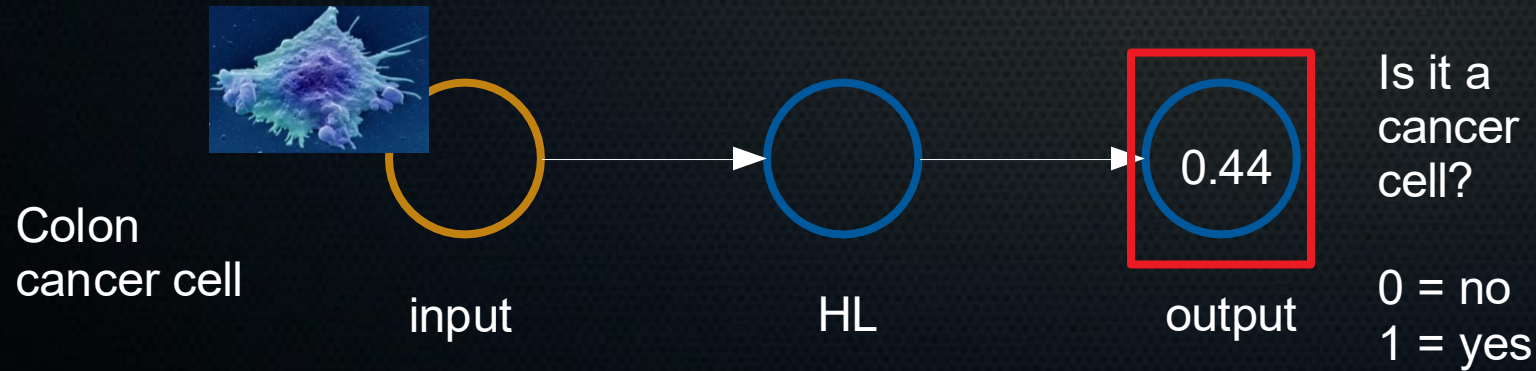
$$\delta_j^{(l)} = -\frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \left| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)} \right.$$



How do we use this cost to update parameters?

- This is a lot of math. Why are we doing this again?

$$\delta_j^{(l)} = -\frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \left| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)} \right.$$

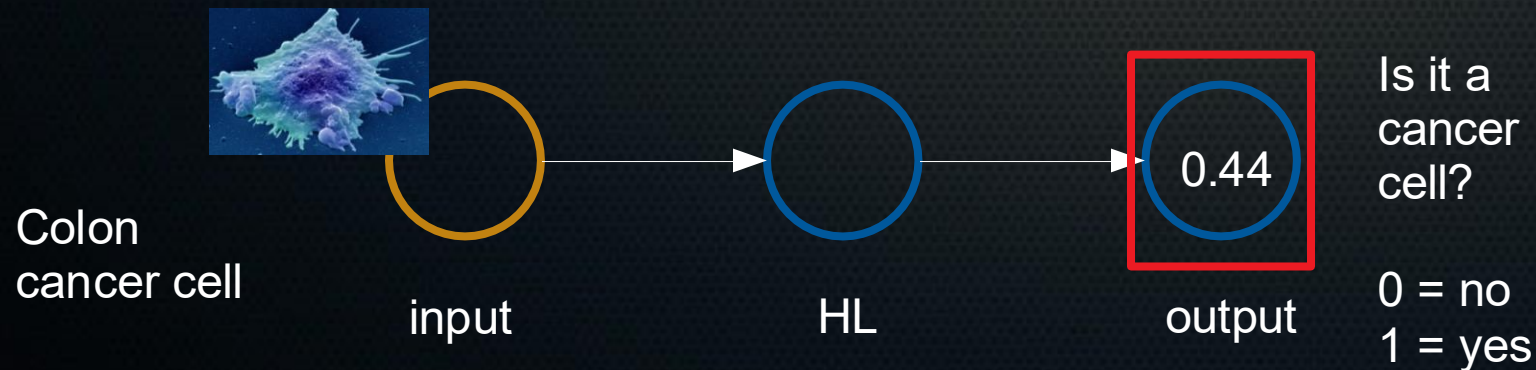


We want it to give 1!

How do we use this cost to update parameters?

- This is a lot of math. Why are we doing this again?

$$\delta_j^{(l)} = -\frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \left| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)} \right.$$



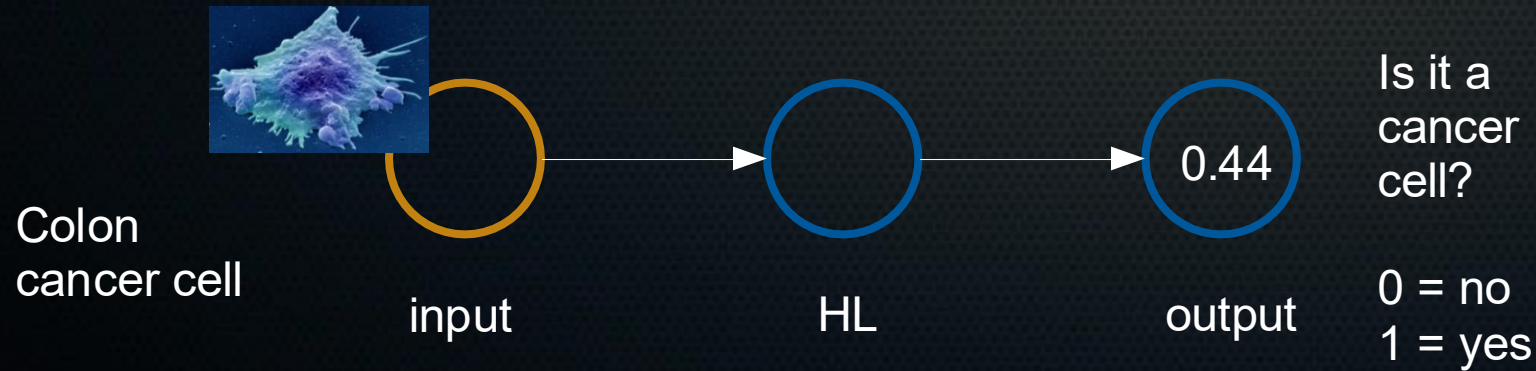
So we calculate a cost.
(For ease of calculation I use MSE, but in reality you would use the cost function explained earlier)

$$MSE = \frac{1}{2} (y - a^{(L)})^2$$
$$MSE = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$

How do we use this cost to update parameters?

- This is a lot of math. Why are we doing this again?

$$\delta_j^{(l)} = -\frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \left| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)} \right.$$



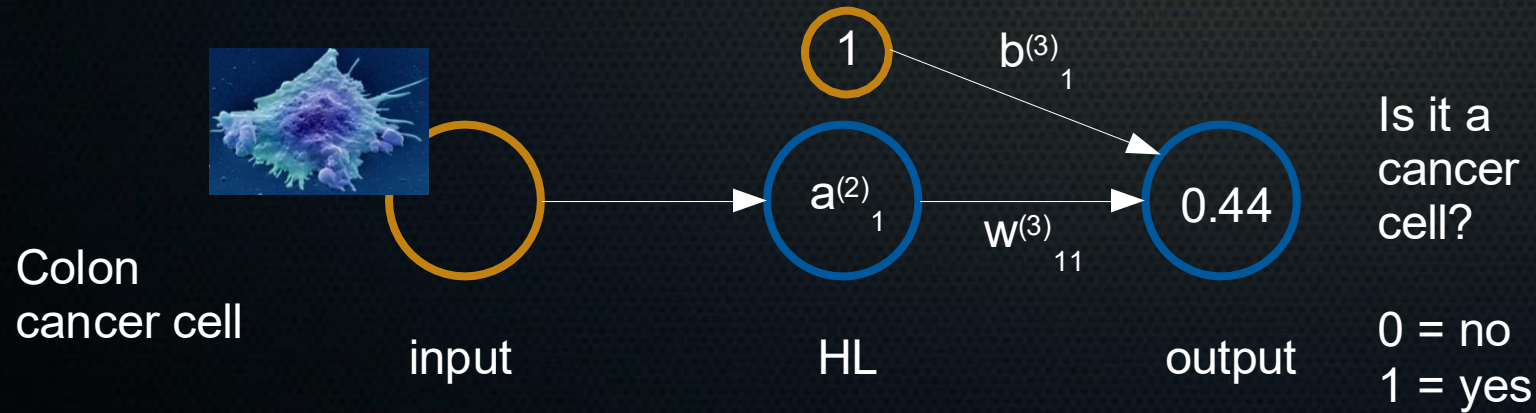
Okay, we know how wrong we are (for this one example). What next?

$$MSE = \frac{1}{2} (y - a^{(L)})^2$$
$$MSE = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$

How do we use this cost to update parameters?

- This is a lot of math. Why are we doing this again?

$$\delta_j^{(l)} = \frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \right| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)}$$



$$MSE = \frac{1}{2} (y - a^{(L)})^2$$

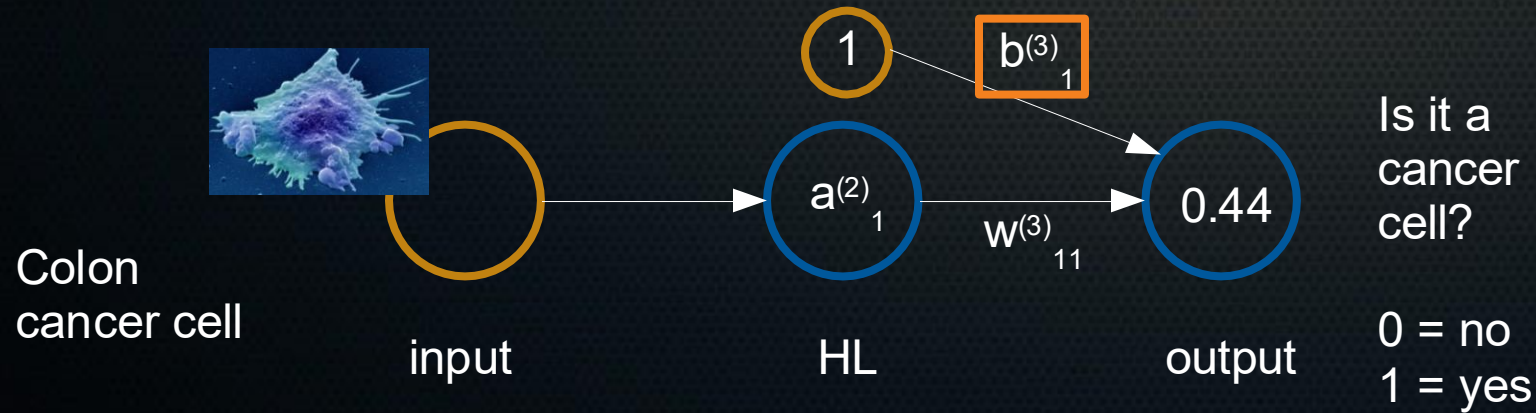
$$MSE = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$

Well, the cost hinges on the activation (0.44). The activation hinges on the weight(s) and bias(es), and on the incoming activations from the previous layer.

How do we use this cost to update parameters?

- This is a lot of math. Why are we doing this again?

$$\delta_j^{(l)} = \frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \right| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)}$$



$$MSE = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$

We can calculate the partial derivative of the cost with respect to:

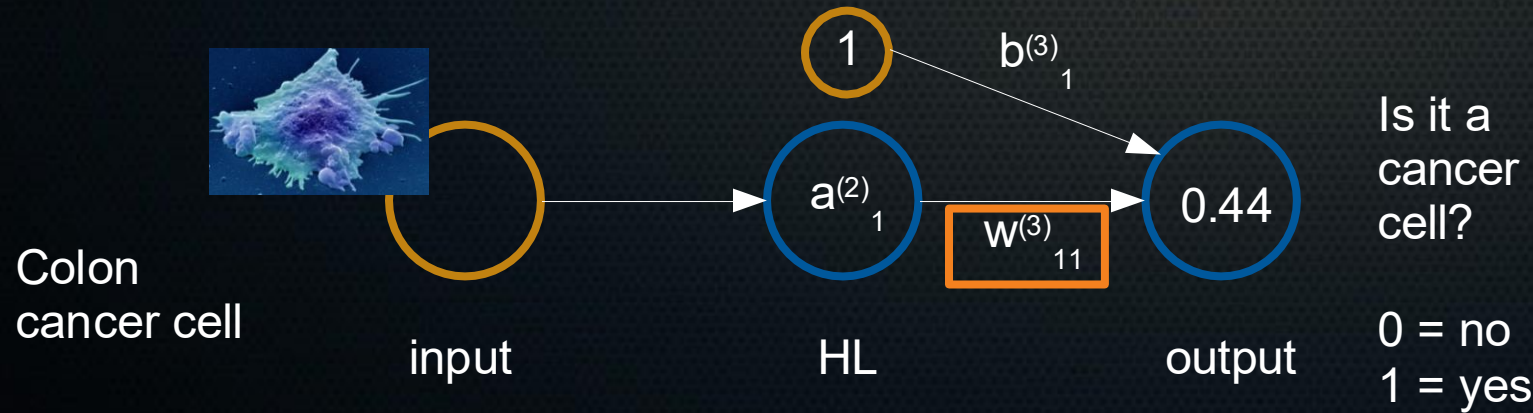
- the bias of the final neuron
- the weight of the final neuron
- the activation of the previous layer that the final neuron takes in

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial a_j^{(l)}} \cdot \frac{\partial a_j^{(l)}}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}}$$

How do we use this cost to update parameters?

- This is a lot of math. Why are we doing this again?

$$\delta_j^{(l)} = \frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \right| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)}$$



$$MSE = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$

$$\frac{\partial C}{\partial w_{jk}^{(l)}} = \frac{\partial C}{\partial a_j^{(l)}} \cdot \frac{\partial a_j^{(l)}}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial w_{jk}^{(l)}}$$

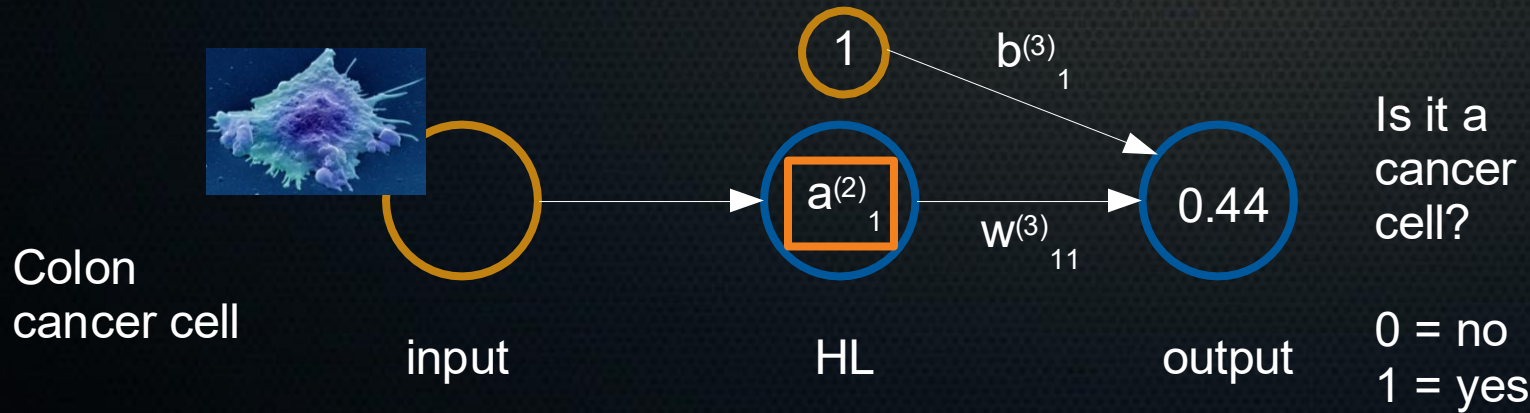
We can calculate the partial derivative of the cost with respect to:

- the bias of the final neuron
- the weight of the final neuron
- the activation of the previous layer that the final neuron takes in

How do we use this cost to update parameters?

- This is a lot of math. Why are we doing this again?

$$\delta_j^{(l)} = \frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \right| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)}$$



$$MSE = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$

$$\frac{\partial C}{\partial a_j^{(l-1)}} = \frac{\partial C}{\partial a_j^{(l)}} \cdot \frac{\partial a_j^{(l)}}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial a_j^{(l-1)}}$$

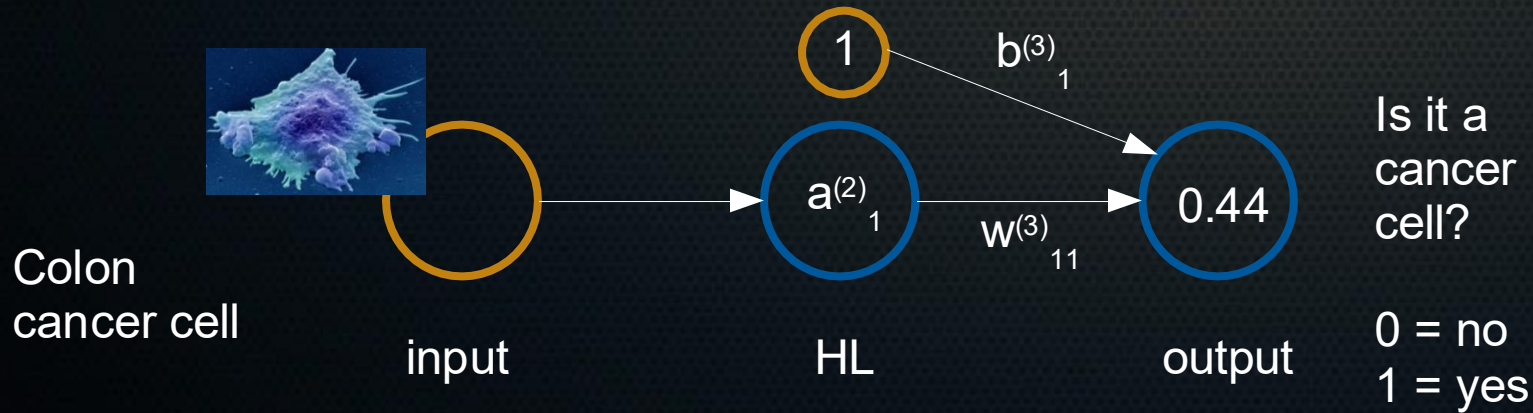
We can calculate the partial derivative of the cost with respect to:

- the bias of the final neuron
- the weight of the final neuron
- the activation of the previous layer that the final neuron takes in

How do we use this cost to update parameters?

- This is a lot of math. Why are we doing this again?

$$\delta_j^{(l)} = \frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \right| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)}$$



$$MSE = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$

$$\frac{\partial C}{\partial a_j^{(l-1)}} = \frac{\partial C}{\partial a_j^{(l)}} \cdot \frac{\partial a_j^{(l)}}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial a_j^{(l-1)}}$$

We can calculate the partial derivative of the cost with respect to:

- the bias of the final neuron
 - the weight of the final neuron
 - the activation of the previous layer that the final neuron takes in
- Here we can take a small step in the direction that minimises it

Gradient descent

Calculate cost of the current network

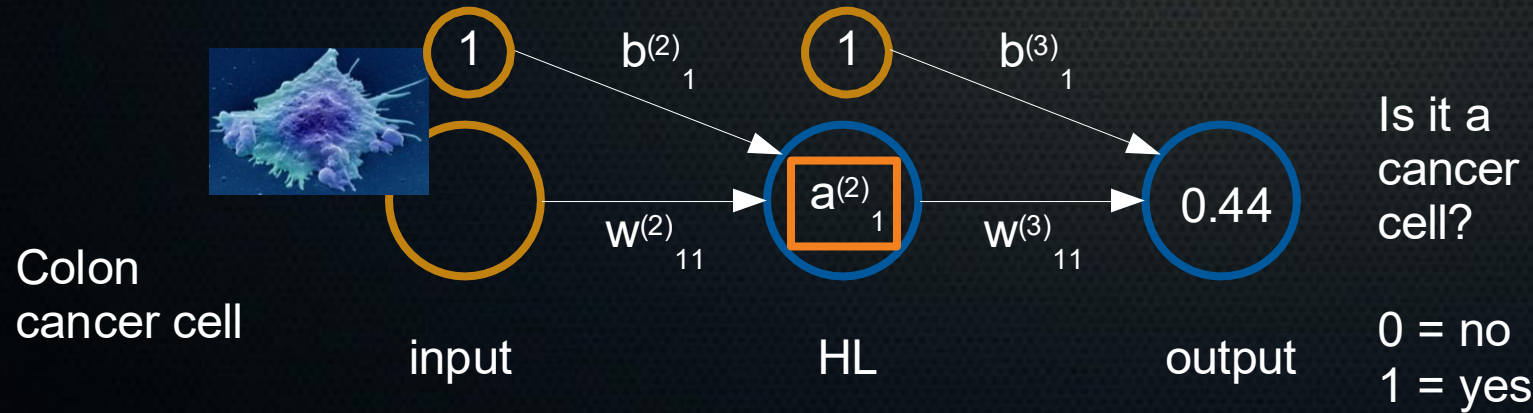
Calculate partial derivative of the cost w.r.t. each parameter

Take a small step in this direction

How do we use this cost to update parameters?

- This is a lot of math. Why are we doing this again?

$$\delta_j^{(l)} = \frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \right| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)}$$



$$MSE = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$

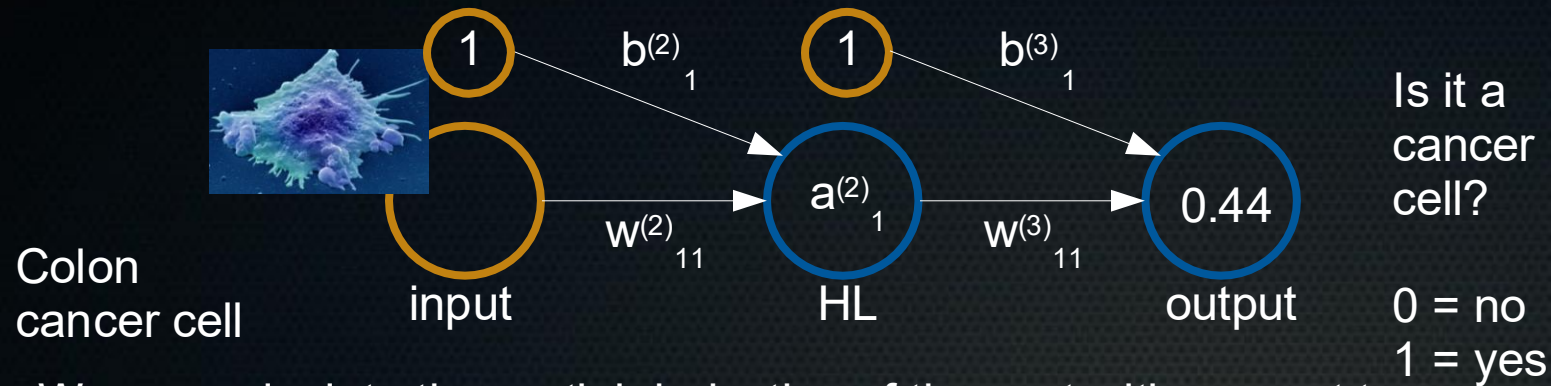
This term we can propagate back!

We know *that* neuron's activation hinges on *its* weights and biases and the previous layer's activation.

We can calculate the partial derivative of the cost with respect to:

- the bias of the final neuron
- the weight of the final neuron
- the activation of the previous layer that the final neuron takes in

How do we use this cost to update parameters?



$$MSE = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$

We can calculate the partial derivative of the cost with respect to:

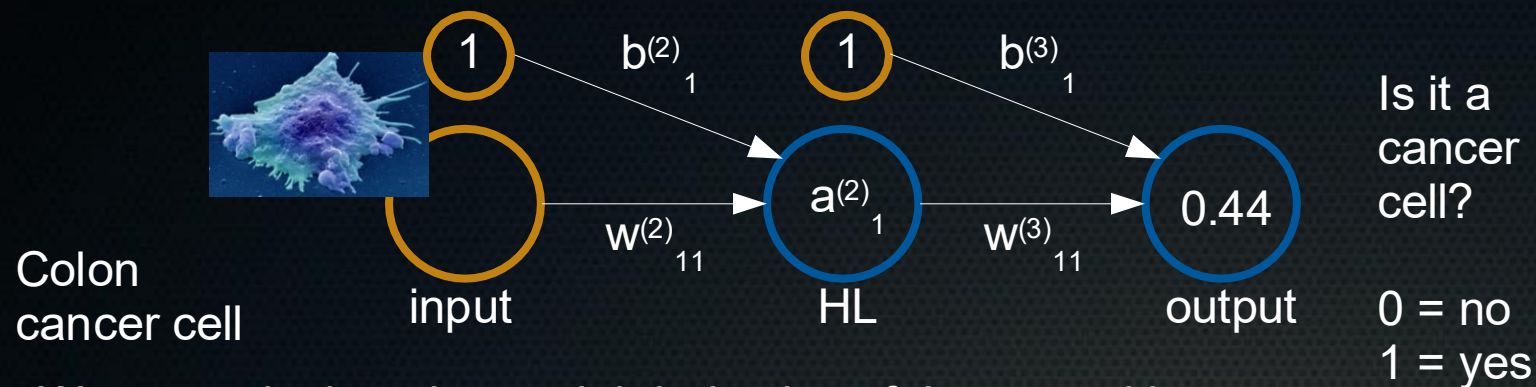
-the bias of the HL neuron

-the weight of the HL neuron

~~-the activation of the previous layer that the HL neuron takes in (but we don't do this for the input!)~~

$$\frac{\partial C}{\partial w_{11}^{(2)}} = \frac{\partial C}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}}$$

How do we use this cost to update parameters?



$$MSE = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$

We can calculate the partial derivative of the cost with respect to:

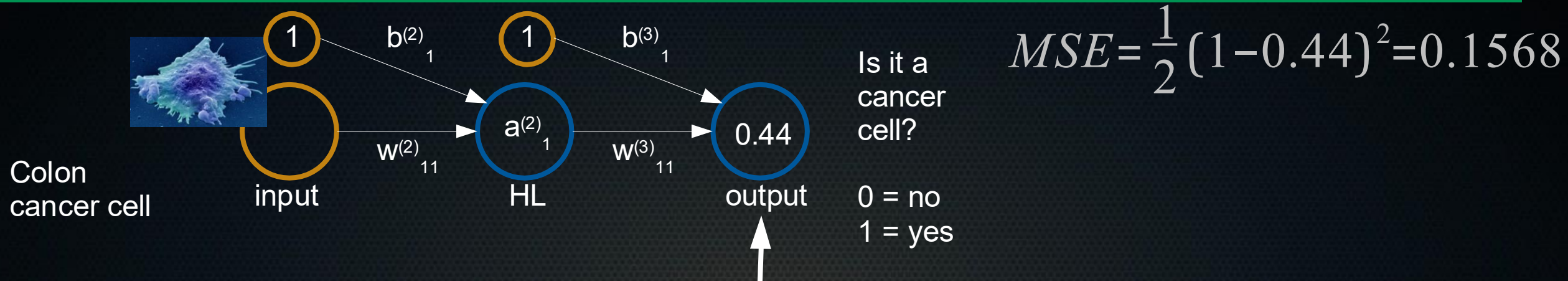
- the bias of the HL neuron
- the weight of the HL neuron

$$\frac{\partial C}{\partial w_{11}^{(2)}} = \frac{\partial C}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}}$$

In words: how the total cost depends on the weight of the HL neuron is:

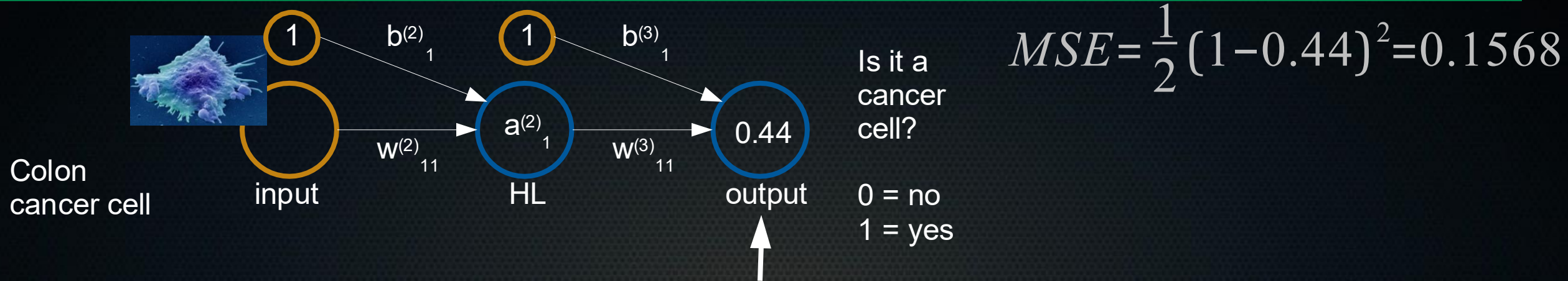
- how the cost depends on the activation of the final layer (0.44), times
- how the activation of the final layer depends on the weighted sum + bias ($z_1^{(3)}$), times
- how that weighted sum depends on the activation of the HL ($a_1^{(2)}$), times
- how the activation of the HL depends on the weighted sum + bias ($z_1^{(2)}$), times
- how that weighted sum depends on the weight from the input.

How do we use this cost to update parameters?



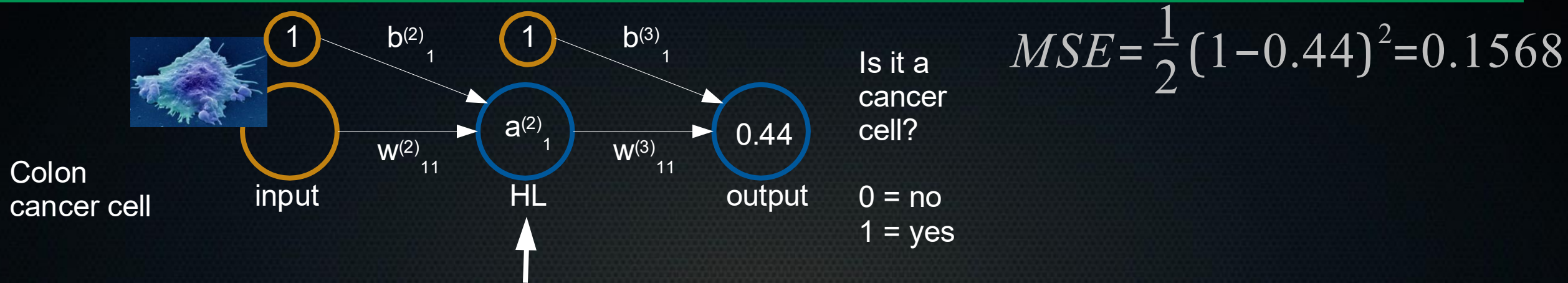
- It might seem complex, but the core idea is simple:
 - We can calculate the cost in the final layer
 - We can calculate how that cost hinges on the parameters of that layer
 - → remember, for gradient descent we want to take a small step in *every* parameter of the network to decrease the cost and make it perform better.

How do we use this cost to update parameters?



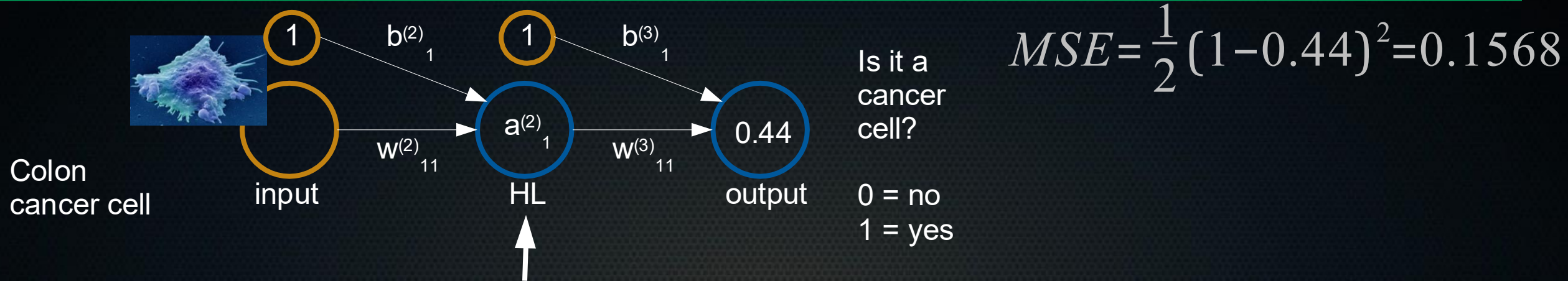
- It might seem complex, but the core idea is simple:
 - We can calculate the cost in the final layer
 - We can calculate how that cost hinges on the parameters of that layer
 - We can also calculate how the cost hinges on the activations of the previous layer

How do we use this cost to update parameters?



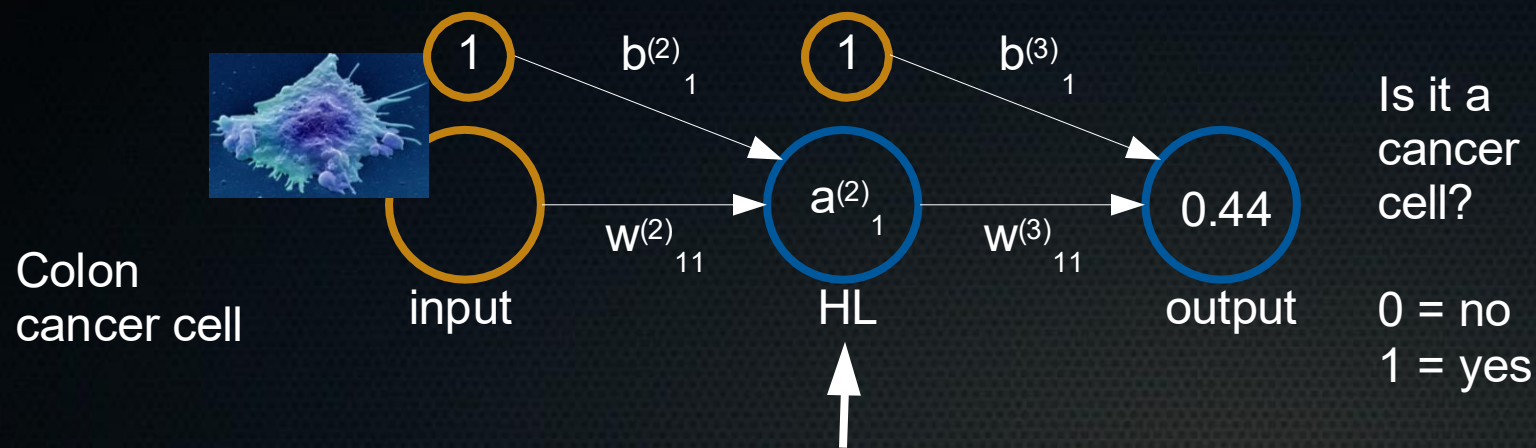
- It might seem complex, but the core idea is simple:
 - We can calculate the cost in the final layer
 - We can calculate how that cost hinges on the parameters of that layer
 - We can also calculate how the cost hinges on the activations of the previous layer
 - We can use *that* to go one step back: once we know how to decrease the cost for the activations of the previous layer, we can figure out how to tweak the parameters there to make that activation better.

How do we use this cost to update parameters?



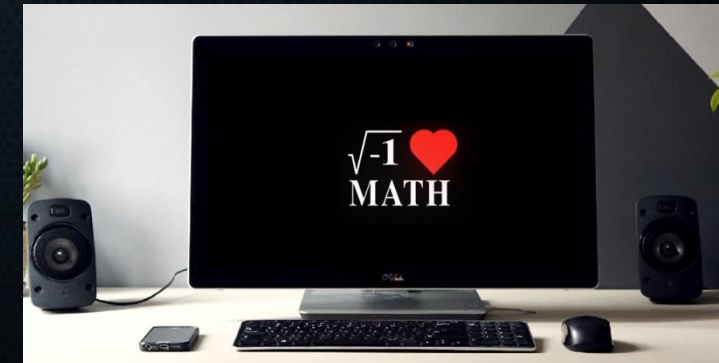
- It might seem complex, but the core idea is simple:
 - We can calculate the cost in the final layer
 - We can calculate how that cost hinges on the parameters of that layer
 - We can also calculate how the cost hinges on the activations of the previous layer
 - We can use *that* to go one step back: once we know how to decrease the cost for the activations of the previous layer, we can figure out how to tweak the parameters there to make that activation better.
 - For this toy network, it ends there, at the input. With more layers, you can keep on chaining derivatives until you know how to change every network parameter for the better!

How do we use this cost to update parameters?



$$MSE = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$

- It might seem complex, but the core idea is simple.
- It's a lot of number-crunching, but computers love that!



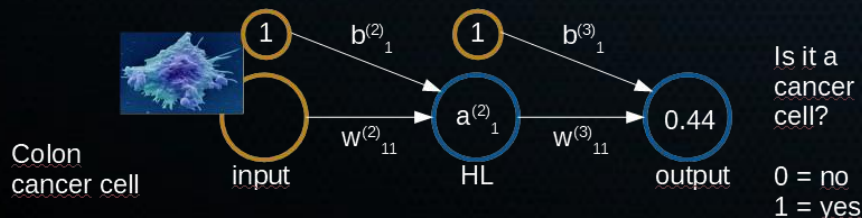
How do we use this cost to update parameters?

- We have almost all the ingredients:

$$\delta_j^{(l)} = \frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \right| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)}$$

- What we still need is how the cost hinges on the activation, and how the activation hinges on its inputs (it was hidden in the left term):

$$\frac{\partial C}{\partial w_{11}^{(2)}} = \boxed{\frac{\partial C}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}}} \cdot \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}}$$

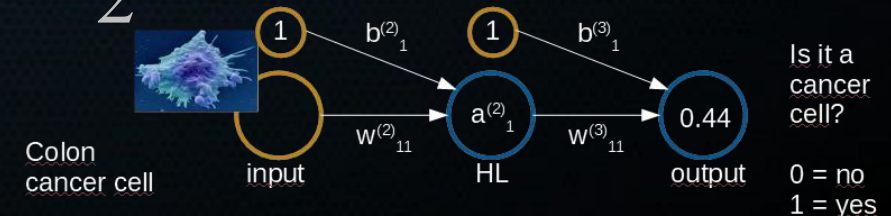


How do we use this cost to update parameters?

- What we still need is how the cost hinges on the activation, and how the activation hinges on its inputs (it was hidden in the left term):

$$\frac{\partial C}{\partial w_{11}^{(2)}} = \boxed{\frac{\partial C}{\partial a_1^{(3)}}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}} \quad C = \frac{1}{2} (y - a_1^{(L)})^2$$

$$C = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$



How do we use this cost to update parameters?

- What we still need is how the cost hinges on the activation, and how the activation hinges on its inputs (it was hidden in the left term):

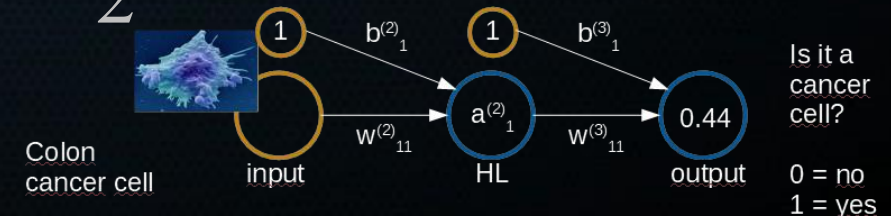
$$\frac{\partial C}{\partial w_{11}^{(2)}} = \boxed{\frac{\partial C}{\partial a_1^{(3)}}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}}$$

$$C = \frac{1}{2} (y - \boxed{a_1^{(L)}})^2$$

$$C = \frac{1}{2} (y - \boxed{\text{sigmoid}(z_j^{(l)})})^2$$

$$C = \frac{1}{2} (y - \boxed{\text{sigmoid}(w_{11}^{(3)} \cdot a_1^{(2)} + b_1^{(3)})})^2$$

$$C = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$



How do we use this cost to update parameters?

- What we still need is how the cost hinges on the activation, and how the activation hinges on its inputs (it was hidden in the left term):

$$\frac{\partial C}{\partial w_{11}^{(2)}} = \frac{\partial C}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}}$$

$$\frac{\partial C}{\partial a_1^{(3)}} = y - a^{(L)} = 1 - 0.44 = 0.56$$

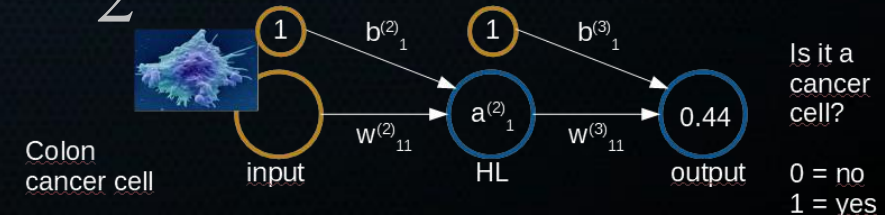
Take the derivative

$$C = \frac{1}{2} (y - a_1^{(L)})^2$$

$$C = \frac{1}{2} (y - \text{sigmoid}(z_j^{(l)}))^2$$

$$C = \frac{1}{2} (y - \text{sigmoid}(w_{11}^{(3)} \cdot a_1^{(2)} + b_1^{(3)}))^2$$

$$C = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$



How do we use this cost to update parameters?

$$f(g(x)) \quad \frac{dy}{dx} = f'(g(x)) \times g'(x)$$

- What we still need is how the cost hinges on the activation, and how the activation hinges on its inputs (it was hidden in the left term):

$$\frac{\partial C}{\partial w_{11}^{(2)}} = \frac{\partial C}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}}$$

$$\frac{\partial C}{\partial a_1^{(3)}} = y - a^{(L)} = 1 - 0.44 = 0.56$$

$$\frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} = \sigma'(z_1^{(3)}) = \sigma(z_1^{(3)}) \cdot (1 - \sigma(z_1^{(3)}))$$

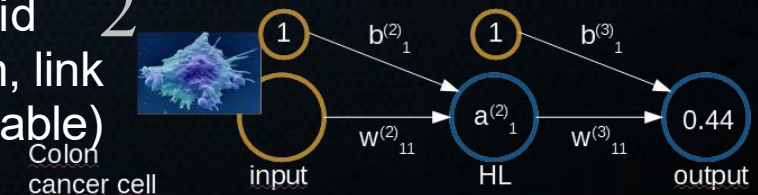
Derivative of sigmoid
(complex derivation, link
with extra info available)

$$C = \frac{1}{2} (y - a_1^{(L)})^2$$

$$C = \frac{1}{2} (y - \text{sigmoid}(z_j^{(l)}))^2$$

$$C = \frac{1}{2} (y - \text{sigmoid}(w_{11}^{(3)} \cdot a_1^{(2)} + b_1^{(3)}))^2$$

$$C = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$



Is it a
cancer
cell?

0 = no
1 = yes

How do we use this cost to update parameters?

- What we still need is how the cost hinges on the activation, and how the activation hinges on its inputs (it was hidden in the left term):

$$\frac{\partial C}{\partial w_{11}^{(2)}} = \frac{\partial C}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}}$$

$$\frac{\partial C}{\partial a_1^{(3)}} = y - a^{(L)} = 1 - 0.44 = 0.56$$

$$\frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} = \sigma'(z_1^{(3)}) = \sigma(z_1^{(3)}) \cdot (1 - \sigma(z_1^{(3)}))$$

$$\sigma(-0.241) \approx 0.44$$

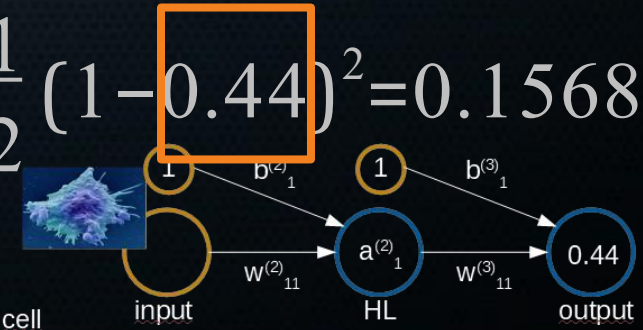
$$C = \frac{1}{2} (y - a_1^{(L)})^2$$

$$C = \frac{1}{2} (y - \text{sigmoid}(z_j^{(l)}))^2$$

$$C = \frac{1}{2} (y - \text{sigmoid}(w_{11}^{(3)} \cdot a_1^{(2)} + b_1^{(3)}))^2$$

$$C = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$

Derivative of sigmoid
(complex derivation, link
with extra info available)



Is it a
cancer
cell?

0 = no
1 = yes

How do we use this cost to update parameters?

- What we still need is how the cost hinges on the activation, and how the activation hinges on its inputs (it was hidden in the left term):

$$\frac{\partial C}{\partial w_{11}^{(2)}} = \frac{\partial C}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}}$$

$$\frac{\partial C}{\partial a_1^{(3)}} = y - a^{(L)} = 1 - 0.44 = 0.56$$

$$\frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} = \sigma'(z_1^{(3)}) = \sigma(z_1^{(3)}) \cdot (1 - \sigma(z_1^{(3)}))$$

$$\sigma'(-0.241) = \sigma(-0.241) \cdot (1 - \sigma(-0.241)) \approx 0.2464$$

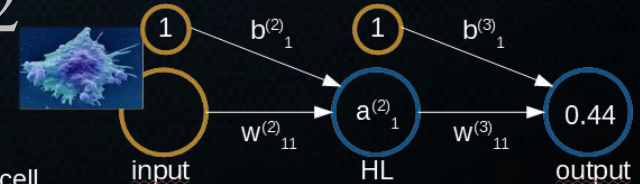
$$C = \frac{1}{2} (y - a_1^{(L)})^2$$

$$C = \frac{1}{2} (y - \text{sigmoid}(z_j^{(l)}))^2$$

$$C = \frac{1}{2} (y - \text{sigmoid}(w_{11}^{(3)} \cdot a_1^{(2)} + b_1^{(3)}))^2$$

$$C = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$

Derivative of sigmoid
(complex derivation, link
with extra info available)



Is it a
cancer
cell?
0 = no
1 = yes

How do we use this cost to update parameters?

- What we still need is how the cost hinges on the activation, and how the activation hinges on its inputs (it was hidden in the left term):

$$\frac{\partial C}{\partial w_{11}^{(2)}} = \frac{\partial C}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}}$$

$$\frac{\partial C}{\partial a_1^{(3)}} = y - a^{(L)} = 1 - 0.44 = 0.56$$

$$\frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} = 0.2464$$

$$\sigma(-0.241) \approx 0.44$$

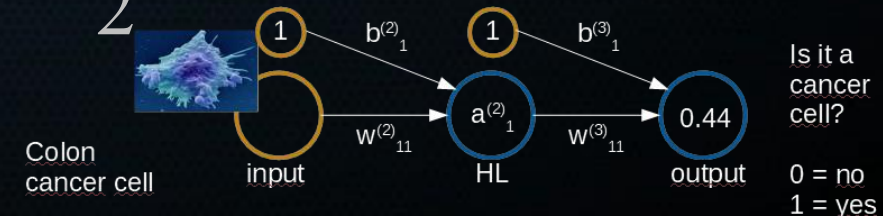
So we can keep on calculating, calculating partial derivatives for weights and biases in the current layer, and propagating the error back to get them for previous layers iteratively

$$C = \frac{1}{2} (y - a_1^{(L)})^2$$

$$C = \frac{1}{2} (y - \text{sigmoid}(z_j^{(l)}))^2$$

$$C = \frac{1}{2} (y - \text{sigmoid}(w_{11}^{(3)} \cdot a_1^{(2)} + b_1^{(3)}))^2$$

$$C = \frac{1}{2} (1 - 0.44)^2 = 0.1568$$



How do we use this cost to update parameters?

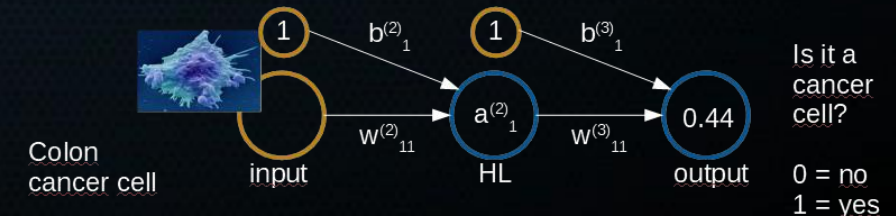
- Taken together:

$$\delta_j^{(l)} = \frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \left| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)} \right.$$

For the
final layer

$$\delta_j^{(L)} = \frac{\partial C}{\partial z_j^{(L)}} = (a_j^{(L)} - y_j) \sigma(z_j^{(L)}) \cdot (1 - \sigma(z_j^{(L)}))$$

* (For cost function == MSE $(a^{(L)}_j - y_j)$, just $a^{(L)}_j - y_j$ for binary cross-entropy)



How do we use this cost to update parameters?

- Taken together:

$$\delta_j^{(l)} = \frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \right| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a_k^{(l-1)}$$

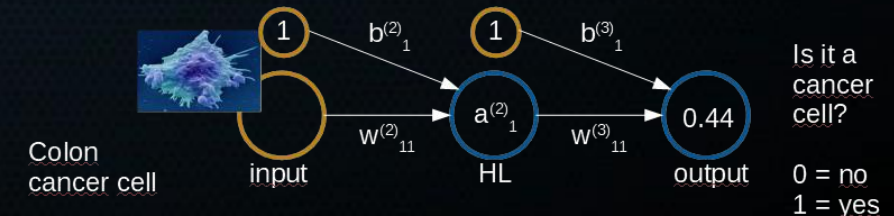
For the
final layer

For the
previous layer
(applied iteratively)

$$\delta_j^{(L)} = \frac{\partial C}{\partial z_j^{(L)}} = (a_j^{(L)} - y_j) \sigma'(z_j^{(L)}) \quad *$$

$$\delta_j^{(L-1)} = \frac{\partial C}{\partial z_j^{(L-1)}} = \sum_{k=1}^{n_{weights}} w_{kj}^{(L)} \cdot \delta_k^{(L)} \cdot \sigma'(z_j^{(L-1)}) \quad \text{Etc.}$$

* (For cost function == MSE $(a^{(L)}_j - y_j)$, slightly more complex for binary crossentropy)



How do we use this cost to update parameters?

- Taken together:

$$\delta_j^{(l)} = \frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \right| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a_k^{(l-1)}$$

For the
final layer

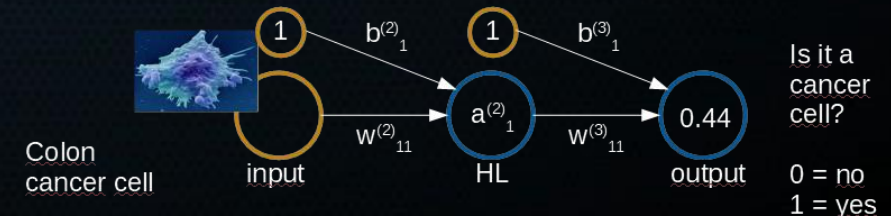
For the
previous layer
(applied iteratively)

$$\delta_j^{(L)} = \frac{\partial C}{\partial z_j^{(L)}} = (a_j^{(L)} - y_j) \sigma'(z_j^{(L)})$$

Moves the error from the inputs of neurons in layer $l+1$ to the outputs of a neuron in the previous layer

$$\delta_j^{(L-1)} = \frac{\partial C}{\partial z_j^{(L-1)}} = \sum_{k=1}^{n_{weights}} w_{kj}^{(L)} \cdot \delta_k^{(L)} \cdot \sigma'(z_j^{(L-1)})$$

* (For cost function == MSE $(a^{(L)} - y_j)$, slightly more complex for binary crossentropy)



How do we use this cost to update parameters?

- Taken together:

$$\delta_j^{(l)} = \frac{\partial C}{\partial z_j^{(l)}} \quad \left| \quad \frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad \right| \quad \frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a^{(l-1)}$$

For the
final layer

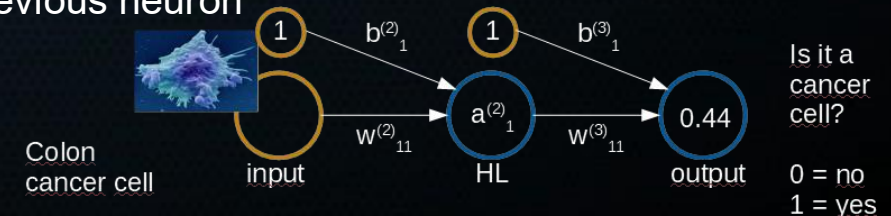
For the
previous layer
(applied iteratively)

$$\delta_j^{(L)} = \frac{\partial C}{\partial z_j^{(L)}} = (a_j^{(L)} - y_j) \sigma'(z_j^{(L)})$$

Moves error back through the activation function, so that we can calculate the partial derivative w.r.t. weights, bias, and input of previous neuron

$$\delta_j^{(L-1)} = \frac{\partial C}{\partial z_j^{(L-1)}} = \sum_{k=1}^{n_{weights}} w_{kj}^{(L)} \cdot \delta_k^{(L)} \cdot \sigma'(z_j^{(L-1)})$$

* (For cost function == MSE $(a^{(L)}_j - y_j)$, slightly more complex for binary crossentropy)



How do we use this cost to update parameters?

- Taken together:

$$\delta_j^{(l)} = -\frac{\partial C}{\partial z_j^{(l)}}$$

For the final layer

$$\delta_j^{(L)} = \frac{\partial C}{\partial z_j^{(L)}} = (a_j^{(L)} - y_j) \sigma'(z_j^{(L)})$$

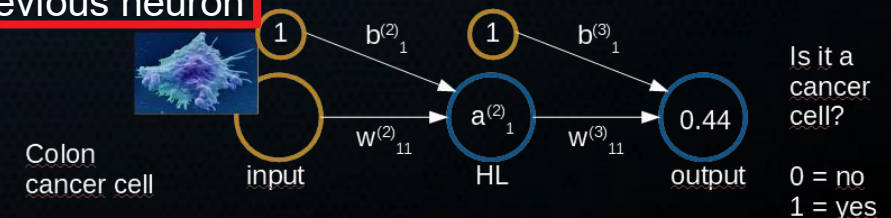
For the previous layer (applied iteratively)

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial C}{\partial z_j^{(l)}} \cdot \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} = \delta_j^{(l)}$$

$$\frac{\partial C}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} \cdot a_k^{(l-1)}$$

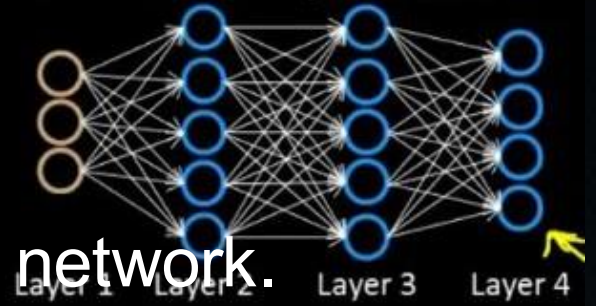
Moves error back through the activation function, so that we can calculate the partial derivative w.r.t. **weights**, **bias**, and **input of previous neuron**

$$\delta_j^{(L-1)} = \frac{\partial C}{\partial z_j^{(L-1)}} = \sum_{k=1}^{n_{weights}} w_{kj}^{(L)} \cdot \delta_k^{(L)} \cdot \sigma'(z_j^{(L-1)})$$



How do we use this cost to update parameters?

- The procedure we just discussed leads to calculating the partial derivatives of the cost function with respect to all the weights and biases in the network.

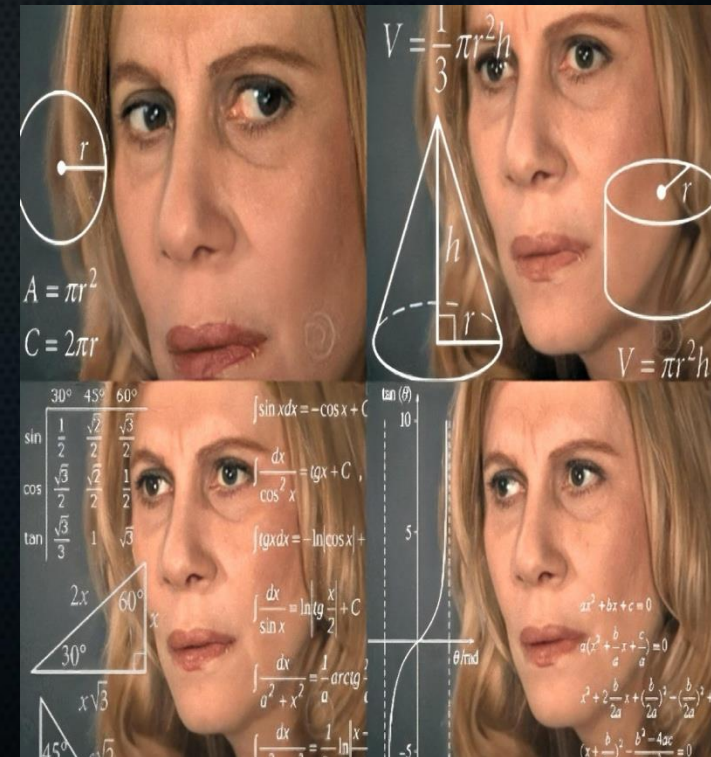


And now?

- Most equations I gave here were not vectorised, but they can all be vectorised. That makes them fast(er) to compute.
- In the example here, we had one training sample. In reality you have more, so you compute an average cost over all training samples and *then* update the network using backpropagation.

And now?

- You might feel like the lady in this distinctly non-spicy stale meme, though without trigonometrical or geometric concerns.
- I for sure did not get this quickly. It's not simple!
- Modern libraries do all the heavy lifting for you.



And now?

- Let this sink in, and continue where you left off yesterday

