

# Backpropagation Guided Notes

Backpropagation is a critical piece of modern deep learning. To really get a grasp of how backpropagation works, there's nothing quite like deriving the equations for yourself. Let's do it.

Data	Architecture	Forward Equations
		$z^{(2)} = XW^{(1)} \quad (1)$
		$a^{(2)} = f(z^{(2)}) \quad (2)$
		$z^{(3)} = a^{(2)}W^{(2)} \quad (3)$
		$\hat{y} = f(z^{(3)}) \quad (4)$
		$J = \sum \frac{1}{2}(y - \hat{y})^2 \quad (5)$

Code Symbol	Math Symbol	Definition	Dimensions
X	X	Input Data, each row in an example	(numExamples, inputLayerSize)
y	y	target data	(numExamples, outputLayerSize)
W1	W <sup>(1)</sup>	Layer 1 weights	(inputLayerSize, hiddenLayerSize)
W2	W <sup>(2)</sup>	Layer 2 weights	(hiddenLayerSize, outputLayerSize)
z2	z <sup>(2)</sup>	Layer 2 activation	(numExamples, hiddenLayerSize)
a2	a <sup>(2)</sup>	Layer 2 activity	(numExamples, hiddenLayerSize)
z3	z <sup>(3)</sup>	Layer 3 activation	(numExamples, outputLayerSize)
J	J	Cost	(1, outputLayerSize)
dJdz3	$\frac{\partial J}{\partial z^{(3)}} = \delta^{(3)}$	Partial derivative of cost with respect to z <sup>(3)</sup>	
dJdW2	$\frac{\partial J}{\partial W^{(2)}}$	Partial derivative of cost with respect to W <sup>(2)</sup>	
dz3dz2	$\frac{\partial z^{(3)}}{\partial z^{(2)}}$	Partial derivative of z <sup>(3)</sup> with respect to z <sup>(2)</sup>	
dJdW1	$\frac{\partial J}{\partial W^{(1)}}$	Partial derivative of cost with respect to W <sup>(1)</sup>	
delta2	$\delta^{(2)}$	Backpropagating Error 2	
delta3	$\delta^{(3)}$	Backpropagating Error 1	

For you to figure out!

# Your Mission

$$\frac{\partial J}{\partial W^{(1)}} = ? \quad \frac{\partial J}{\partial W^{(2)}} = ?$$

1. The dimension of  $\frac{\partial J}{\partial W^{(1)}}$  is \_\_\_\_\_.

2. The dimension of  $\frac{\partial J}{\partial W^{(2)}}$  is \_\_\_\_\_.

3. Using (5), we can write  $\frac{\partial J}{\partial W^{(2)}} = \frac{\partial \sum \frac{1}{2}(y - \hat{y})^2}{\partial W^{(2)}}.$

Use the sum rule for differentiation to move the summation outside the gradient:

$$\frac{\partial J}{\partial W^{(2)}} =$$

4. Let's temporarily remove the summation, and consider  $\frac{\partial J}{\partial W^{(2)}}$  in terms of just one example (numExamples = 1). Using the chain rule, derive an expression for  $\frac{\partial J}{\partial W^{(2)}}$  in terms of  $y, \hat{y}, \frac{\partial \hat{y}}{\partial W^{(2)}}$ .

$$\frac{\partial J}{\partial W^{(2)}} =$$

5. Now, use the chain rule again to express  $\frac{\partial J}{\partial W^{(2)}}$  in terms of  $y, \hat{y}, \frac{\partial \hat{y}}{\partial z^{(3)}}, \frac{\partial z^{(3)}}{\partial W^{(2)}}.$

$$\frac{\partial J}{\partial W^{(2)}} =$$

6.  $\hat{y}$  and  $z^{(3)}$  are connected by our sigmoid activation function  $f(z) = \frac{1}{1 + e^{-z}}.$

$$\frac{\partial \hat{y}}{\partial z^{(3)}} = f'(z) =$$

You should now have an equation that looks something like this:

$$\frac{\partial J}{\partial W^{(2)}} = -(y - \hat{y}) f'(z^{(3)}) \frac{\partial z^{(3)}}{\partial W^{(2)}}$$

To simplify our equations a little, let's introduce a new term, the "backpropogating error":

$$\delta^{(3)} = -(y - \hat{y}) f'(z^{(3)})$$

7. What is the dimension of  $\delta^{(3)}$  ?

8. Now we need to work on  $\frac{\partial z^{(3)}}{\partial W^{(2)}}$ . To get started, write out the full matrix equation for (3), using numExamples = 1, and inputLayerSize = 2, hiddenLayerSize = 3, and outputLayerSize = 1.

9. Now, using your calculus skills:

$$\frac{\partial z^{(3)}}{\partial W^{(2)}} = \begin{bmatrix} \frac{\partial z^{(3)}}{\partial W_{11}^{(2)}} \\ \frac{\partial z^{(3)}}{\partial W_{21}^{(2)}} \\ \frac{\partial z^{(3)}}{\partial W_{31}^{(2)}} \end{bmatrix} = \begin{bmatrix} \quad \\ \quad \\ \quad \end{bmatrix}$$

10. Now, write  $\frac{\partial z^{(3)}}{\partial W^{(2)}}$  in terms of the vector  $a^{(2)}$ :

$$\frac{\partial z^{(3)}}{\partial W^{(2)}} =$$

11.  $W^{(1)} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$     $W^{(2)} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$     $X = [0.3 \ 1.0]$     $\frac{\partial J}{\partial W^{(2)}} = ?$

12. Next, let's deal with the `numExamples > 1` case. Back in question 4 we temporarily took away the summation, we'll figure out how to re-introduce it now. To get started, write out the full matrix equation for (3), using `numExamples = 3`, and `inputLayerSize = 2`, `hiddenLayerSize = 3`, and `outputLayerSize = 1`.

What do the rows and columns of your "a" matrix represent?

13. Now that we've let `numExamples=3`, what is the dimension of  $\delta^{(3)}$ ?

14. Almost there! Now, sum across our examples in terms of the individual elements of  $\delta^{(3)}$  and  $a^{(2)}$ :  
 (Hint: to compute the gradients with respect to each of our 3 weights in  $W_2$ , **we need to add across our examples.**)

$$\frac{\partial J}{\partial W^{(2)}} = \begin{bmatrix} \frac{\partial J}{\partial W_{11}^{(2)}} \\ \frac{\partial J}{\partial W_{21}^{(2)}} \\ \frac{\partial J}{\partial W_{31}^{(2)}} \end{bmatrix} = \begin{bmatrix} \end{bmatrix}$$

15. Now express the above operation in terms of the matrix  $a^{(2)}$  and the vector  $\delta^{(3)}$ .

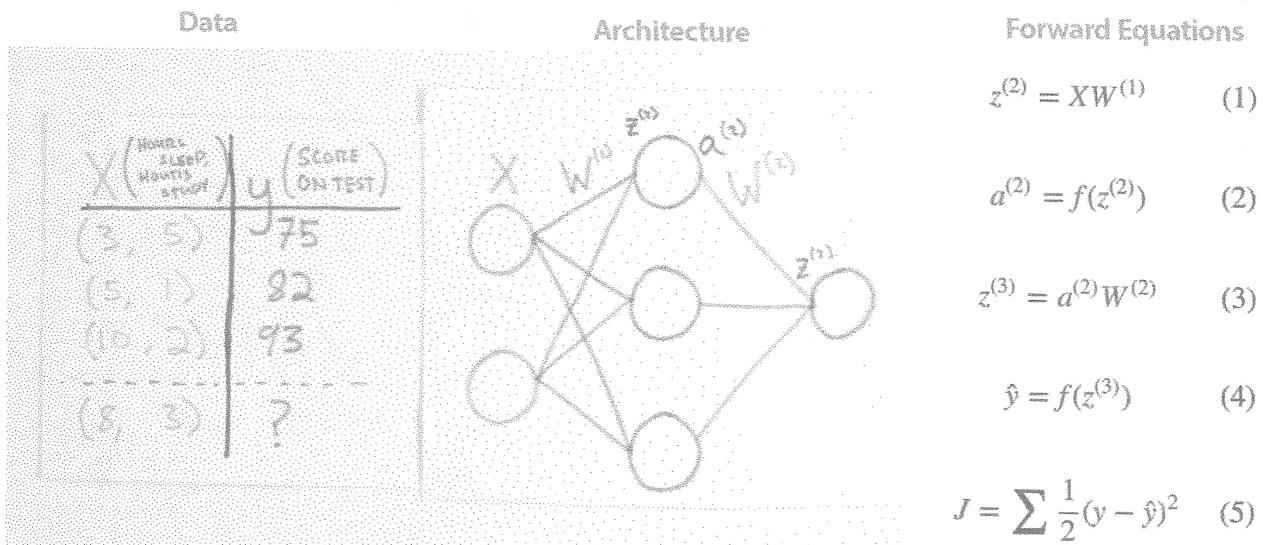
$$\frac{\partial J}{\partial W^{(2)}} =$$

$$16. \quad W^{(1)} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad W^{(2)} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad X = \begin{bmatrix} 0.3 & 1.0 \\ 0.5 & 0.2 \\ 1.0 & 0.4 \end{bmatrix} \quad \frac{\partial J}{\partial W^{(2)}} = ?$$

17. Derive an expression for  $\frac{\partial J}{\partial W^{(1)}}$  by continuing to propagate errors backwards through our network.

# Backpropagation Guided Notes

Propagation is a critical piece of modern deep learning. To really get a grasp of how backpropagation here's nothing quite like deriving the equations for yourself. Let's do it.



Code Symbol	Math Symbol	Definition	Dimensions
X	X	Input Data, each row in an example	(numExamples, inputLayerSize)
y	y	target data	(numExamples, outputLayerSize)
W1	W <sup>(1)</sup>	Layer 1 weights	(inputLayerSize, hiddenLayerSize)
W2	W <sup>(2)</sup>	Layer 2 weights	(hiddenLayerSize, outputLayerSize)
z2	z <sup>(2)</sup>	Layer 2 activation	(numExamples, hiddenLayerSize)
a2	a <sup>(2)</sup>	Layer 2 activity	(numExamples, hiddenLayerSize)
z3	z <sup>(3)</sup>	Layer 3 activation	(numExamples, outputLayerSize)
J	J	Cost	(1, outputLayerSize)
dJdz3	$\frac{\partial J}{\partial z^{(3)}} = \delta^{(3)}$	Partial derivative of cost with respect to z <sup>(3)</sup>	(3x1) <i>one for each example.</i>
dJdW2	$\frac{\partial J}{\partial W^{(2)}}$	Partial derivative of cost with respect to W <sup>(2)</sup>	(3x1)
dz3dz2	$\frac{\partial z^{(3)}}{\partial z^{(2)}}$	Partial derivative of z <sup>(3)</sup> with respect to z <sup>(2)</sup>	(3x3)
dJdW1	$\frac{\partial J}{\partial W^{(1)}}$	Partial derivative of cost with respect to W <sup>(1)</sup>	(2x3)
delta2	$\delta^{(2)}$	Backpropagating Error 2	
delta3	$\delta^{(3)}$	Backpropagating Error 1	(3x1) numExamples $\nwarrow$ output layer size.

For you to figure out!

## Your Mission

$$\frac{\partial J}{\partial W^{(1)}} = ? \quad \frac{\partial J}{\partial W^{(2)}} = ?$$

1. The dimension of  $\frac{\partial J}{\partial W^{(1)}}$  is  $(2, 3) \leftarrow \text{SAME AS } W^{(1)}$ .

2. The dimension of  $\frac{\partial J}{\partial W^{(2)}}$  is  $(3, 1) \leftarrow \text{SAME AS } W^{(2)}$ .

3. Using (5), we can write  $\frac{\partial J}{\partial W^{(2)}} = \frac{\partial \sum \frac{1}{2}(y - \hat{y})^2}{\partial W^{(2)}}$ .

Use the sum rule for differentiation to move the summation outside the gradient:

$$\frac{\partial J}{\partial W^{(2)}} = \sum \frac{1}{2} \frac{\partial(y - \hat{y})^2}{\partial W^{(2)}}$$

4. Let's temporarily remove the summation, and consider  $\frac{\partial J}{\partial W^{(2)}}$  in terms of just one example (numExamples = 1). Using the chain rule, derive an expression for  $\frac{\partial J}{\partial W^{(2)}}$  in terms of  $y, \hat{y}, \frac{\partial \hat{y}}{\partial W^{(2)}}$ .

$$\frac{\partial J}{\partial W^{(2)}} = -(y - \hat{y}) \cdot \frac{\partial \hat{y}}{\partial W^{(2)}} \quad (\hat{y} \text{ is constant}).$$

5. Now, use the chain rule again to express  $\frac{\partial J}{\partial W^{(2)}}$  in terms of  $y, \hat{y}, \frac{\partial \hat{y}}{\partial z^{(3)}}, \frac{\partial z^{(3)}}{\partial W^{(2)}}$ .

$$\frac{\partial J}{\partial W^{(2)}} = -(y - \hat{y}) \cdot \frac{\partial \hat{y}}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial W^{(2)}}$$

6.  $\hat{y}$  and  $z^{(3)}$  are connected by our sigmoid activation function  $f(z) = \frac{1}{1 + e^{-z}}$ .

$$\begin{aligned} \frac{\partial \hat{y}}{\partial z^{(3)}} &= f'(z) = ((1 + e^{-z})^{-1})' = -(1 + e^{-z})^{-2} \cdot -e^{-z} \\ &= \frac{e^{-z}}{(1 + e^{-z})^2} \end{aligned}$$

You should now have an equation that looks something like this:

3

$$\frac{\partial J}{\partial W^{(2)}} = -(y - \hat{y}) f'(z^{(3)}) \frac{\partial z^{(3)}}{\partial W^{(2)}} \quad \checkmark$$

To simplify our equations a little, let's introduce a new term, the "backpropogating error":

$$\delta^{(3)} = -(y - \hat{y}) f'(z^{(3)})$$

$\begin{pmatrix} 1 \\ 1 \end{pmatrix}$  NUMBER OF EXAMPLES.

7. What is the dimension of  $\delta^{(3)}$ ?

8. Now we need to work on  $\frac{\partial z^{(3)}}{\partial W^{(2)}}$ . To get started, write out the full matrix equation for (3), using numExamples = 1, and inputLayerSize = 2, hiddenLayerSize = 3, and outputLayerSize = 1.

$$z^{(3)} = a^{(2)} W^{(2)} \Rightarrow z^{(3)} = \begin{bmatrix} a_{11}^{(2)} & a_{12}^{(2)} & a_{13}^{(2)} \end{bmatrix} \begin{bmatrix} W_{11}^{(2)} \\ W_{21}^{(2)} \\ W_{31}^{(2)} \end{bmatrix}$$

9. Now, using your calculus skills:

$$\frac{\partial z^{(3)}}{\partial W^{(2)}} = \begin{bmatrix} \frac{\partial z^{(3)}}{\partial W_{11}^{(2)}} \\ \frac{\partial z^{(3)}}{\partial W_{21}^{(2)}} \\ \frac{\partial z^{(3)}}{\partial W_{31}^{(2)}} \end{bmatrix} = \begin{bmatrix} a_{11}^{(2)} \\ a_{12}^{(2)} \\ a_{13}^{(2)} \end{bmatrix}$$

10. Now, write  $\frac{\partial z^{(3)}}{\partial W^{(2)}}$  in terms of the vector  $a^{(2)}$ :

$$\frac{\partial z^{(3)}}{\partial W^{(2)}} = (a^{(2)})^T$$

This is just like:  
 $f(x) = ax \Rightarrow f'(x) = a$   
 (We're assuming a is constant, which is true here b/c our training data is fixed).

$$\hat{y} = f(z^{(3)}) = 0.92$$

11.  $W^{(1)} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad W^{(2)} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$   $y = 0.75$   
 $X = [0.3 \ 1.0]$   $\frac{\partial J}{\partial W^{(2)}} = ?$

$J = \frac{1}{2} (0.75 - 0.92)^2$   
 $J = 0.015$  ← Don't need this by compute  $\partial J / \partial$

FORWARD PASS FIRST:

$$\begin{aligned} a^{(1)} &= [0.3 \ 1.0] \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} & a^{(2)} &= f([1.3, 1.3, 1.3]) \\ &= [1.3 \ 1.3 \ 1.3] & &= [0.74, 0.74, 0.74] \\ && z^{(3)} &= a^{(2)} W^{(2)} = 2.37 \end{aligned}$$

$$\frac{\partial J}{\partial W^{(2)}} = (a^{(2)})^T = \begin{bmatrix} 0.74 \\ 0.74 \\ 0.74 \end{bmatrix}$$

12. Next, let's deal with the  $\text{numExamples} > 1$  case. Back in question 4 we temporarily took away the summation, we'll figure out how to re-introduce it now. To get started, write out the full matrix equation for (3), using  $\text{numExamples} = 3$ , and  $\text{inputLayerSize} = 2$ ,  $\text{hiddenLayerSize} = 3$ , and  $\text{outputLayerSize} = 1$ .

$$\begin{bmatrix} z_{11}^{(3)} \\ z_{21}^{(3)} \\ z_{31}^{(3)} \end{bmatrix} = \begin{bmatrix} a_{11}^{(2)} & a_{12}^{(2)} & a_{13}^{(2)} \\ a_{21}^{(2)} & a_{22}^{(2)} & a_{23}^{(2)} \\ a_{31}^{(2)} & a_{32}^{(2)} & a_{33}^{(2)} \end{bmatrix} \begin{bmatrix} W_{11}^{(2)} \\ W_{21}^{(2)} \\ W_{31}^{(2)} \end{bmatrix}$$

What do the rows and columns of your "a" matrix represent?

Rows = examples  
cols = hidden units/neurons.

13. Now that we've let  $\text{numExamples}=3$ , what is the dimension of  $\delta^{(3)}$ ?

$3 \times 1$

$$\begin{bmatrix} \delta_{11}^{(3)} \\ \delta_{21}^{(3)} \\ \delta_{31}^{(3)} \end{bmatrix} = \begin{bmatrix} f(y_{11} - \hat{y}_{11}) f'(z_{11}^{(3)}) \\ f(y_{21} - \hat{y}_{21}) f'(z_{21}^{(3)}) \\ f(y_{31} - \hat{y}_{31}) f'(z_{31}^{(3)}) \end{bmatrix}$$

14. Almost there! Now, sum across our examples in terms of the individual elements of  $\delta^{(3)}$  and  $a^{(2)}$ :

SUMMING ACROSS EXAMPLES

~~Summing across examples~~

$$\frac{\partial J}{\partial W^{(2)}} = \left( \begin{bmatrix} \frac{\partial J}{\partial W_{11}^{(2)}} \\ \frac{\partial J}{\partial W_{21}^{(2)}} \\ \frac{\partial J}{\partial W_{31}^{(2)}} \end{bmatrix} \right) = \begin{bmatrix} a_{11}^{(2)} \delta_{11}^{(3)} + a_{21}^{(2)} \delta_{21}^{(3)} + a_{31}^{(2)} \delta_{31}^{(3)} \\ a_{12}^{(2)} \delta_{11}^{(3)} + a_{22}^{(2)} \delta_{21}^{(3)} + a_{32}^{(2)} \delta_{31}^{(3)} \\ a_{13}^{(2)} \delta_{11}^{(3)} + a_{23}^{(2)} \delta_{21}^{(3)} + a_{33}^{(2)} \delta_{31}^{(3)} \end{bmatrix}$$

15. Now express the above operation in terms of the matrix  $a^{(2)}$  and the vector  $\delta^{(3)}$ .

$$\frac{\partial J}{\partial W^{(2)}} = (a^{(2)})^T \delta^{(3)} = -(a^{(2)})^T (y - \hat{y}) \frac{\partial z^{(3)}}{\partial W^{(2)}}$$