

BACKPROPAGATION

Gabriella Miles

bristol.ac.uk



Recap

Some kind
of input

Weights (or
parameters)

$$f(x, W)$$

Loss

Loss
Function, L

56	89
10	43

$$\begin{bmatrix} 56 & 89 \\ 10 & 43 \end{bmatrix}$$

$$Wx$$

$$\begin{bmatrix} 145.8 \\ -259.6 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$



$$\begin{bmatrix} 0.2 & 0.4 & 1.3 & 2.0 \\ 0 & -2.3 & 0.1 & -1.3 \end{bmatrix} \begin{bmatrix} 56 \\ 89 \\ 10 \\ 43 \end{bmatrix} = \begin{bmatrix} 145.8 \\ -259.6 \end{bmatrix}$$





$\nabla_w L$

What next?

- Let's come up with a way of efficiently finding the parameters that minimise the loss function.

Backpropagation

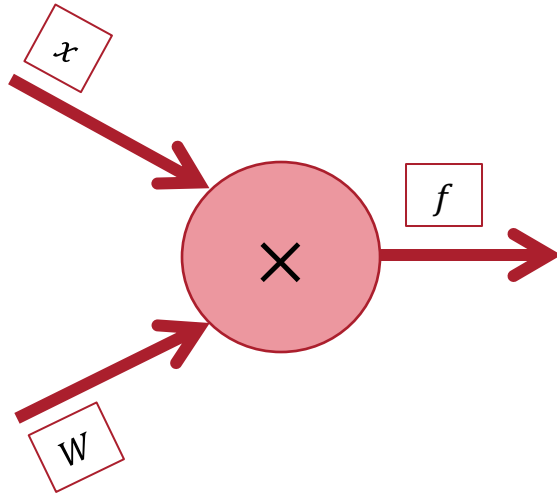
Objectives

Understand:

- What backpropagation *is*
- How it *works*
- Why it's *useful* to us

Computational Graphs

$$f(x, W) = Wx$$



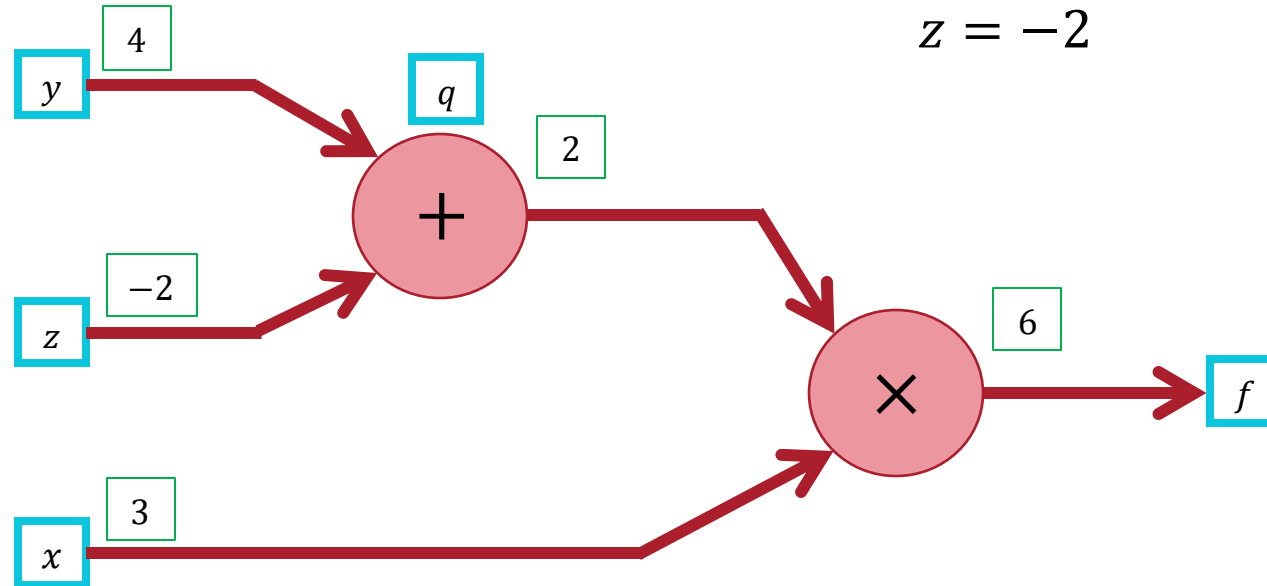
$$f(x, y, z) = x(y + z)$$

Hint: you will need 2 nodes!

The Forward Pass

$$f(x, y, z) = x(y + z)$$

$$\begin{aligned}x &= 3 \\y &= 4 \\z &= -2\end{aligned}$$

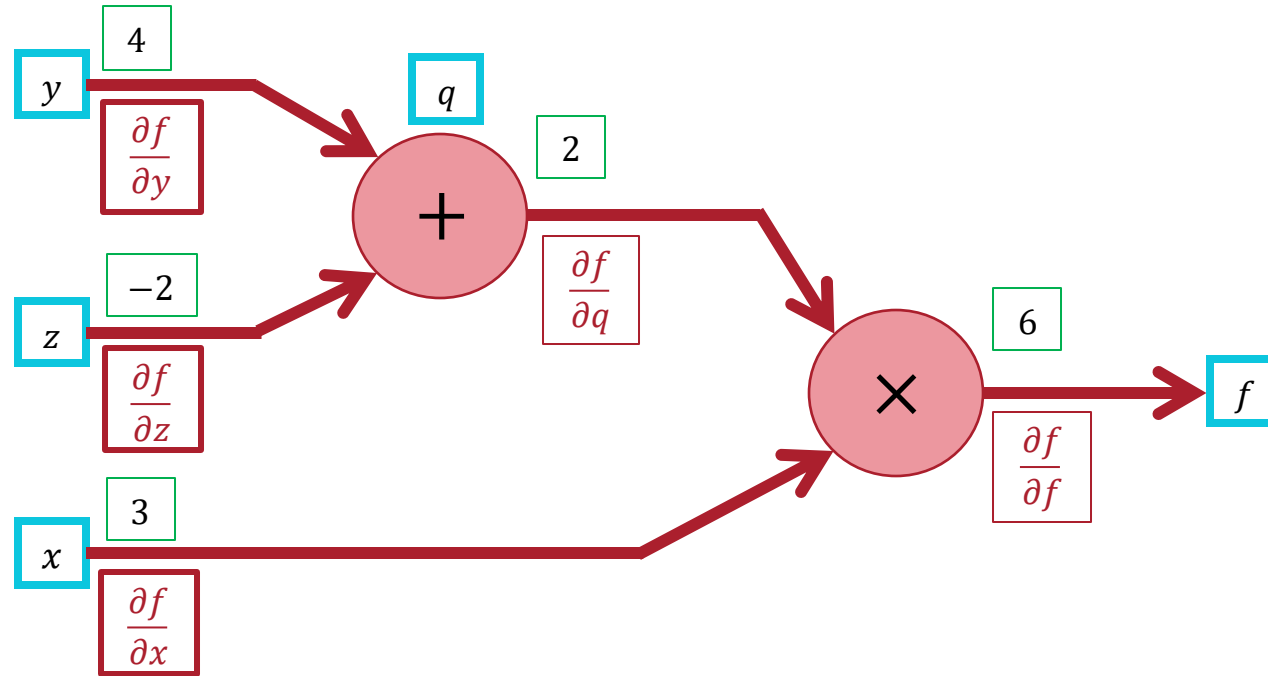


Check out the Computational Graph Worksheet
if you'd like some extra practice!



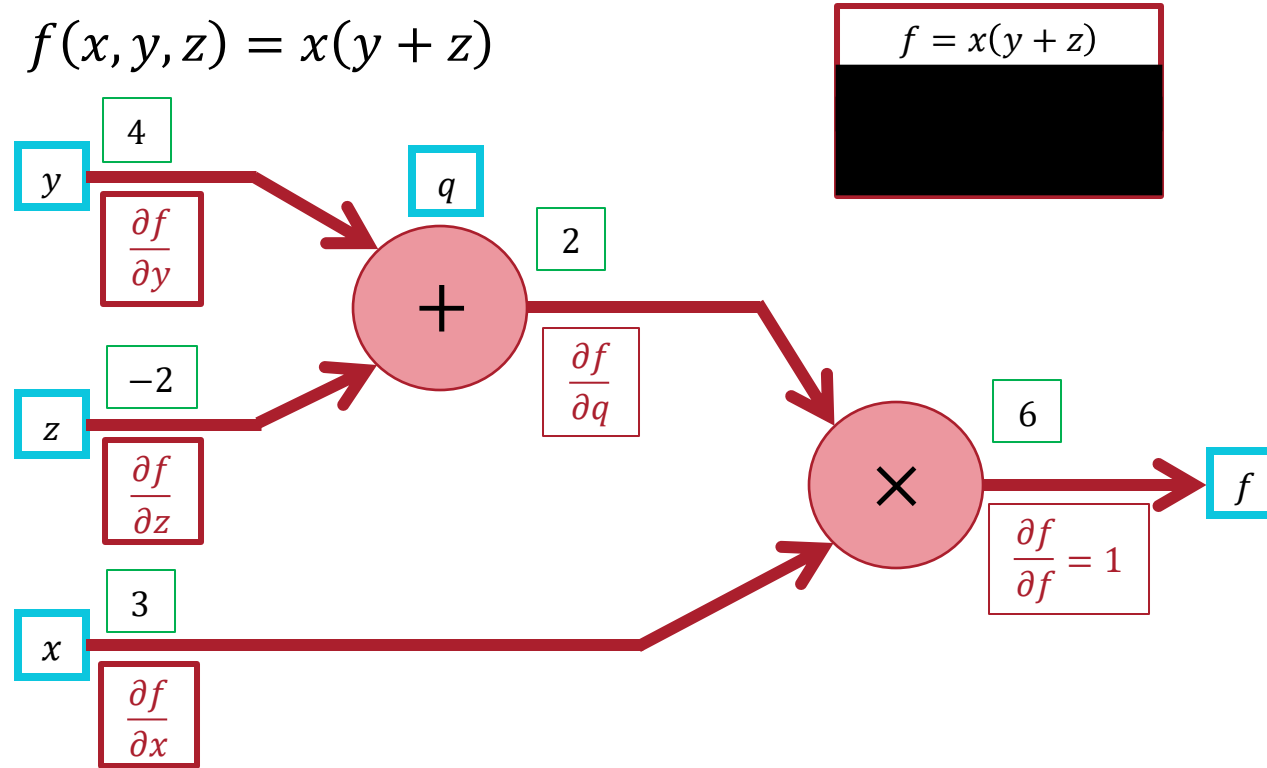
Backpropagation

$$f(x, y, z) = x(y + z)$$

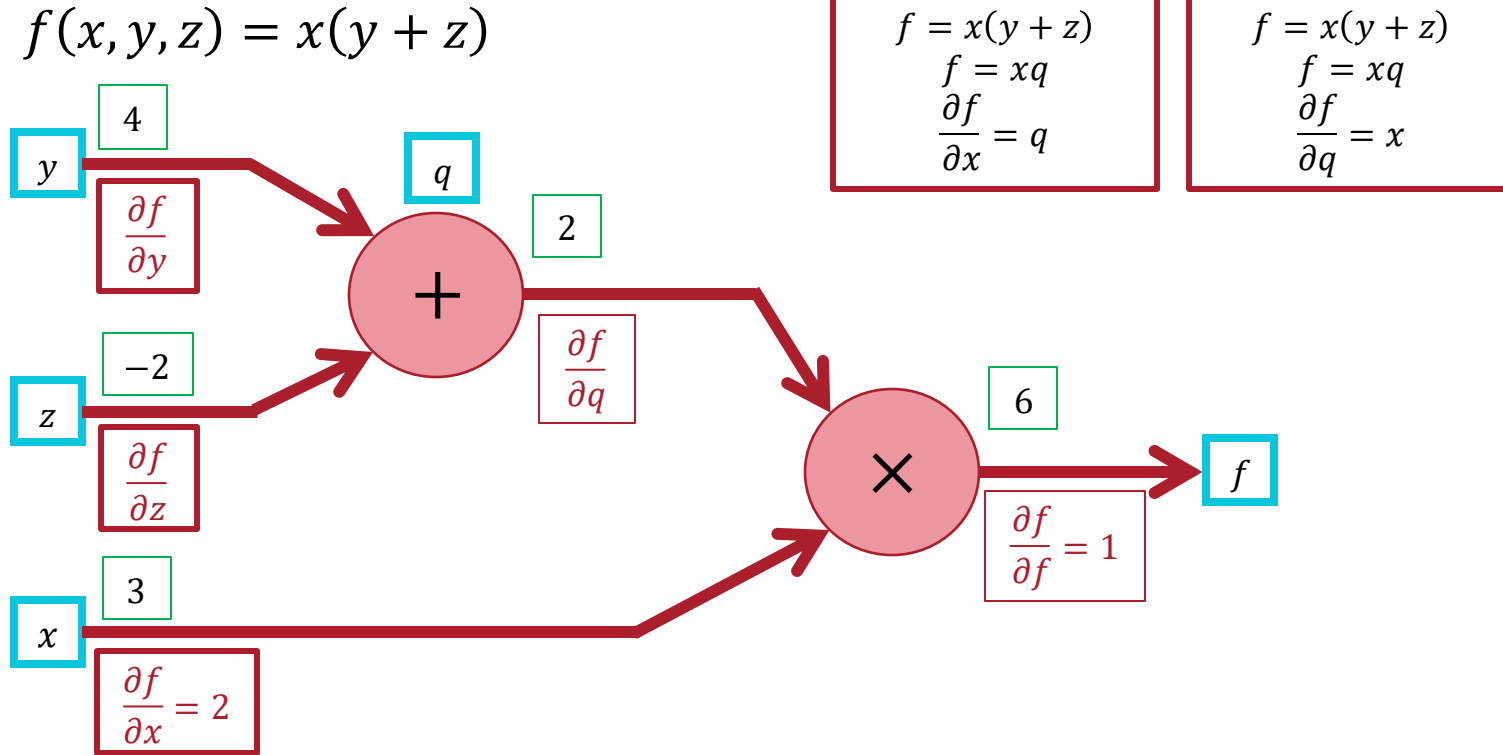


Backpropagation

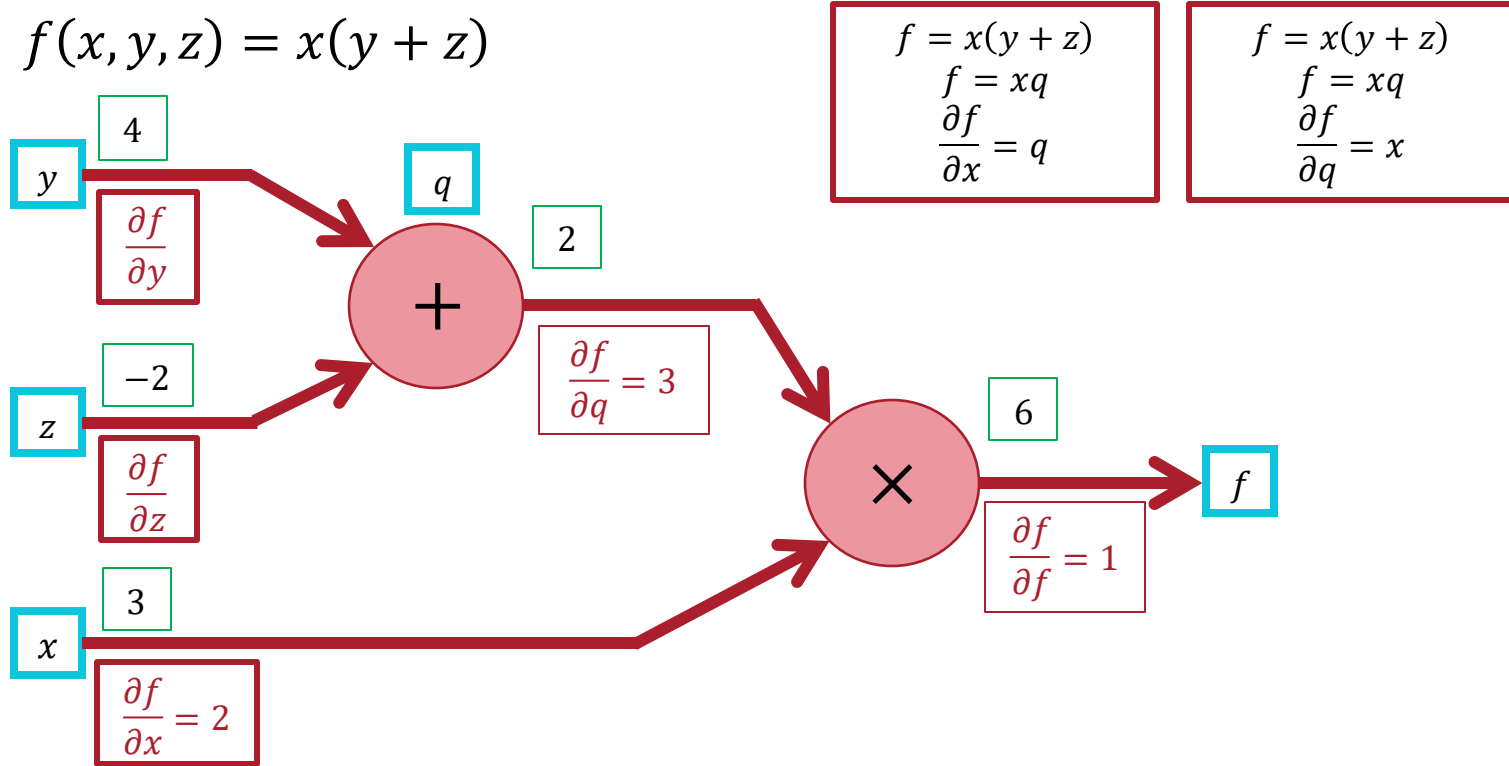
$$f(x, y, z) = x(y + z)$$



Backpropagation

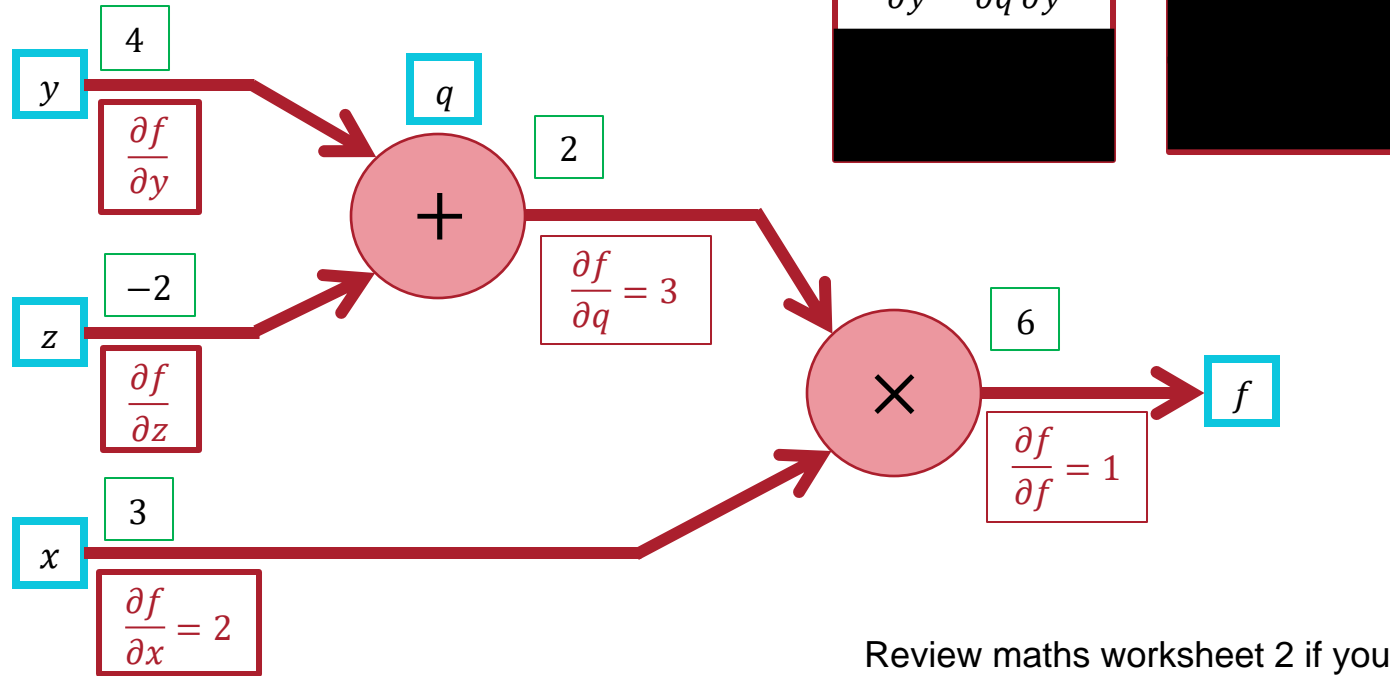


Backpropagation



Backpropagation

$$f(x, y, z) = x(y + z)$$

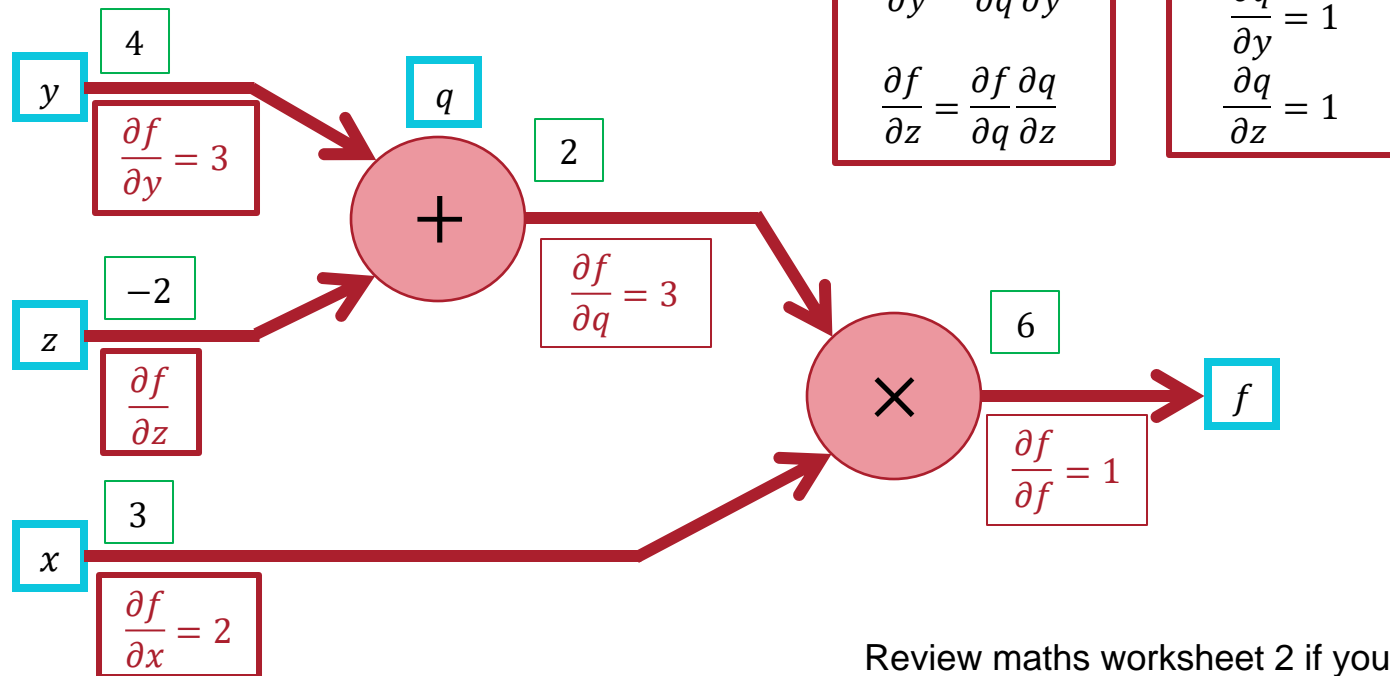


Review maths worksheet 2 if you are feeling rusty on the chain rule!



Backpropagation

$$f(x, y, z) = x(y + z)$$

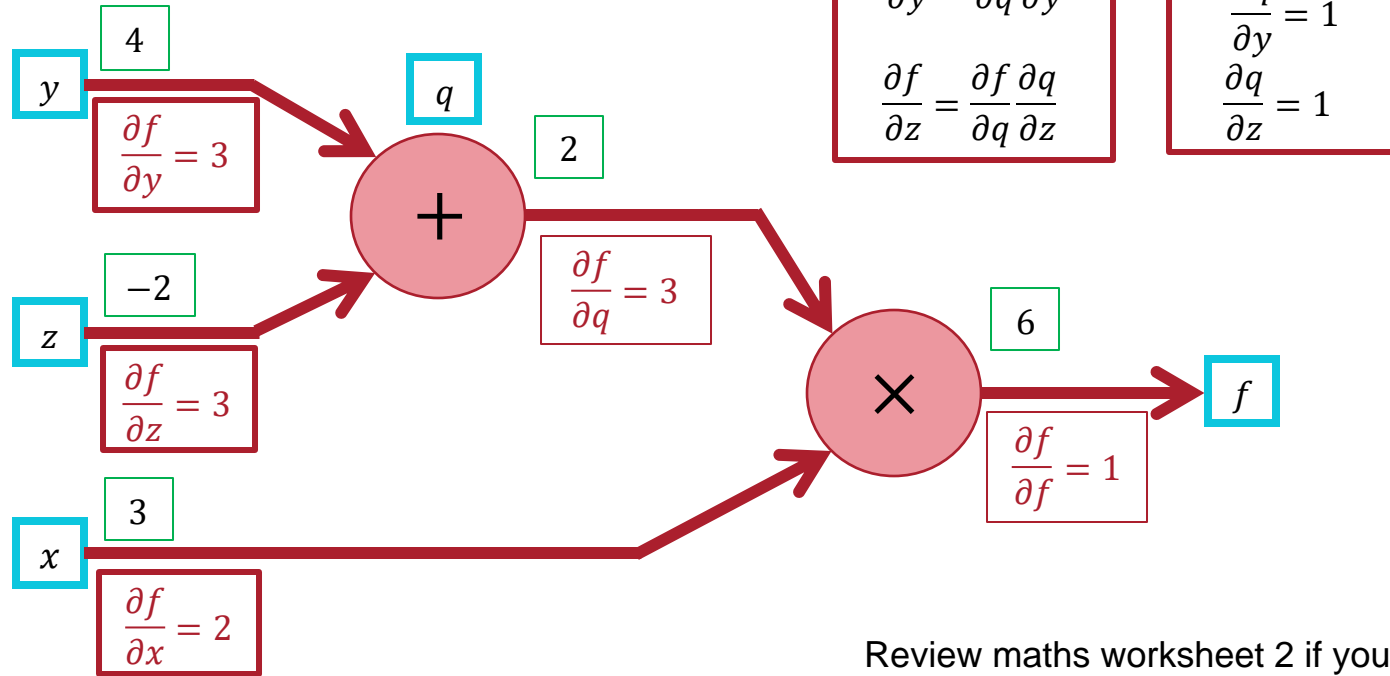


Review maths worksheet 2 if you are feeling rusty on the chain rule!



Backpropagation

$$f(x, y, z) = x(y + z)$$



Review maths worksheet 2 if you are feeling rusty on the chain rule!



Your turn!

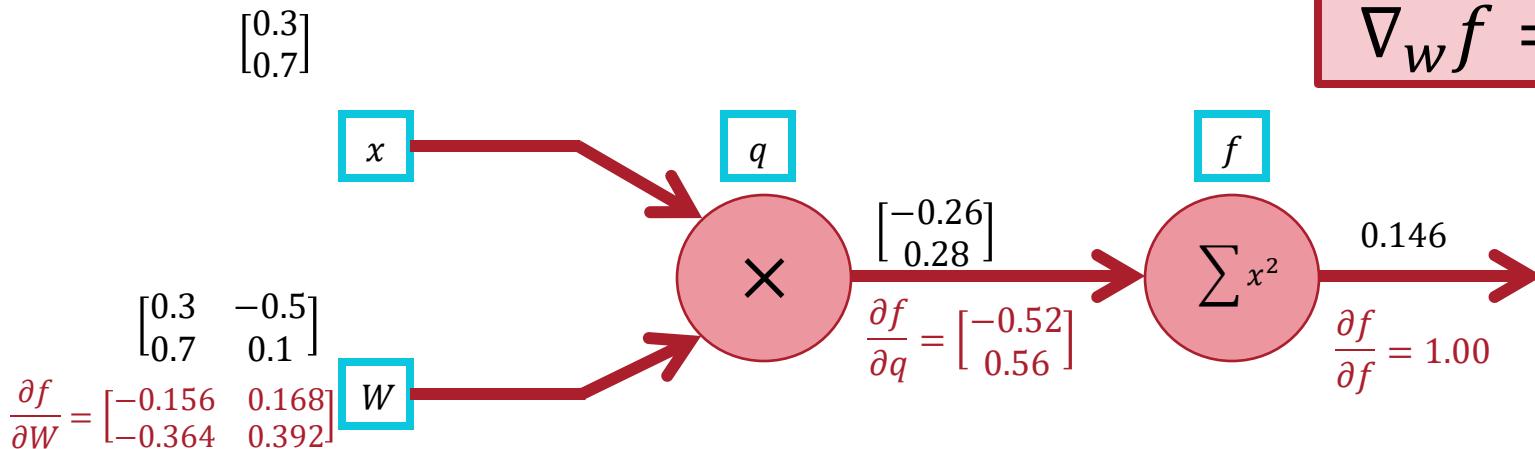
- Present your problem and answer to the class.

Problem Number	Groups
2	1, 4, 7, 10, 13
3	2, 5, 8, 11, 14
4	3, 6, 9, 12, 15

Vectorised...

$$f(x, W) = \sum_{i=1}^n (W \cdot x)_i^2$$

$$\nabla_w f = 2q \cdot x^T$$



$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \dots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \dots + W_{n,n}x_n \end{pmatrix} \longrightarrow \frac{\partial q_k}{\partial W_{i,j}} = x_j$$

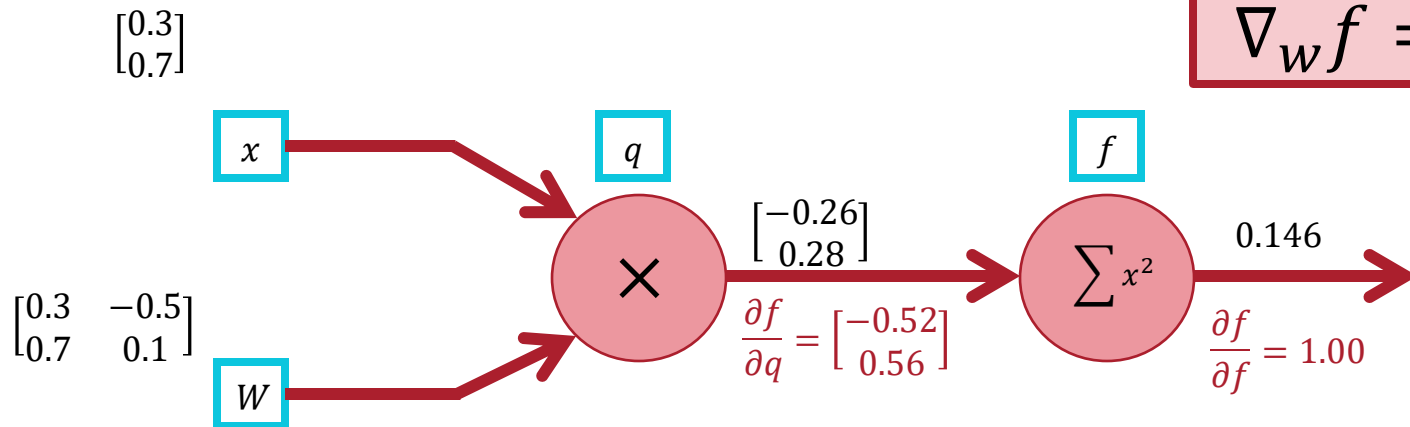
$$f(q) = \sum_{i=1}^n (q)_i^2 = q_1^2 + \dots + q_n^2 \longrightarrow \frac{\partial f}{\partial q_i} = 2q_i$$

$$\frac{\partial f}{\partial W_{i,j}} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial W_{i,j}} = \sum_k (2q_k) x_j = 2q_i x_j$$

Vectorised...

$$f(x, W) = \sum_{i=1}^n (W \cdot x)_i^2$$

$$\nabla_w f = 2q \cdot x^T$$



$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \dots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \dots + W_{n,n}x_n \end{pmatrix} \longrightarrow \frac{\partial q_k}{\partial W_{i,j}} = x_j$$

$$f(q) = \sum_{i=1}^n (q)_i^2 = q_1^2 + \dots + q_n^2 \longrightarrow \frac{\partial f}{\partial q_i} = 2q_i \qquad \frac{\partial f}{\partial W_{i,j}} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial W_{i,j}} = \sum_k (2q_k) x_j = 2q_i x_j$$

Vectorised...

$$f(x, W) = \sum_{i=1}^n (W \cdot x)_i^2$$

$$\nabla_x f = W^T \cdot 2q$$

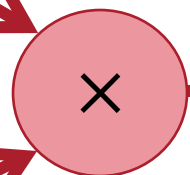
$$\begin{bmatrix} 0.3 \\ 0.7 \end{bmatrix}$$

$$\frac{\partial f}{\partial x} = \begin{bmatrix} 0.236 \\ 0.316 \end{bmatrix}$$

x

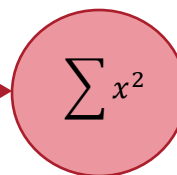
q

f



$$\begin{bmatrix} -0.26 \\ 0.28 \end{bmatrix}$$

$$\frac{\partial f}{\partial q} = \begin{bmatrix} -0.52 \\ 0.56 \end{bmatrix}$$



$$0.146$$

$$\frac{\partial f}{\partial f} = 1.00$$

$$\begin{bmatrix} 0.3 & -0.5 \\ 0.7 & 0.1 \end{bmatrix}$$

$$\frac{\partial f}{\partial W} = \begin{bmatrix} -0.156 & 0.168 \\ -0.364 & 0.392 \end{bmatrix}$$

W

$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \dots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \dots + W_{n,n}x_n \end{pmatrix}$$

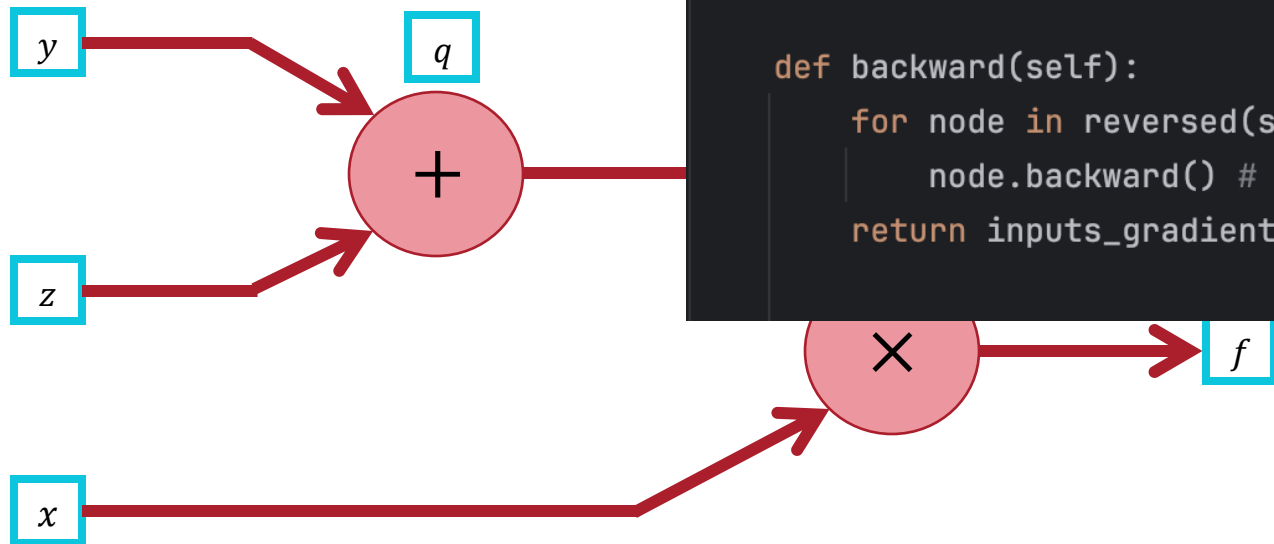
$$\frac{\partial q_k}{\partial x_i} = W_{k,i}$$

$$f(q) = \sum_{i=1}^n (q)_i^2 = q_1^2 + \dots + q_n^2 \rightarrow \frac{\partial f}{\partial q_i} = 2q_i$$

$$\frac{\partial f}{\partial x_i} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial x_i} = \sum_k (2q_k) W_{k,i}$$

Coding

$$f(x, y, z) = x(y + z)$$



```
class ComputationalGraph:
```

```
    def forward(self, inputs):
```

```
        for node in self.graph.nodes():
```

```
            node.forward()
```

```
        return loss
```

```
    def backward(self):
```

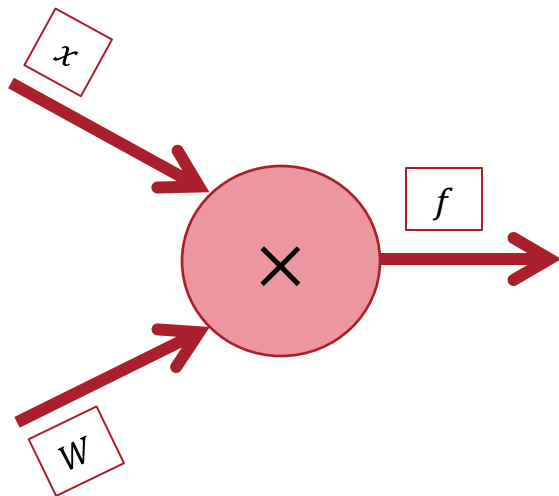
```
        for node in reversed(self.graph.nodes()):
```

```
            node.backward() # backpropagation
```

```
        return inputs_gradients
```

Coding

$$f(x, W) = Wx$$



```
class MultiplyGate(object):  
    def forward(self, x, W):  
        self.x = x  
        self.W = W  
        output = self.x * self.W  
        return output  
  
    def backward(self, dz):  
        # df/dx  
        dx = self.W * dz # [dz/dx * df/dz]  
  
        # df/dW  
        dW = self.x * dx # [dx/dW * df/dz]  
  
        return [dx, dy]
```


Application

- Optimisation

$$w_x^* = w_x - \alpha \left(\frac{\partial L}{\partial w_x} \right)$$

56	89
10	43



Loss
Function, L

Loss

w_x

$\begin{bmatrix} 145.8 \\ -259.6 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

bristol.ac.uk



Backpropagation: where did it come from?

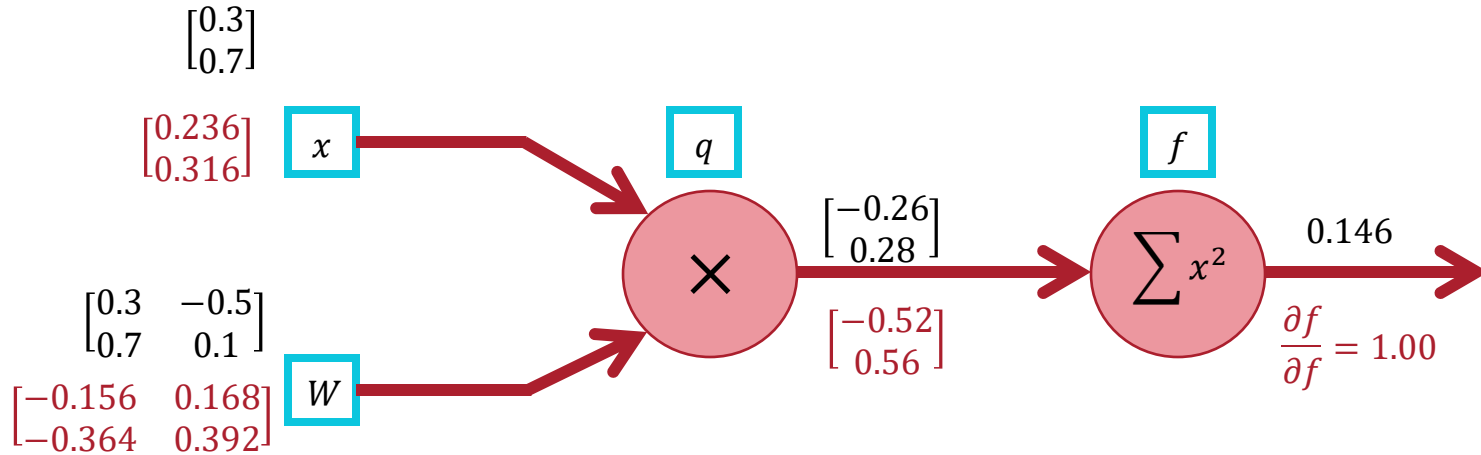
Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323(6088), 533-536.

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.

Recommended Reading List

- Deep Learning, Section 6.5 Back-Propagation and Other Differentiation Algorithms (see Blackboard for pdf)
- Hinton, G. (2022). The forward-forward algorithm: Some preliminary investigations. *arXiv preprint arXiv:2212.13345*.
- Atilim Gunes Baydin, Barak A. Pearlmutter, Alexey Andreyevich Radul, Jeffrey Mark Siskind. Automatic differentiation in machine learning: a survey. *The Journal of Machine Learning Research*, 18(153):1--43, 2018
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *nature*, 323(6088), 533-536.

Drawbacks



Drawbacks



Vanishing gradients: this becomes a real problem in deeper networks (i.e. those with more layers)



What if we can't differentiate the function?



Not really biologically plausible

Hinton, G. (2022). The forward-forward algorithm: Some preliminary investigations. arXiv preprint [arXiv:2212.13345](https://arxiv.org/abs/2212.13345).

A Look Back At Our Objectives

Understand:

- What backpropagation *is*
 - *An algorithm that allows us (in combination with gradient descent) to efficiently train neural networks by minimising the error*
- How it *works*
 - *Compute loss, backpropagate layer by layer (then use gradient descent to update the weights)*
- Why it's *useful* to us
 - *Allows us to efficiently update the weights of the network*

BEFORE YOU GO

Please could you indicate on the menti poll how well you feel you understand backpropagation on a scale of 1-5.



Scan the QR code to book a drop-in if you have struggled to understand this material.



ANY QUESTIONS?

bristol.ac.uk

