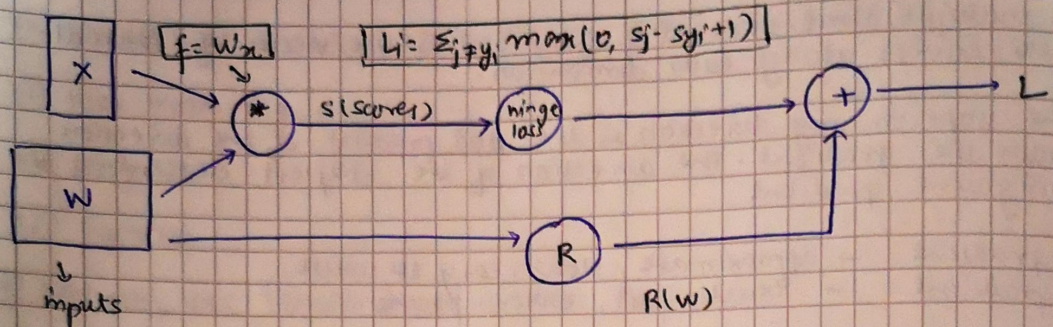


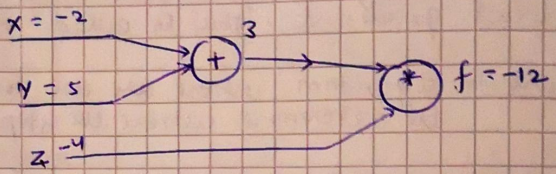
Backpropagation and Neural Networks

① Calculating analytic gradient with computation graph.



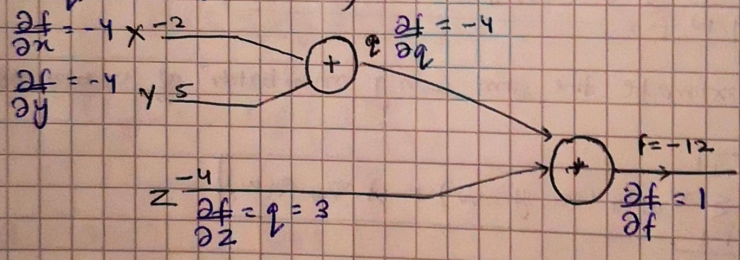
② Backpropagation.

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



$q = x + y$ $\frac{\partial q}{\partial x} = 1$ $\frac{\partial q}{\partial y} = 1$
 $f = qz$ $\frac{\partial f}{\partial q} = z$ $\frac{\partial f}{\partial z} = q$

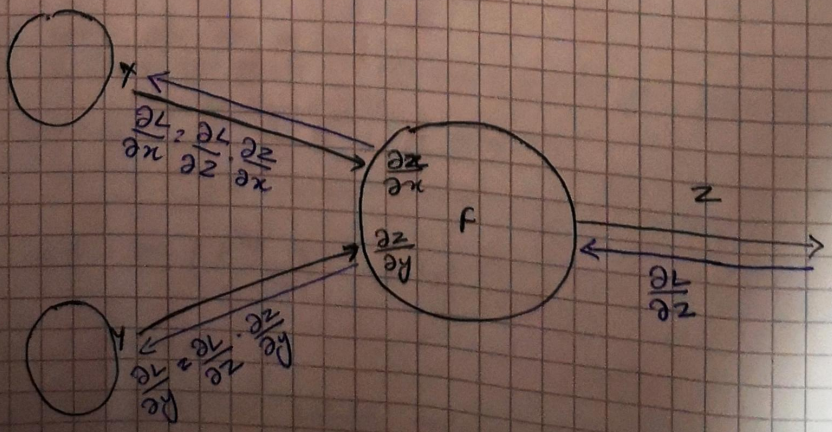
we want $\frac{\partial f}{\partial x}$, $\frac{\partial f}{\partial y}$ and $\frac{\partial f}{\partial z}$



using chainrule for calculating $\frac{\partial f}{\partial x}$ and $\frac{\partial f}{\partial y}$

$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \cdot \frac{\partial q}{\partial x} = -4 \times 1 = -4$

$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \cdot \frac{\partial q}{\partial y} = -4 \times 1 = -4$



② Patterns in Backward flow.

→ Just as an add gate acts as a "gradient distributor" i.e. distributes the upstream gradient equally

→ A max gate acts as a "gradient router" meaning one of the nodes will get the full value of upstream gradient and the other will not receive anything. We can use this to route our desired output

→ A multiplier gate acts as a gradient switcher.

③ Vectorized example.

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

\mathbb{R}^n

$\mathbb{R}^{n \times n}$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} W$$

x

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} x$$

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$\begin{bmatrix} 0.44 \\ 0.56 \end{bmatrix}$$

0.116

1

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

$$f(q) = \|q\|^2 = q_1^2 + \dots + q_n^2$$

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

$$\text{Eg. } \begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix} \begin{bmatrix} 0.2 & 0.4 \end{bmatrix}$$

$$= \begin{bmatrix} 0.088 & 0.104 \\ 0.176 & 0.208 \end{bmatrix} = W$$

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

$$\text{Eg. } 2 \begin{bmatrix} 0.1 & -0.3 \\ 0.5 & 0.8 \end{bmatrix} \begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$= \begin{bmatrix} 0.2 & -0.6 \\ 1.0 & 1.6 \end{bmatrix} \begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$= \begin{bmatrix} -0.112 \\ 0.636 \end{bmatrix} = x$$

$$\frac{\partial f}{\partial q_i} = 2q_i$$

$$\nabla_q f = 2q$$

$$\frac{\partial q_k}{\partial W_{i,j}} = 1_{k=i} x_j$$

$$\frac{\partial f}{\partial W_{i,j}} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial W_{i,j}}$$

$$= \sum_k (2q_k) (1_{k=i} x_j) = 2q_i x_j$$

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

$$\frac{\partial f}{\partial x_i} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial x_i}$$

$$= \sum_k 2q_k W_{k,i}$$