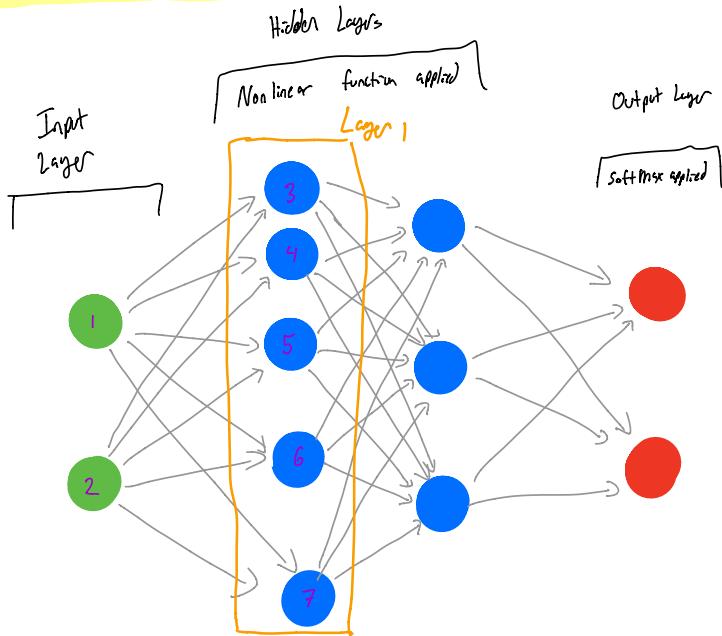


## Neural Networks from Scratch



The more nodes here  
the more complex functions  
we can fit. Careful,  
you can also overfit data.

Forward Propagation → - Calculating subsequent layer values.

Each path has a weight,  $W_{i,j}$ , and each node has a bias,  $b_i$ .

So for example at Layer 1

$$\text{Values} = \text{activation\_function} \cdot \left( \underbrace{\text{weight}}_{\text{matrix of weights}} \cdot \underbrace{\text{inputs}}_{\text{Vector of current node layer}} + \underbrace{\text{bias}}_{\text{Vector of values from previous layer}} \right)$$

Element wise application of activation-function.

$$\begin{bmatrix} \text{Value}_3 \\ \text{Value}_4 \\ \text{Value}_5 \\ \text{Value}_6 \\ \text{Value}_7 \end{bmatrix} = \text{activation function} \left( \begin{bmatrix} W_{1,3} & W_{2,3} \\ W_{1,4} & W_{2,4} \\ W_{1,5} & W_{2,5} \\ W_{1,6} & W_{2,6} \\ W_{1,7} & W_{2,7} \end{bmatrix} \begin{bmatrix} \text{Value}_{-1} \\ \text{Value}_{-2} \end{bmatrix} + \begin{bmatrix} b_3 \\ b_4 \\ b_5 \\ b_6 \\ b_7 \end{bmatrix} \right)$$

Activation functions - Nonlinear function for our hidden layer applied at each node so our hypothesis can take on nonlinear functions. -  $\tanh(x)$ ,  $\frac{e^x}{e^x + 1}$ , etc...

- Softmax function for output layer to transform to probabilities. Turns a vector into a proportional one that sums to 1.  $\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}$

Learning - How do we find weights and biases that minimize error?

A nice choice for a loss function of Softmax outputs:

Cross-entropy loss, AKA negative log likelihood.

$$L(y, \hat{y}) = \frac{-1}{N} \sum_{n \in N} \sum_{i \in C} y_{n,i} \log(\hat{y}_{n,i})$$

↪ prediction value  
 ↪ true value  
 ↪  $C$  classes (output nodes)  
 ↪  $N$  training examples

## Backwards Propagation (Gradient descent)

We use gradient descent to find minimum of loss function with respect to the parameters.

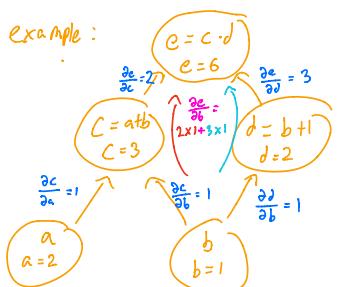
Thus we need the derivatives.  $\frac{\partial L}{\partial b_1}, \frac{\partial L}{\partial b_2}, \dots, \frac{\partial L}{\partial w_{1,3}}, \frac{\partial L}{\partial w_{2,3}}, \dots$

We do this by starting from the output layer and working backwards via the backpropagation algorithm.

Backwards Propagation, formally - reverse-mode differentiation

Derivatives on computational graphs

Example:



$C$  is directly effected by  $a$ , how does it change if we change  $a$ ?

$$\frac{\partial C}{\partial a} = \frac{\partial}{\partial a}(ab) = \frac{\partial a}{\partial a} + \frac{\partial b}{\partial a} = 1 + 0$$

Also need product rule: ( $U$  &  $V$  are independent)

$$\frac{\partial}{\partial U} UV = U \frac{\partial V}{\partial U} + V \frac{\partial U}{\partial V} = 0 + V$$

General Rule - sum over all possible paths, multiplying the derivatives on each edge of a path together.

But just summing over all paths is extremely inefficient with large graphs.

Consider another graph

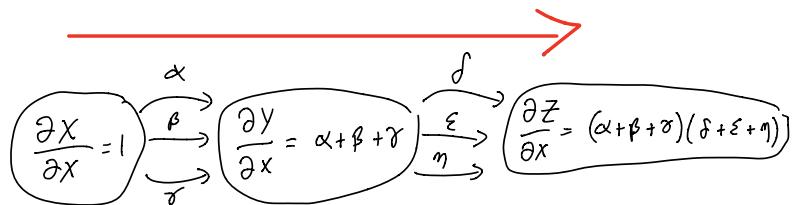


Naive way  $\rightarrow \frac{\partial Z}{\partial x} = \alpha f + \alpha e + \alpha m + \beta f + \beta e + \beta m + \gamma d + \gamma e + \gamma m$   
 how can we speed this up?

But we can just factor...

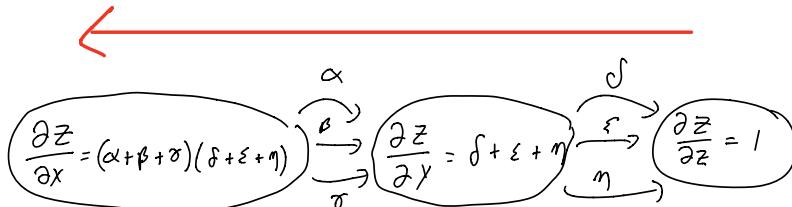
$$(\alpha + \beta + \gamma)(f + e + m)$$

### Forward mode Differentiation



Calculates how one input effects every node  $\left( \frac{\partial}{\partial x} \right)$

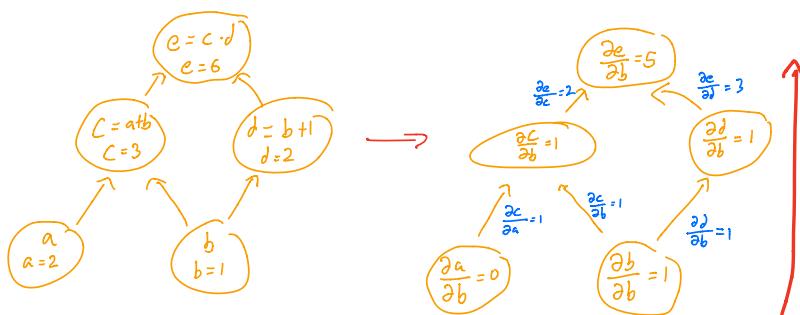
### Reverse mode Differentiation



Calculates how every node effects one output  $\left( \frac{\partial}{\partial} \right)$

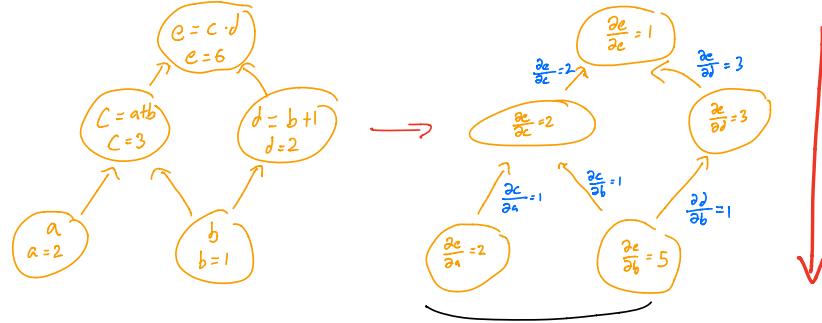
Ok back to the original graph... How is reverse mode computationally efficient?

With forward mode differentiation:



But, then we need to  
 re-traverse the graph for a.  
 What if we have 1000 inputs?  
 That's 1000x traversing the  
 same graph.

With reverse-mode differentiation



We get all the derivatives we need in one graph traversal!!

Forward diff gives derivative of our output with respect to one input.

Reverse diff gives all of them!

Forward diff however gives us all of the derivatives of outputs with respect to our input.

If we have a large amount of outputs forward mode can make more sense.

## Overlook

1. Forward propagation, calculate weights and biases. Use auto diff to keep track of partial derivatives on paths
2. Work your way back to calculate gradient using reverse-mode differentiation.
3. Finally, apply gradient descent updating weights and biases.

$$\hat{y}_{out} = \hat{y}_n - \gamma \nabla L(\hat{y}_n(W, b))$$

$x \in \mathbb{R}$

$b$  is irrational

$\sqrt{a} \in \mathbb{Q}^+$

$x \in \mathbb{R}$

$x$  is rational