

Spelled out intro to neural nets and back propagation
building micrograd - Andrej Karpathy

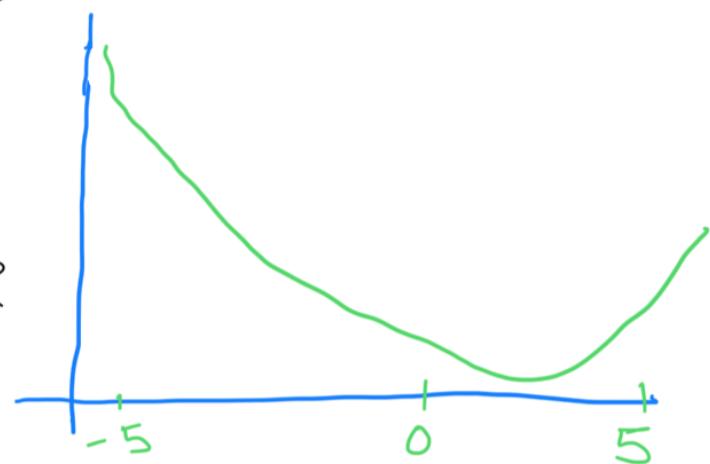
$$x = \text{np.arange}(-5, 5, 0.25)$$

$$y = 3x^2 - 4x + 5$$

what is $\frac{dy}{dx}$ at every point?

usual way \rightarrow find $\frac{dy}{dx}$ using calculus and put in the

value of x .



$$\text{Derivative} = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}$$

$$f(x) = 3x^2 - 4x + 5$$

$$d = \frac{f(x+h) - f(x)}{h}, \quad h = 0.0001, \quad x = 3$$

Now,

consider,

$$d = a \cdot b + c, \quad a = 2.0 \\ b = -3.0 \Rightarrow d = 4 \\ c = 10.0$$

How, to calculate derivative of d wrt a ?

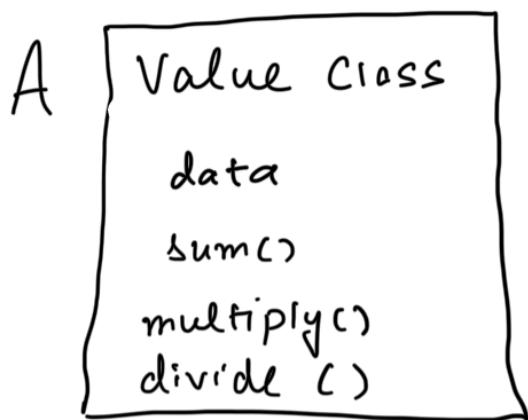
$$a_1 = a, \quad a_2 = a + h$$

$$d_1 = a_1 \cdot b + c, \quad d_2 = a_2 \cdot b + c$$

$$\frac{\nabla d}{\nabla a} = \frac{d_2 - d_1}{h} \Rightarrow \frac{(a+h)b + c - (a \cdot b + c)}{h}$$
$$\Rightarrow \frac{ab + hb + c - ab - c}{h} = \frac{hb}{h} = b$$

Similarly, we do it for other variables.

Wrapping data into objects.



In order to overload `+`, `-`, `*` in python, you have special class of function `--add--(self, other)` where `other` is another object that you want to add with.

Other operations are

Addition	<code>--add--</code>	<code>+</code>
Subtraction	<code>--sub--</code>	<code>-</code>
Multiplication	<code>--mul--</code>	<code>*</code>
Division (float)	<code>--truediv--</code>	<code>/</code>

We also have reverse methods which gets called when the main object is on the right hand side of the operand, just add `r` before the name, `--rsub--`

In place methods, when we don't want to return a new object, just add `i` before the name, `--isub--()`.

If the operation is not defined by special function class,

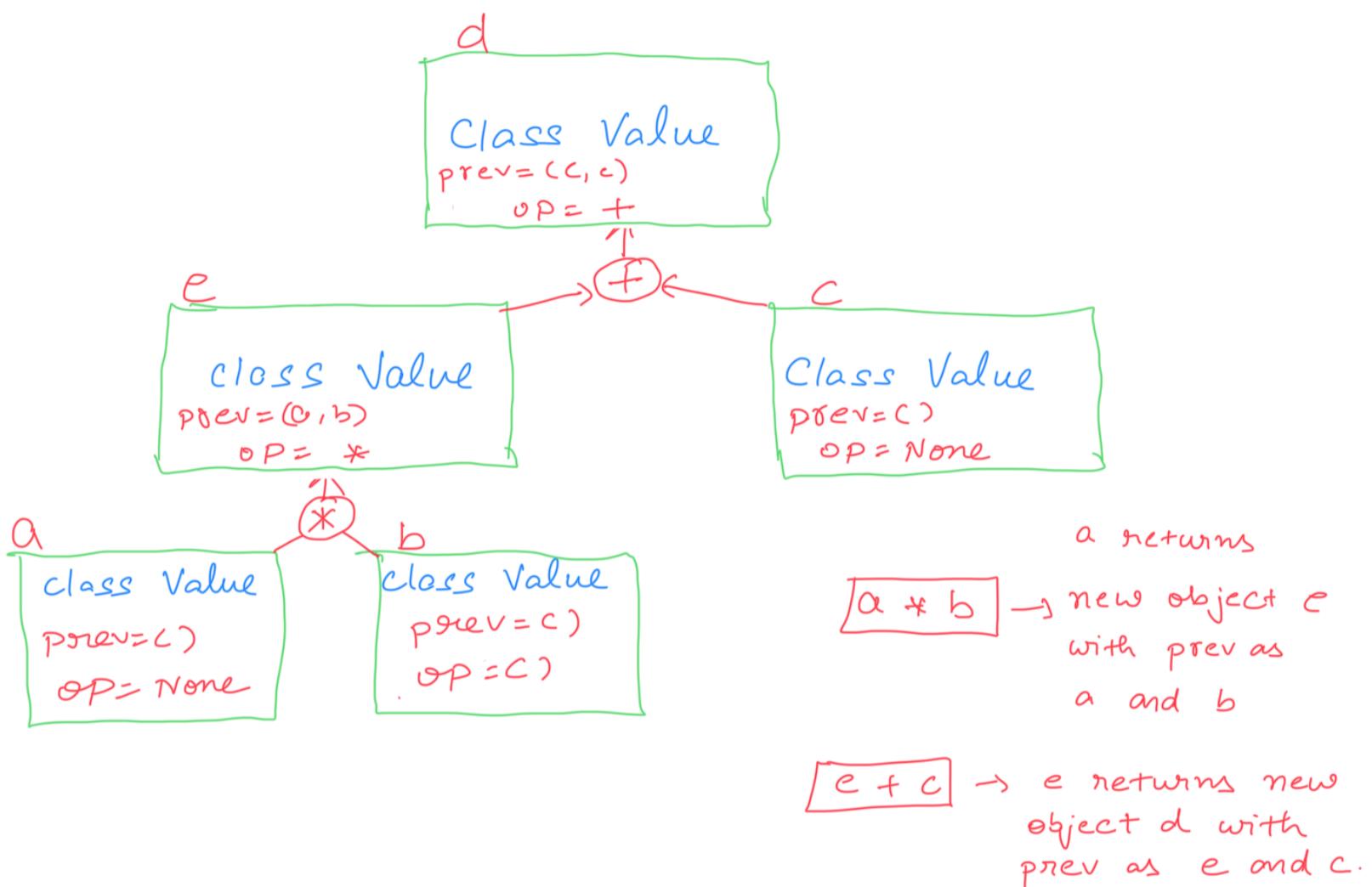
- ① Piggybank on existing operation not used much.
- ② Use functions instead of operation.

Now, that we have object representation of values

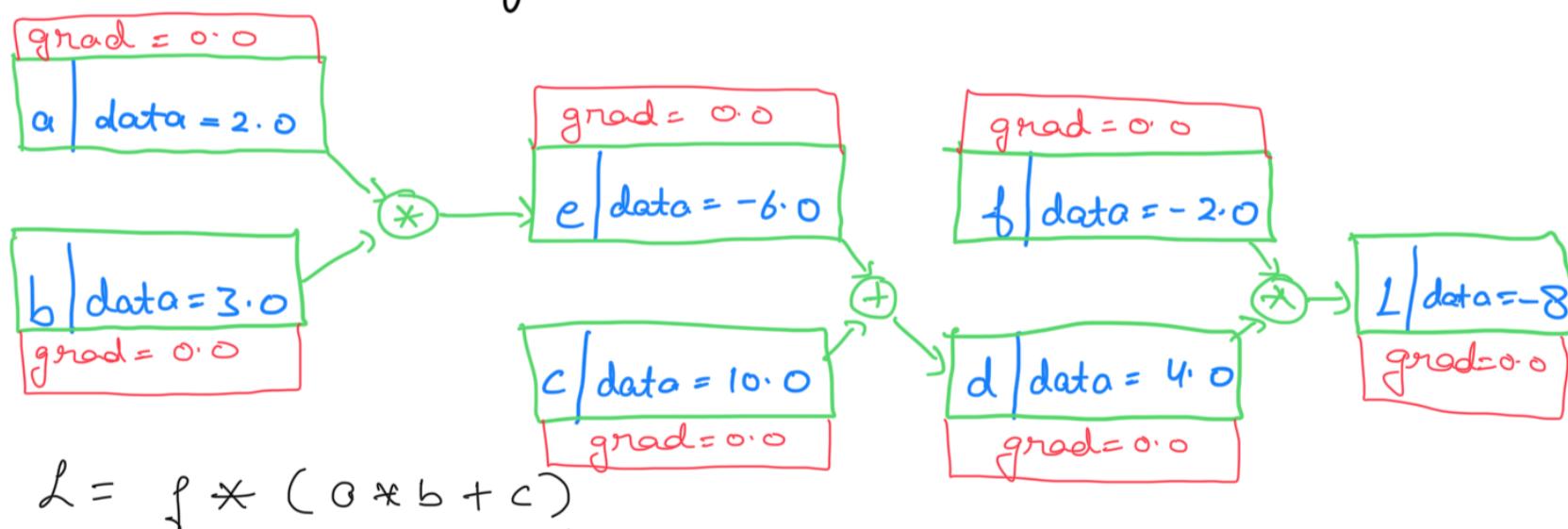
$$d = \text{Value}(a) * \text{Value}(b) + \text{Value}(c) \quad \begin{array}{l} \uparrow \\ \text{Same} \end{array}$$

$$u = (\text{a} \dots \text{mul} \dots (\text{b})) \dots \text{add} \dots (\text{c})$$

Now, we have the object notations of numbers. we now want to create a dependency graph to define a function and what values are dependent on what other values.



Consider a loss function such as:



how do we find the derivative of L with all the variables?

We define a new variable grad for each value which will represent the gradient of output w.r.t that value.

Gradient Chunks -

numerical way = calculate the change in output
with a small change in input.

$$\frac{dL}{da} = \frac{L(a+h) - L(a)}{h}$$

$$= \frac{f * ((a+h)*b+c) - f * (a*b+c)}{h}$$

by this method we can find the derivative of output with each dependent values.

problems →

for this method to implement, we have to calculate a new value $f(x+h)$ for every dependent value, hence will take $O(N^2)$ complexity for n variables as each calculation will take $O(N)$ which will be repeated N times.

This method is not suitable to be used in production but can be used for correctness of our algorithm.

Auto Differentiation

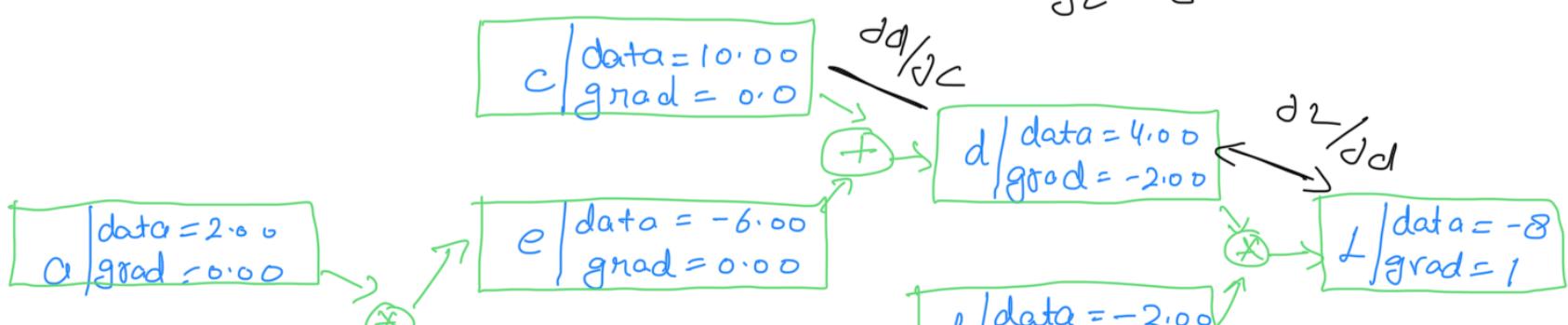
$$L = d * f \quad , \quad \frac{\partial L}{\partial d} = f \quad , \quad \frac{\partial L}{\partial f} = d$$

$$L = f * (a * b + c)$$

how do we calculate $\frac{\partial L}{\partial c}$ without symbolic differentiation and without numerical differentiation?

let's draw the graph again.

$$\frac{\partial L}{\partial c} = \frac{\partial L}{\partial d} * \frac{\partial d}{\partial c}$$



b | data = 3.00
grad = 0.00

f | grad = 4.00

in order to calculate $\frac{dL}{dd}$, we can just calculate $\frac{dL}{da}$ and $\frac{dL}{dc}$ and apply the chain rule.

$$\text{Chain rule} = \frac{dL}{da} = \frac{dL}{db} \times \frac{db}{dc} \times \frac{dc}{dd} \times \frac{dd}{da} \dots$$

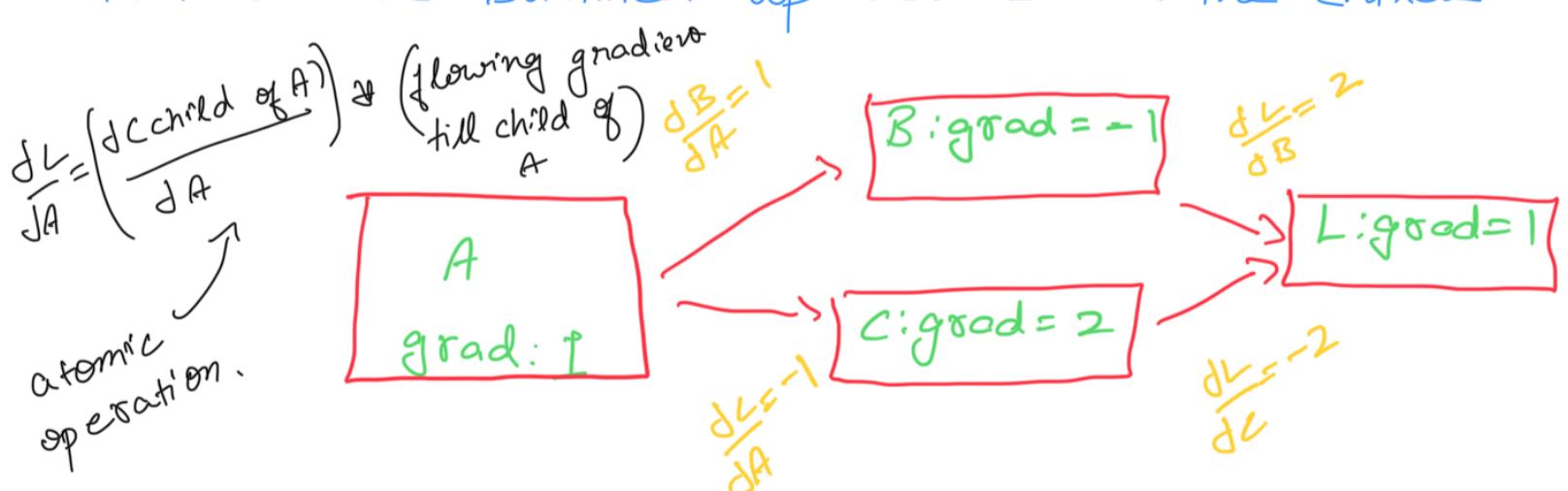
So,

in order to calculate the derivate of output with any variable, we move along the edges to that input and accumulate the gradients to that node via multiplication.

Hence,

- 1) The graph will always be a DAG.
- 2) The edge will store the derivative of Child w.r.t parent.

The value at edge which is a derivative of child wrt parent can be stored in parent class as a grad that can be summed up across all the childs.

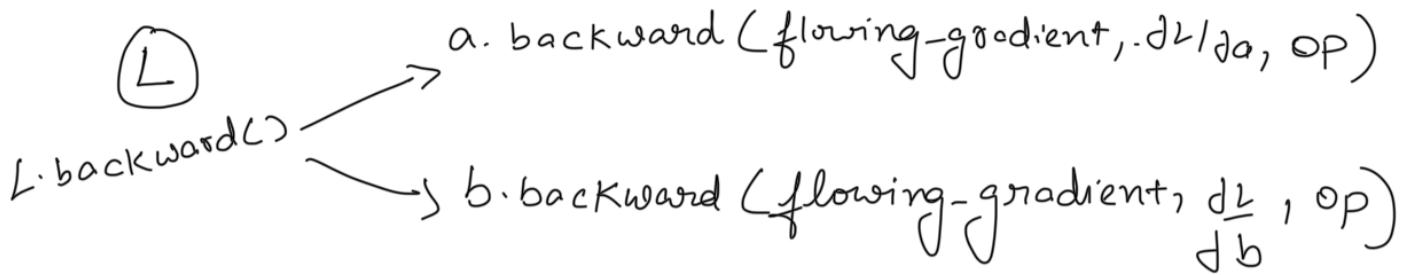


$$B \cdot \text{grad} + L \cdot \text{grad} * \frac{dL}{dB}, C \cdot \text{grad} = L \cdot \text{grad} * \frac{dL}{dC}$$

$$A \cdot \text{grad} + B \cdot \text{grad} * \frac{dB}{dA} + C \cdot \text{grad} * \frac{dC}{dA}$$

thought:-

pseudocode.



a. backward (flowing-gradient, $\frac{\partial \text{child}}{\partial \text{parent}}$, operation);

$$\text{grad} + = \text{flowing gradient} \times \frac{\partial \text{child}}{\partial \text{parent}}$$

$$\text{flowing-gradient} \cdot \dagger = \frac{\partial \text{child}}{\partial \text{parent}}$$

if op = '†':

child1.backward (flowing-gradient, child2.data, '+')

child2.backward (flowing-gradient, child1.data, '†')

Similarly for other atomic operations.

—————

In summary, Auto differentiation is

- ① Splitting the output expression into atomic operations of +, -, ÷, * etc.
- ② Creating a DAG of our function using these atomic expressions.
- ③ Doing a forward pass on our DAG using our atomic expressions.
- ④ Recursively applying chain rule and calculating the gradient of parent w.r.t child node.
- ⑤ after gradient accumulation, changing the value of variables according to gradients.

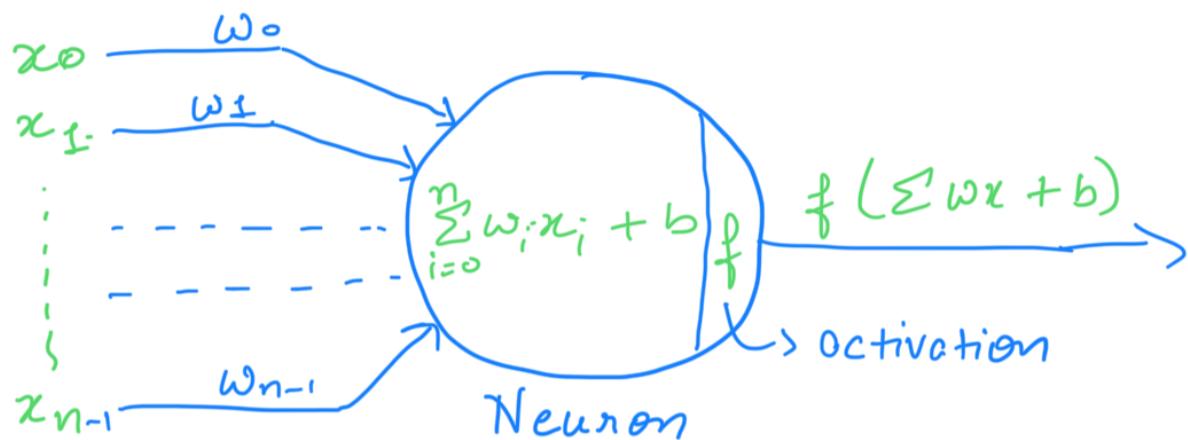
Total complexity = $O(\underline{N})$, N = number of variables.

backpropagation = Recursive application of chain rule backwards in computational graph

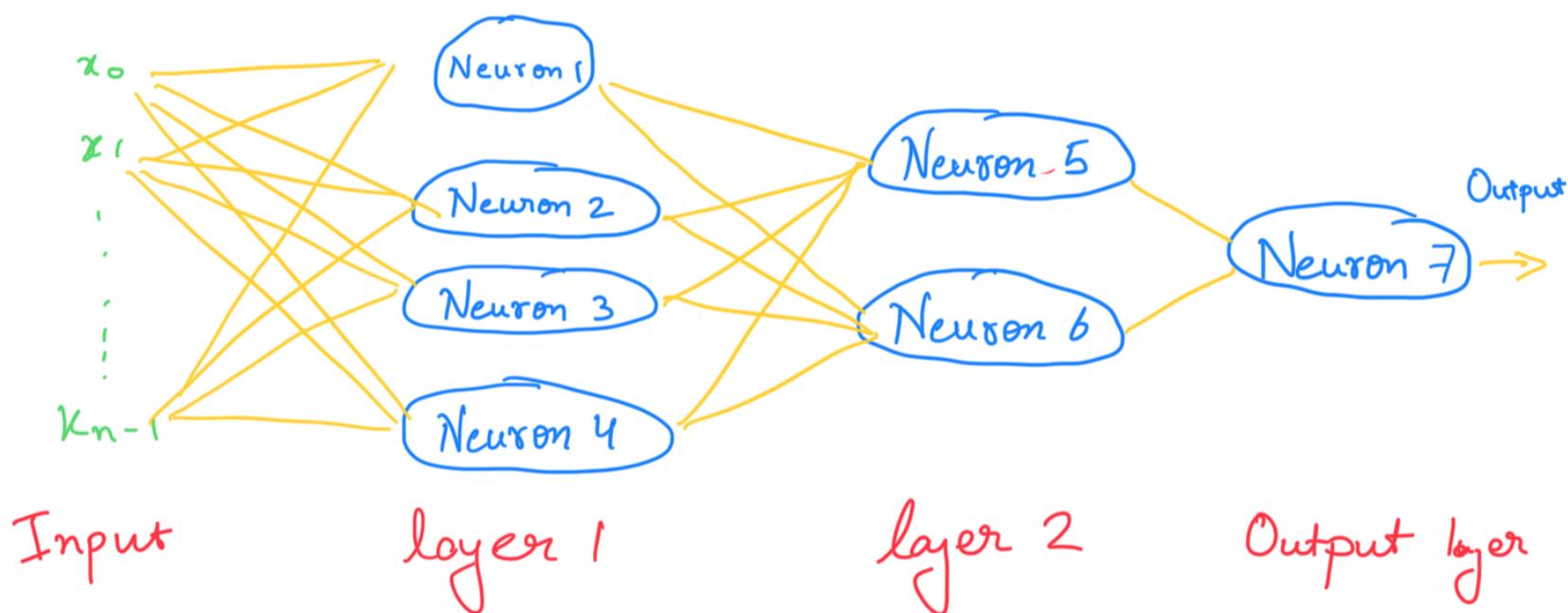
—————

Now, consider a more realistic and complex scenario.

How to design a neuron?



2 Layer Neural Network.



given input $x = [x_0, x_1, x_2, x_3, \dots, x_{n-1}]$

how to prepare a neuron which can do a forward pass as well as backward pass?

$$\text{Neuron } n = \sum_{i=0}^n w_i x_i + b$$

activation functions :

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^{2x} - 1}{e^{2x} + 1}$$

designing a neuron:

$$x_1 = \text{Value}(2.0)$$

$$x_1 w_1 = x_1 * w_1$$

$$x_2 = \text{Value}(0.0)$$

$$x_2 w_2 = x_2 * w_2$$

$$w_1 = \text{Value}(-3.0)$$

$$x_1 w_1 x_2 w_2 =$$

$$O = n \cdot \tanh(L)$$

$$w_2 = \text{Value}(1.0)$$

$$x_1 w_1 + x_2 * w_2$$

$$b = \text{Value}(6.7)$$

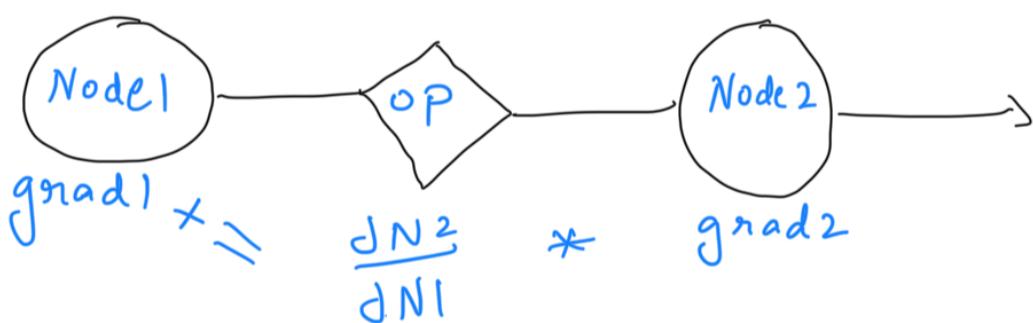
$$n = x_1 w_1 x_2 w_2 + b$$

You don't have to break the local derivative into atomic pieces, like the case of \tanh .

$$\tanh = \frac{e^{2x} - 1}{e^{2x} + 1}$$

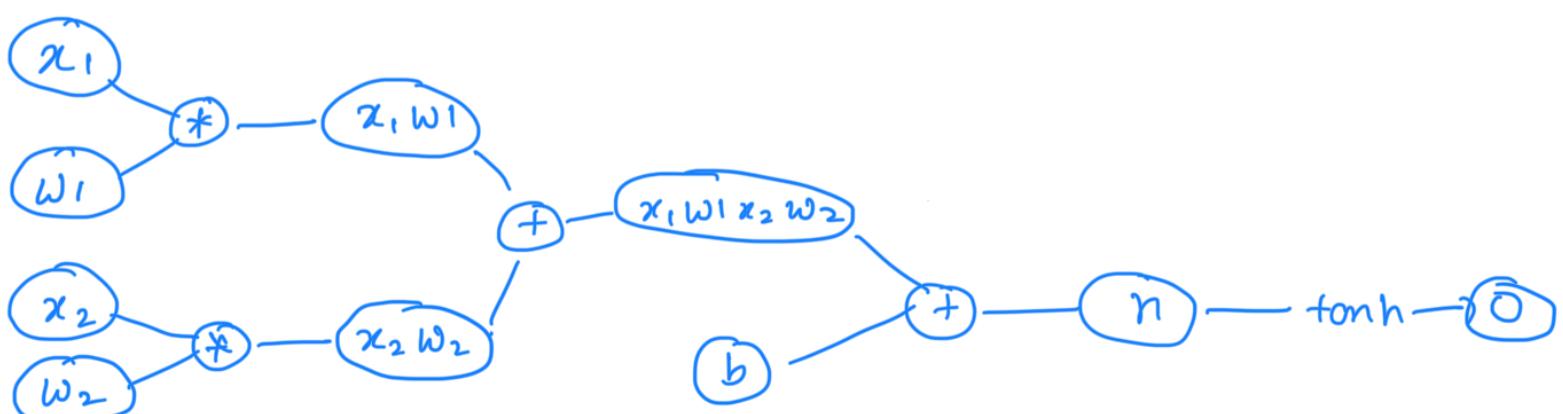
\tanh can be broken down further into atomic operations of $2x$, e^{2x} , $e^{2x} - 1$, $e^{2x} + 1$ but we can choose to not do that and just compute the gradient of \tanh directly as it is commonly used.

general rule of backprop:



$$\frac{\partial \tanh}{\partial x} = 1 - (\tanh)^2$$

$$\frac{\partial \text{sigmoid}}{\partial x} = (\text{sigmoid})(1 - \text{sigmoid})$$



$$x_1 \cdot \text{grad} = x_1 w_1 \cdot \text{grad} * w_1 \cdot \text{data}$$

$$w_1 \cdot \text{grad} = x_1 w_1 \cdot \text{grad} * x_1 \cdot \text{data}$$

$$x_2 \cdot \text{grad} = x_2 w_2 \cdot \text{grad} * w_2 \cdot \text{data}$$

$$w_2 \cdot \text{grad} = x_2 w_2 \cdot \text{grad} * x_2 \cdot \text{data}$$

$$x_1 w_1 \cdot \text{grad} = x_1 w_1 x_2 w_2 \cdot \text{grad} * 1$$

$$x_2 w_2 \cdot \text{grad} = x_1 w_1 x_2 w_2 \cdot \text{grad} * 1$$

$$x_1 w_1 x_2 w_2 \cdot \text{grad} = n \cdot \text{grad} * 1$$

$$b \cdot \text{grad} = n \cdot \text{grad} * 1$$

$$n \cdot \text{grad} = (1 - (o \cdot \text{data})^2) \cdot o \cdot \text{grad} \cdot$$

The NN graph is a DAG so while doing backprop, we apply topological sort.

so that the parents gradients are computed only after child's gradient has been computed.