

* Backpropagation

- Like the dependencies for ^{loss on} earlier parameters w and b were found and used to update them, in the same way the dependency of error/loss on each parameter of the net is evaluated and used to tune the parameters to increase accuracy. for example

$$\Delta w_{11} = \frac{\partial L(\theta)}{\partial w_{11}}$$

unit in the layer 1
unit in the layer 2
corresponding layer

$$w_{11}' = w_{11} - \eta \Delta w_{11}$$

$$\Delta b_{12} = \frac{\partial L(\theta)}{\partial b_{12}}$$

bias involved in layer 1 connected to neuron on layer 2.

- similarly for the final outputs and the bias as well.
- As stated earlier, θ is the vector of the parameters, which earlier were w , and b . the derivative of loss wrt θ , were the vector containing the partial derivatives of loss wrt the specific parameters of the vector.

↓
The same argument can be extended here, wherein the parameters are segmented layer wise and compose a 2D matrix.

↓
but $\Delta \theta$, still remains the same : vector of the partial derivatives of the constituent elements.

- An example of derivatives.
ex 1. e^{x^2}

$$x \rightarrow (s_1) \rightarrow (\exp) \rightarrow (y)$$

$h_1 = x^2$ $h_2 = e^{x^2}$ $y = e^{e^{x^2}}$

$$\frac{dy}{dx} = \frac{dy}{dh_2} \cdot \frac{dh_2}{dh_1} \cdot \frac{dh_1}{dx}$$

(intuition for autodifferentiation)
chain rule

Now, $y = e^{h_2}$

$$\frac{dy}{dh_2} = e^{h_2}$$

$$\Rightarrow h_2 = e^{h_1}$$

$$\frac{dh_2}{dh_1} = e^{h_1}$$

$$\Rightarrow h_1 = x^2$$

$$\frac{dh_1}{dx} = 2x$$

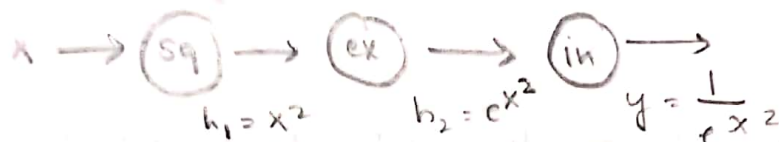
$$\frac{dy}{dx} = \frac{dy}{dh_2} \cdot \frac{dh_2}{dh_1} \cdot \frac{dh_1}{dx}$$

$$= e^{h_2} \cdot e^{h_1} \cdot 2x$$

↓ substituting values of h_1 and h_2

$$= e^{e^{x^2}} \cdot e^{x^2} \cdot 2x$$

ex 2. $\frac{1}{e^{x^2}}$



$$\frac{dy}{dx} = \frac{dy}{dh_2} \cdot \frac{dh_2}{dh_1} \cdot \frac{dh_1}{dx}$$

$$y = \frac{1}{h_2} \Rightarrow \frac{dy}{dh_2} = -\frac{1}{(h_2)^2}$$

$$h_2 = e^{h_1} \Rightarrow \frac{dh_2}{dh_1} = e^{h_1}$$

$$h_1 = x^2 \Rightarrow \frac{dh_1}{dx} = 2x$$

↓ substituting values

$$\frac{dy}{dh_2} = -\frac{1}{(e^{x^2})^2}$$

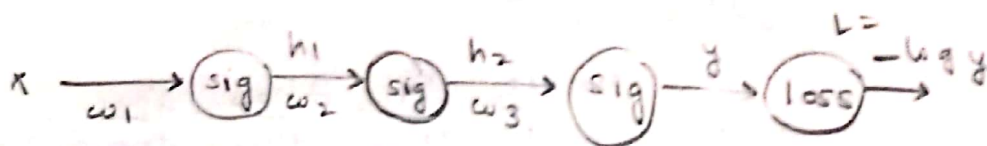
$$\frac{dh_2}{dh_1} = e^{x^2}$$

$$\frac{dh_1}{dx} = 2x$$

↓

$$\frac{dy}{dx} = -\frac{2x}{e^{x^2}}$$

-) The chain rule based computation can be visualised as the reverse of the computation done in forward direction.



$$L = -\log y$$

~~dh/dx~~

$$\frac{dL}{dx} = \frac{dL}{dy} \cdot \frac{dy}{dh_2} \cdot \frac{dh_2}{dh_1} \cdot \frac{dh_1}{dx}$$

$$\frac{dL}{dw_1} = \frac{dL}{dy} \cdot \frac{dy}{dh_2} \cdot \frac{dh_2}{dh_1} \cdot \frac{dh_1}{dw_1}$$

$$\frac{dL}{dw_2} = \frac{dL}{dy} \cdot \frac{dy}{dh_2} \cdot \frac{dh_2}{dw_2}$$

$$h_1 = \frac{1}{1 + e^{-w_1 x}}$$

$$h_2 = \frac{1}{1 + e^{-w_2 h_1}}$$

$$y = \frac{1}{1 + e^{-w_3 h_2}}$$

$$L = -\log y$$

So we keep propagating till we reach the variable that is directly dependent on the differentiating variable.

*) Applying chain rule across multiple paths

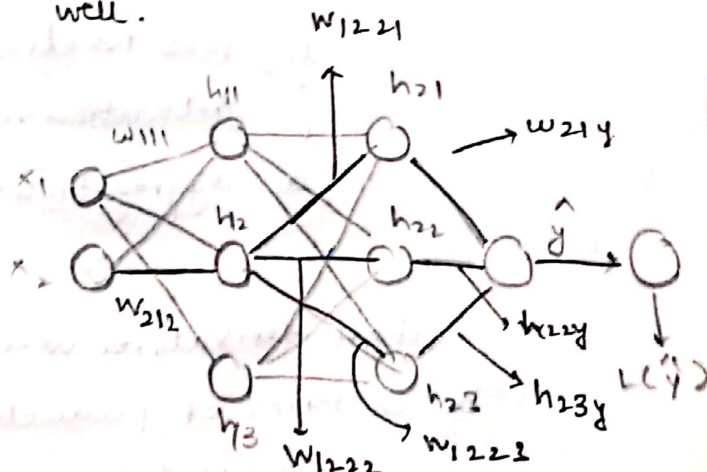


chain rule is applied across each path separately and then all derivatives are added to get the complete derivative term



since the weight that decides the value of ~~these~~ one neuron, might affect other neurons in deeper layers indirectly using the same neuron. i.e. the final loss might be impacted by one weight in multiple ways

This is a very simple network, but the basic concept remains the same for branched networks as well, and generalise really well.



$$\frac{\partial L}{\partial w_{212}} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial h_{22}} \cdot \frac{\partial h_{22}}{\partial a_{21}} \cdot \frac{\partial a_{21}}{\partial w_{212}} + \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial h_{21}} \cdot \frac{\partial h_{21}}{\partial a_{21}} \cdot \frac{\partial a_{21}}{\partial w_{212}}$$

↓
partial derivatives

$$+ \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial h_{22}} \cdot \frac{\partial h_{22}}{\partial a_{21}} \cdot \frac{\partial a_{21}}{\partial w_{212}}$$

↓

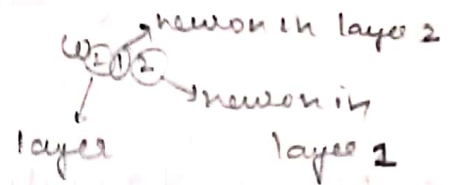
Now there were 3 paths wherein w_{212} was encountered

↓

hij here are the intermediate o/p's here which themselves contribute to the i/p of the next layer neurons which are then activated.

← The chain rule concept involves all paths wherein w_{212} or any weight was encountered to operate upon.

↓
These are themselves a composition of two ~~versatile~~ functions, but that should be pretty trivial.



$$\frac{\partial L}{\partial w_{212}} = \frac{\partial L}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial a_{21}} \cdot \frac{\partial a_{21}}{\partial w_{212}}$$

$$\hat{y}_1 = \frac{1}{1 + e^{-a_{21}}}$$

$$a_{21} = \sum w_{ij} x_{ij}$$

L : loss function

\hat{y} : activation function

a : aggregation function

↓

These derivatives when computed come out to be in terms of products of elementary values computed during forward propagation

Building in the dependency order or path might be more intuitive & safe in some scenarios.

↓ just for intuition

$$\frac{\partial L}{\partial w_{131}} = \frac{\partial L}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial a_{21}} \cdot \frac{\partial a_{21}}{\partial h_{13}} \cdot \frac{\partial h_{13}}{\partial a_{13}} \cdot \frac{\partial a_{13}}{\partial w_{131}}$$

$$+ \frac{\partial L}{\partial \hat{y}_2} \cdot \frac{\partial \hat{y}_2}{\partial a_{22}} \cdot \frac{\partial a_{22}}{\partial h_{13}} \cdot \frac{\partial h_{13}}{\partial a_{13}} \cdot \frac{\partial a_{13}}{\partial w_{131}}$$

↓

So the concept mainly is that the derivative for any weight gets broken into a chain of elementary computations.