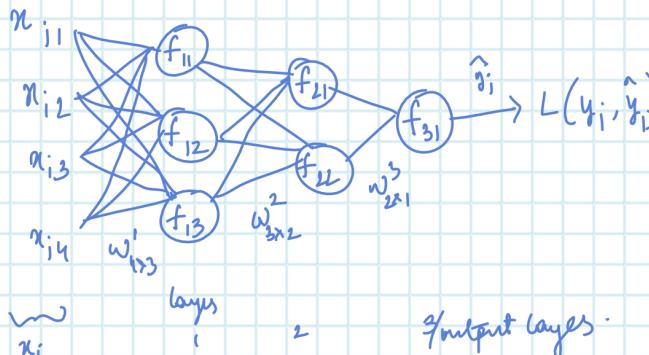


## 1.8 Training an MLP (Chain Rule) :-

$D = \{x_i, y_i\}$   $x_i \in \mathbb{R}^4$ ,  $y_i \in \mathbb{R}$  (example), loss is squared loss in regressions.



→ By the end we need to determine  $w_{1x3}^1, w_{2x2}^2, w_{2x1}^3 = 20$  weights (this is the training in MLP)

Step ① :-  $L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \text{reg}$  → L<sub>2</sub> reg equivalent would be  $\sum_{i,j,k} (|w_{ij}^k|)^2$   
 $\underbrace{\quad}_{\text{squared loss}}$       L<sub>1</sub> reg equivalent would be  $\sum_{i,j,k} (|w_{ij}^k|)$

$$L_i = (y_i - \hat{y}_i)^2$$

$$L = \sum_{i=1}^n L_i + \text{reg}$$

Step ② :- Optimization problem :- squared loss  $\leftrightarrow$  regularizations

$$\min_{w^1, w^2, w^3} L$$

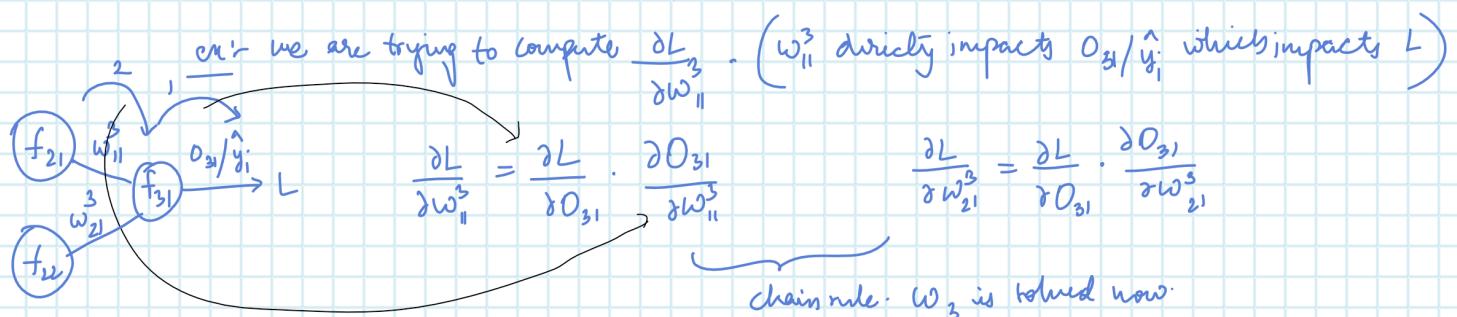
$$= \min_{\substack{w^1 \\ w^2 \\ w^3}} L$$

Step ③ :- SGD (or) GD

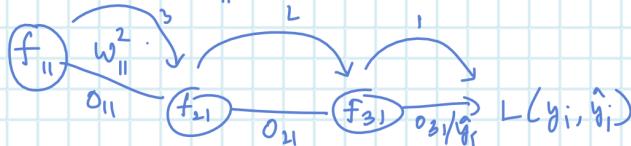
ⓐ Initialization of  $w_{i,j}^k$  randomly (there are lot of techniques to initialize them)

ⓑ  $(w_{i,j}^k)_{\text{new}} = (w_{i,j}^k)_{\text{old}} - \eta \left( \frac{\partial L}{\partial w_{i,j}^k} \right)$  → More important part  
 $\underbrace{\quad}_{\text{learning rate}}$

ⓒ Perform updates till convergence.



→ Now we calculate  $\frac{\partial L}{\partial w_{11}^2}$ .  $w_{11}^2$  impacts  $O_{11}, O_{21} \propto O_{31}/\hat{y}_1$

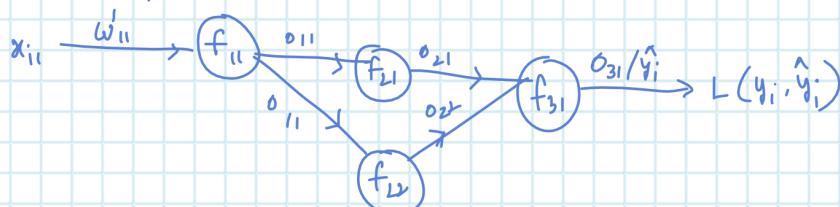


$$\Rightarrow \frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial w_{11}^2}$$

Similarly

$$\frac{\partial L}{\partial w_{21}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial w_{21}^2}$$

→ Calculating  $\frac{\partial L}{\partial w_{11}^1}$



In calculus,

$$x \rightarrow f(x) \rightarrow h(f(x), g(x)) \text{, then } \frac{\partial h}{\partial x} = \frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x}$$

We have

$$x \rightarrow f(x) \rightarrow h(f(x), g(x)) \rightarrow K(x) \cdot \frac{\partial K}{\partial x} = \frac{\partial K}{\partial h} \left( \frac{\partial h}{\partial x} \right)$$

$$\Rightarrow \frac{\partial K}{\partial x} = \frac{\partial h}{\partial h} \left[ \frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x} \right]$$

$$\rightarrow \frac{\partial L}{\partial w_{11}^1} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial w_{11}^1}$$

$$= \frac{\partial L}{\partial O_{31}} + \left[ \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w_{11}^1} + \frac{\partial O_{31}}{\partial O_{22}} \cdot \frac{\partial O_{22}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w_{11}^1} \right]$$

### Training an MLP: Memoization (1.1)

→ As seen in Dynamic Programming

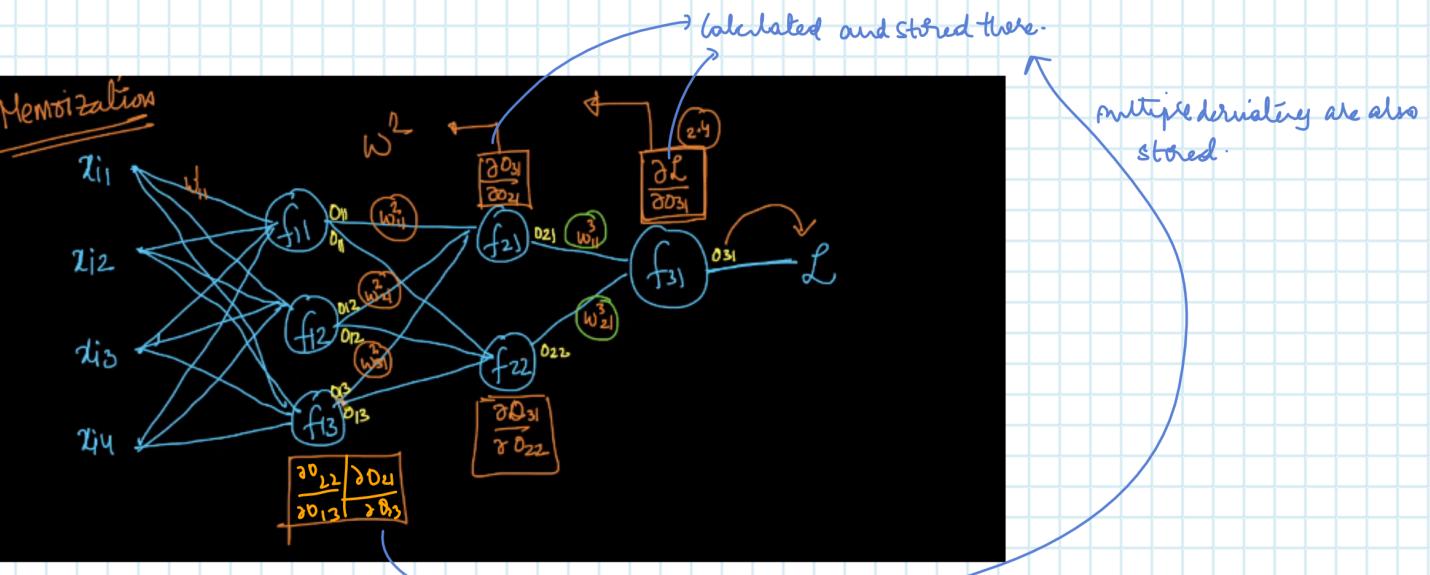
$$\rightarrow \frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial w_{11}^3}$$

same - core idea is to compute and reuse.

$$\frac{\partial L}{\partial w_{21}^3} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial w_{21}^3}$$

→ if there is any operations performed repeatedly, compute it once & reuse it -

→ takes slightly more memory though -



→ Chain Rule + Memoization = Backpropagation; with some optimization.

### Backpropagation (1 to 10) :-

Given  $D = \{x_i, y_i\}$

Step ① :- Initialize  $w_{ij}^k$

Step ② :- for each  $x_i$  in  $D$ , pass  $x_i$  forward through the network. The end result would be  $\hat{y}_i \rightarrow L(\hat{y}_i, y_i)$ .  
This step is called forward propagation.

Then compute the  $L(\hat{y}_i, y_i)$

Compute all derivatives using chain rule and memoizations.

Update weights from end of the network to the start  $\rightarrow (w_{ii}^3)_{\text{new}} = (w_{ii}^3)_{\text{old}} - \eta \frac{dL}{d w_{ii}^3}$

$$\begin{aligned}
 &\text{Backward propagation.} \\
 &(w_{ii}^2)_{\text{new}} = \dots \\
 &\text{We are using} \\
 &\text{the loss to} \\
 &\text{update the weights.} \\
 &(w_{ii}^1)_{\text{new}} = \dots \dots \dots
 \end{aligned}$$

Step ③ :- repeat step ② till convergence. i.e.,  $(w_{ii}^k)_{\text{new}} \approx (w_{ii}^k)_{\text{old}}$

epoch = Passed all the datapoints to the neural network once. We do multiple epochs in Real World.

Backpropagation = Multi epoch training methodology that uses chain rule & memoization to update weights.

→ Backpropagation works iff activation functions are differentiable. If it is easily differentiable, then training of NN can be sped up.

→ Instead of sending single point in backprop we can send a subset  $S \subseteq D$ , (similar to mini batch SGD)

$|S|$  = mini batch size - (typically powers of 2 like 64, 128, ...)

keeping all datapoints in RAM & computing gradients is hard.

∴ Mini batch SGD is most popular.

## Activation function (1-11) :-

→ In 1980 & 1990's, sigmoid and tanh were the most popular.

$$\text{let } z = \sum_{i=1}^n w_i x_i = w^T x$$

$$\text{if } f \text{ is Sigmoid then } f(z) = \frac{c}{1+e^{-z}} = \frac{1}{1+e^{-z}}$$

(This is also used in L<sup>R</sup>)

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

$$\frac{\partial r}{\partial z} = \underbrace{\sigma(z)}_{\downarrow} (1 - \underbrace{\sigma(z)}_{\downarrow})$$

Derivative of  $r(z)$  is expressed in terms of  $\sigma(z)$ .  
→ Can be helpful in backpropagation.

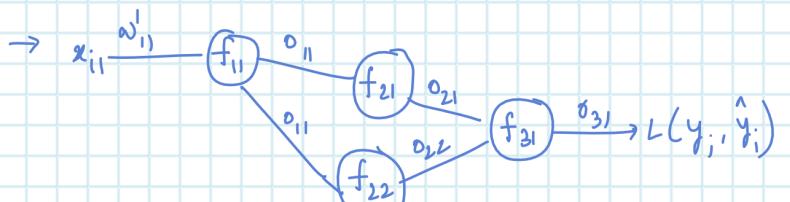
$$\rightarrow a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\frac{d \tanh}{dz} = 1 - \tanh^2(z)$$

→ ReLU are more popular activation function these days.

## Vanishing Gradient Problem:-

→ Most common reason NN lost hype back in 90's



$$\frac{\partial L}{\partial w_{i1}} = \frac{\partial L}{\partial o_{31}} \left[ \underbrace{\frac{\partial o_{31}}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{i1}}}_{0 \leq \dots < 1} + \underbrace{\frac{\partial o_{31}}{\partial o_{22}} \cdot \frac{\partial o_{22}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{i1}}}_{0 \leq \dots < 1} \right]$$

$$\frac{\partial o_{31}}{\partial o_{21}} = \frac{\partial f_{31}}{\partial o_{21}} \text{ since } o_{31} = f(\text{something})$$

All  $f$ 's are sigmoids and lie b/w 0 & 1 and the derivatives are being multiplied. Results in much smaller value.

As a result, in

$$(w'')_{\text{new}} = (w'')_{\text{old}} - \eta \left( \frac{\partial L}{\partial w_{i1}} \right)$$

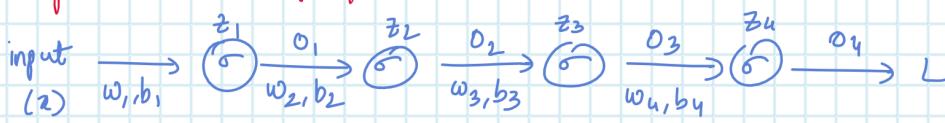
$(w'')_{\text{new}} \leftarrow (w'')_{\text{old}}$  will be very close.

This problem is called vanishing gradient.

$\frac{\partial L}{\partial w_{ij}^k}$  is small → vanishing gradient

$\frac{\partial L}{\partial w_{ij}^k}$  is very large → exploding gradient (happens when activation function's result > 1)

## Exploding Gradient using Sigmoid function:-



$$o_1 = \sigma(z_1)$$

$$z_1 = w_1 x + b_1$$

$$o_2 = \sigma(z_2)$$

$$z_2 = w_2 o_1 + b_2$$

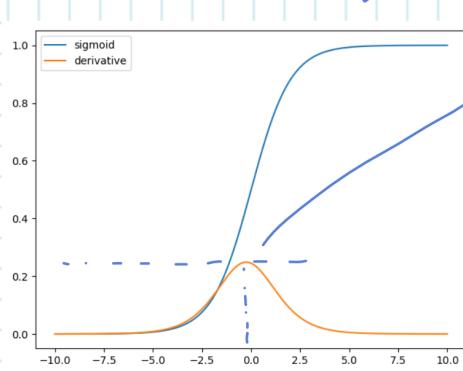
$$o_3 = \sigma(z_3)$$

$$z_3 = w_3 o_2 + b_3$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial o_4} \cdot \underbrace{\frac{\partial o_4}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial o_1}}_{= 0} \cdot \frac{\partial o_1}{\partial w_1}$$

if activation func such that each  $> 1$ , it leads to exploding gradient.

But sigmoid max is (0.25)



$$\text{take } \frac{\partial o_2}{\partial o_1} = \frac{\partial o_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial o_1}$$

$$= \frac{\partial \sigma(z)}{\partial z_2} \cdot \frac{\partial (w_2 o_1 + b_2)}{\partial o_1}$$

$$= \sigma'(z_2) + \underbrace{w_2}_{\downarrow} \quad \text{where } \sigma'(z) = \frac{\partial \sigma(z)}{\partial z}$$

weights are initialized randomly

If it is a large enough value, it leads to exploding gradient.

## Bias Variance Tradeoff (1.13):-

→ # layers ↑ → more weights → higher chance of overfit → high variance.

1 layer → aka log reg → higher chance of underfit → high bias.

→ MLP generally tends to have high variance since no one choose one layer.

→ to reduce variance we use regularization term along with loss in optimizations

$$L = \sum_{i=1}^n \text{loss}(y_i, \hat{y}_i) + \sum_{i,j,k} (w_{ij}^k)^2$$

$$= \sum_{i=1}^n \text{loss} + \lambda L_2$$

larger  $\lambda \Rightarrow$  lesser overfit

using  $L_1$  reg leads to sparsity (connections b/w neurons of various layers don't exist)