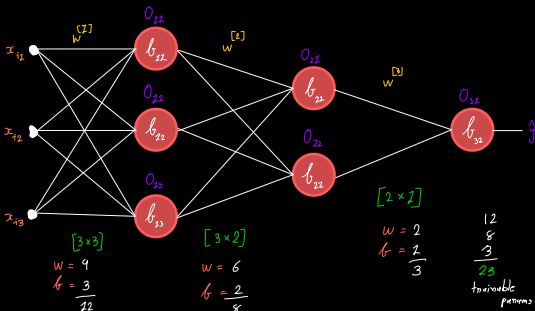


Memorization in Neural Network

Memorization :

in computing, memorization is optimization technique used primarily to speed-up computer programs by speed-up storing the results of expensive function calls & returning the cached results when the same input occurs again.

MLP - Memorization



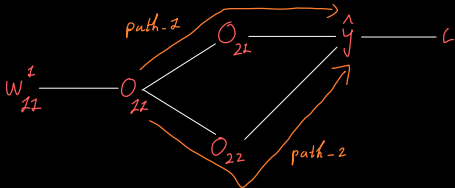
→ Based on above network's structure, we can see that $W^{[2]}$ derivative calculation process will be much more complex than $W^{[1]}$ and $W^{[3]}$.

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_{11}^3}$$

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{21}} \frac{\partial O_{21}}{\partial w_{11}^2}$$

$$L \rightarrow \hat{y} \rightarrow O_{21} \rightarrow O_{22} \rightarrow w_{11}^2$$





$$\frac{\partial L}{\partial w_{11}^1} = \frac{\partial L}{\partial \hat{y}} \times \left[\underbrace{\frac{\partial \hat{y}}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{11}^1}}_{\text{path-1}} + \underbrace{\frac{\partial \hat{y}}{\partial o_{22}} \times \frac{\partial o_{22}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{11}^1}}_{\text{path-2}} \right]$$

→ We can imagine that how things will be scary in big neural networks & how it will be difficult to calculate these values, hence we need to use **memorization technique**.

→ As we will go backward in neural-network, we need those derivatives which are already calculated. that's why we will store those values & reuse it.

Back propagation = Chain Rule + Memorization



$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}} \times \left[\frac{\partial \hat{y}}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial w_{11}^2} + \frac{\partial \hat{y}}{\partial O_{22}} \times \frac{\partial O_{22}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial w_{11}^2} \right]$$

→ As we can see in above formula that in both paths, $\frac{\partial O_{11}}{\partial w_{11}^2}$ is common hence we won't re-calculate it again & again.

→ We will store the o/p of $\frac{\partial O_{11}}{\partial w_{11}^2}$ in a memory and just use it.

