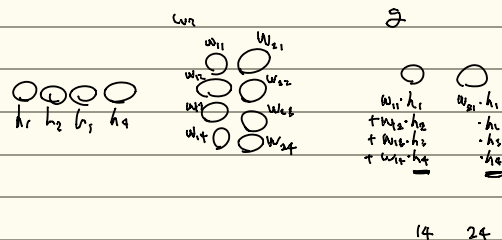
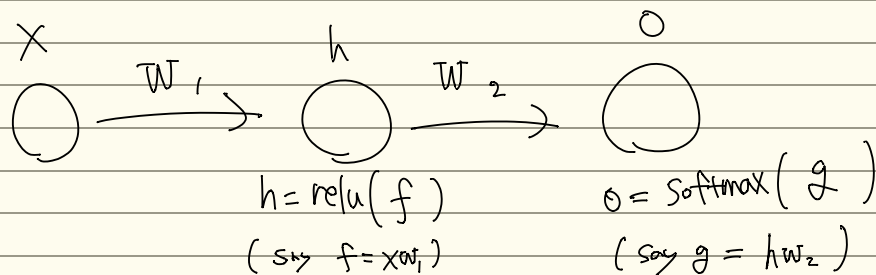


(for simplicity)



$E = - \sum_{j \in D} t_j \cdot \log(o_j)$

Now we want to compute

$\frac{\partial E}{\partial W_1}$  (12 weights),  $\frac{\partial E}{\partial W_2}$  (8 weights), total 20 weight's gradient.

$\frac{\partial E}{\partial W_2} = \frac{\partial E}{\partial o} \frac{\partial o}{\partial g} \frac{\partial g}{\partial W_2}$  (Note we here is not a matrix, just sample notation for any weight in  $W_2$ , such as  $w_{kj}$ ,  $k=1,2$ ,  $j=1,2,3,4$ )

$\frac{\partial E}{\partial W_1} = \frac{\partial E}{\partial o} \frac{\partial o}{\partial g} \frac{\partial g}{\partial h} \frac{\partial h}{\partial f} \frac{\partial f}{\partial W_1}$  (same for  $W_1$ ,  $W_1$  here means any weight in  $W_1$ , such as  $w_{ji}$ ,  $i=1,2,3$ )

Now we can easily see  $\frac{\partial E_d}{\partial o_n} = -\frac{t}{o_n}$ ,  $\frac{\partial o_n}{\partial g_k} = \begin{cases} o_n(1-o_n) & : k=n \\ -o_n \cdot o_k & : k \neq n \end{cases}$  (product of these two may be little bit confusing so let's just skip.)

$\therefore \frac{\partial E_d}{\partial o_n} \cdot \frac{\partial o_n}{\partial g_k} = o_k - t_k$

and for  $\frac{\partial g_k}{\partial w_{kj}} = h_j$   $\therefore \frac{\partial E_d}{\partial w_{kj}} = (o_k - t_k) h_j$  (with  $\sum_o \nabla E$  for  $w_{kj}$ )

Now for  $W_1$ ,  $\frac{\partial g_k}{\partial h_j} = w_{kj}$ ,  $\frac{\partial h}{\partial f} = \begin{cases} 0 & : f \leq 0 \\ 1 & : f > 0 \end{cases}$ ,  $\frac{\partial f_j}{\partial w_{ji}} = x_i$

$\therefore \frac{\partial E_d}{\partial w_{ji}} = \sum_{k=1}^2 (o_k - t_k) \cdot w_{kj} \cdot (\text{if}(f > 0)) \cdot x_i$  (1 = true, 0 = false)  $\Rightarrow \nabla E$  for  $w_{ji}$

for weight in  $W_1$ , we must consider all output nodes ( $k=1 \sim 2$ )