

Now we have found expressions for  $\frac{\partial}{\partial \theta_j} [C(\theta)]$  for each  $\theta_j \in \Theta$ .

To do next time: ① write iterative equation for gradient descent

② try the same exercise for a 2-hidden-layer network

### ① Iterative egn for gradient descent (0-hidden-layer network).

For each  $\theta_j \in \Theta = \{b^{(0)}\} \cup \{w_{i,j}^{(1)} : i \in [0, 783] \cap \mathbb{Z}, j \in [0, 9] \cap \mathbb{Z}\}$ ,

- Initialise  $\theta_j^{[0]} := \text{random} (\text{between } -10 \text{ & } 10 \text{ idk?})$

- Step.  $\theta_j^{[n+1]} := \theta_j^{[n]} - \nu \frac{\partial}{\partial \theta_j} [C(\theta)]$  where  $\nu > 0$  is "learning rate".

Can code logic for Step to repeat until we have hit some max number of iterations OR difference between this step vs previous is negligible.

### ② 2-hidden-layer network

Let  $L_i$  denote the number of neurons in layer  $i$ . We want this:

- $L_0 = 783$  (since 784 pixels in MNIST images)
- $L_1 = 127$
- $L_2 = 63$
- $L_3 = 9$  (10 digits in final layer)

Otherwise, structure is identical to the 0-hidden-layer network I've been working on. So, in particular:

- Each non-zero layer has a bias parameter  $b^{(i)}$  added to the weighted sums evaluated in that layer (i.e.  $b^{(1)}, b^{(2)}, b^{(3)}$ )
- Each possible pairwise connection between layer  $i-1$  and the next layer  $i$  has a weight parameter associated (i.e. for layer  $i$ ,  $W^{(i)} = \{w_{i,j}^{(i)} : i \in [0, L_i] \cap \mathbb{Z}, j \in [0, L_{i+1}] \cap \mathbb{Z}\}$ .)
- ReLU activation used:  $f(x) : \mathbb{R} \rightarrow \mathbb{R}$  where  $f(x) := \begin{cases} x & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$
- Cross-entropy loss used for training:

$$C(\theta) = \frac{-1}{K} \sum_{k=0}^{K-1} \sum_{j=0}^9 t_{j,k} \log(\hat{p}_{j,k}) \text{ as previously defined!}$$

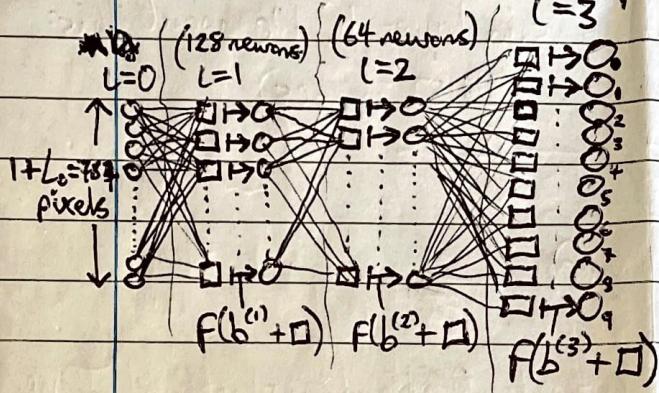
Before differentiating  $C(\underline{\theta})$  by all these new parameters:

$$\underline{\theta} = \{b^{(1)}, b^{(2)}, b^{(3)}\} \quad \text{and} \quad \cancel{UW}^{(1)} \cancel{UTJ}^{(1)} \cancel{UW}^{(2)} \cancel{UTJ}^{(2)} \cancel{UW}^{(3)} \cancel{UTJ}^{(3)}$$

let us write clearly how neuron values in  $C(\underline{\theta})$  are defined.

- $p_{n,k}^{(0)} := n^{\text{th}}$  pixel value in training image  $k$ . (valid for  $n \in \{0, 1, \dots, 783\}$ )
- $p_{n,k}^{(1)} = f(b^{(1)} + \sum_{i=0}^{783} w_{i,n} p_{i,k}^{(0)})$  for  $n \in \{0, 1, \dots, L_1 = 127\}$
- $p_{n,k}^{(2)} = f(b^{(2)} + \sum_{i=0}^{127} w_{i,n} p_{i,k}^{(1)})$  for  $n \in \{0, 1, \dots, L_2 = 63\}$
- $p_{n,k}^{(3)} = f(b^{(3)} + \sum_{i=0}^{63} w_{i,n} p_{i,k}^{(2)})$  for  $n \in \{0, 1, \dots, L_3 = 9\}$
- Then,  $\hat{p}_{n,k} = \frac{\exp(p_{n,k}^{(3)})}{\sum_{z=0}^9 \exp(p_{n,k}^{(z)})}$

So, in general we have  $p_{n,k}^{(l)} = f(b^{(l)} + \sum_{i=0}^{L_{l-1}} w_{i,n} p_{i,k}^{(l-1)}) \quad \forall l \in \{0, 1, 2, 3\}$



then the digit from 0-9 with the highest neuron-value in this final layer is our network's prediction.

(diagram, where  $\square$  represents each weighted sum of neurons)

Now we have everything we need, on the next sheet I'll find  $\frac{\partial}{\partial \theta_j} [C(\underline{\theta})]$  for all  $\theta_j \in \underline{\theta} \dots$

(start with biases, then do weights. Work from layer 3 to layer 1). I will use this fact in each case to kick things off:

$$\frac{\partial}{\partial \theta_j} [C(\underline{\theta})] = \frac{-1}{K} \sum_{k=0}^{K-1} \sum_{m=0}^9 t_{m,k} \frac{1}{\hat{p}_{m,k}} \frac{\partial}{\partial \theta_j} [\hat{p}_{m,k}]$$

$$\text{Da Nang} \rightarrow \text{Hanoi cont.} \quad \frac{\partial}{\partial b^{(3)}} [C(\underline{\theta})] = \frac{-1}{K} \sum_{k=0}^{K-1} \sum_{j=0}^q t_{j,k} \frac{\partial}{\partial b^{(3)}} [\log(\hat{p}_{j,k})]$$

$$\boxed{b^{(3)}} \quad \frac{\partial}{\partial b^{(3)}} [C(\underline{\theta})] = \frac{-1}{K} \sum_{k=0}^{K-1} \sum_{m=0}^q \left( \frac{t_{m,k}}{\hat{p}_{m,k}} \frac{\partial}{\partial b^{(3)}} [\underline{C(\underline{\theta})}] \frac{\partial}{\partial b^{(3)}} [\hat{p}_{m,k}] \right)$$

Now,  $\frac{\partial}{\partial b^{(3)}} [\exp\{\hat{p}_{m,k}^{(3)}\}] = \frac{\partial}{\partial b^{(3)}} [\hat{p}_{m,k}] \frac{\partial}{\partial b^{(3)}} [\hat{p}_{m,k}]$

chain rule  
 $\frac{\partial}{\partial x} [e^x] = e^x$

$$= \exp\{\hat{p}_{m,k}\} \frac{\partial}{\partial b^{(3)}} [\hat{p}_{m,k}]$$

$$= \exp\{\hat{p}_{m,k}\} \underbrace{\frac{\partial}{\partial x} [f(x)]}_{= f'(x)} \frac{\partial}{\partial b^{(3)}} \left[ b^{(3)} + \sum_{i=0}^{63} w_{i,m}^{(3)} p_{i,k}^{(2)} \right]$$

using our old friend defined earlier  
 $\sum_{j=0}^q t_{j,k}$  := sum of  $j^{\text{th}}$  bias + weighted sum going into  $j^{\text{th}}$

$\bullet$  neuron in layer  $l'$   
and the fact  $f'(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$

$$= \exp\{\hat{p}_{m,k}\} \prod_{j=0}^q \{ \sum_{i=0}^{63} w_{i,m}^{(3)} p_{i,k}^{(2)} > 0 \}$$

So, using the quotient rule in essentially the same way we did for the  $b^{(1)}$  partial derivative in my 0-hidden-layer network:

$$\hat{p}_{m,k} = \frac{\exp\{\hat{p}_{m,k}^{(3)}\}}{\sum_{z=0}^q \exp\{\hat{p}_{z,k}^{(3)}\}} = \frac{U}{V} \quad U = \exp\{\hat{p}_{m,k}^{(3)}\} \prod_{j=0}^q \{ \sum_{i=0}^{63} w_{i,m}^{(3)} p_{i,k}^{(2)} > 0 \},$$

$$V = \sum_{z=0}^q \exp\{\hat{p}_{z,k}^{(3)}\} \prod_{j=0}^q \{ \sum_{i=0}^{63} w_{i,z}^{(3)} p_{i,k}^{(2)} > 0 \}.$$

(derived earlier so)  
copied here!

$$\Rightarrow \frac{\partial}{\partial b^{(3)}} [\hat{p}_{m,k}] = \hat{p}_{m,k} \left( \prod_{j=0}^q \{ \sum_{i=0}^{63} w_{i,m}^{(3)} p_{i,k}^{(2)} > 0 \} - \sum_{z=0}^q \hat{p}_{z,k} \prod_{j=0}^q \{ \sum_{i=0}^{63} w_{i,z}^{(3)} p_{i,k}^{(2)} > 0 \} \right)$$

$$\Rightarrow \frac{\partial}{\partial b^{(3)}} [C(\underline{\theta})] = \frac{-1}{K} \sum_{k=0}^{K-1} \sum_{m=0}^q t_{m,k} \left( \prod_{j=0}^q \{ \sum_{i=0}^{63} w_{i,m}^{(3)} p_{i,k}^{(2)} > 0 \} - \sum_{z=0}^q \hat{p}_{z,k} \prod_{j=0}^q \{ \sum_{i=0}^{63} w_{i,z}^{(3)} p_{i,k}^{(2)} > 0 \} \right)$$

$$\boxed{b^{(2)}} \quad \frac{\partial}{\partial b^{(2)}} [C(\underline{\theta})] = \frac{-1}{K} \sum_{k=0}^{K-1} \sum_{m=0}^q \left( \frac{t_{m,k}}{\hat{p}_{m,k}} \frac{\partial}{\partial b^{(2)}} [\hat{p}_{m,k}] \right).$$

Now,  $\frac{\partial}{\partial b^{(2)}} [\exp\{\hat{p}_{m,k}^{(3)}\}] = \exp\{\hat{p}_{m,k}^{(3)}\} \frac{\partial}{\partial b^{(2)}} [\hat{p}_{m,k}^{(3)}]$

$$= \exp\{\hat{p}_{m,k}^{(3)}\} \prod_{j=0}^q \{ \sum_{i=0}^{63} w_{i,m}^{(3)} p_{i,k}^{(2)} > 0 \} \frac{\partial}{\partial b^{(2)}} \left[ b^{(2)} + \sum_{i=0}^{63} w_{i,m}^{(3)} p_{i,k}^{(2)} \right]$$

$\prod_{j=0}^q \{ \sum_{i=0}^{63} w_{i,m}^{(3)} p_{i,k}^{(2)} > 0 \}$  term comes from differentiating ReLU.

$$\text{Now, } \frac{\partial}{\partial b^{(2)}} [p_{i,k}^{(2)}] = \frac{\partial}{\partial b^{(2)}} \left[ b^{(2)} + \sum_{j=0}^{127} w_{i,j}^{(2)} p_{j,k}^{(1)} \right] \cdot \frac{\partial}{\partial x} [f(x)] \rightarrow x = \sum_{i,k}^{(2)}$$

$$= 1 \cdot \prod \left\{ \sum_{i,k}^{(2)} > 0 \right\}$$

$$\text{So, } \frac{\partial}{\partial b^{(2)}} [\exp \{ p_{m,k}^{(3)} \}] = \exp \{ p_{m,k}^{(3)} \} \prod \left\{ \sum_{i,k}^{(3)} > 0 \right\} \left( \sum_{i=0}^{63} w_{i,m}^{(3)} \prod \left\{ \sum_{i,k}^{(2)} > 0 \right\} \right)$$

So, use quotient rule on the following:

$$\hat{p}_{m,k} = \frac{u}{v} = \frac{\exp \{ p_{m,k}^{(3)} \}}{\sum_{z=0}^q \exp \{ p_{z,k}^{(3)} \}}$$

$u$  is given here

$$v' = \sum_{z=0}^q \exp \{ p_{z,k}^{(3)} \} \prod \left\{ \sum_{i,k}^{(3)} > 0 \right\} \left( \sum_{i=0}^{63} w_{i,z}^{(3)} \prod \left\{ \sum_{i,k}^{(2)} > 0 \right\} \right)$$

$$\text{So, } \frac{\partial}{\partial b^{(2)}} \left[ \hat{p}_{m,k} \right] = \frac{vu - uv'}{v^2}$$

$$= \frac{1}{\left( \sum_{z=0}^q \exp \{ p_{z,k}^{(3)} \} \right)^2} \left( \left( \sum_{z=0}^q \exp \{ p_{z,k}^{(3)} \} \right) \exp \{ p_{m,k}^{(3)} \} \prod \left\{ \sum_{i,k}^{(3)} > 0 \right\} \left( \sum_{i=0}^{63} w_{i,m}^{(3)} \prod \left\{ \sum_{i,k}^{(2)} > 0 \right\} \right) \right. \\ \left. - \left( \exp \{ p_{m,k}^{(3)} \} \right) \sum_{z=0}^q \exp \{ p_{z,k}^{(3)} \} \prod \left\{ \sum_{i,k}^{(3)} > 0 \right\} \left( \sum_{i=0}^{63} w_{i,z}^{(3)} \prod \left\{ \sum_{i,k}^{(2)} > 0 \right\} \right) \right)$$

$$= \hat{p}_{m,k} \left( \sum_{i=0}^{63} w_{i,m}^{(3)} \prod \left\{ \sum_{i,k}^{(3)} > 0, \sum_{i,k}^{(2)} > 0 \right\} - \sum_{z=0}^q \hat{p}_{z,k} \sum_{i=0}^{63} w_{i,z}^{(3)} \prod \left\{ \sum_{i,k}^{(3)} > 0, \sum_{i,k}^{(2)} > 0 \right\} \right)$$

$$\Rightarrow \frac{\partial}{\partial b^{(2)}} [C(\theta)] = \frac{-1}{K} \sum_{k=0}^{K-1} \sum_{m=0}^q t_{m,k} \left( \sum_{i=0}^{63} w_{i,m}^{(3)} \prod \left\{ \sum_{i,k}^{(3)} > 0, \sum_{i,k}^{(2)} > 0 \right\} - \sum_{z=0}^q \hat{p}_{z,k} \sum_{i=0}^{63} w_{i,z}^{(3)} \prod \left\{ \sum_{i,k}^{(3)} > 0, \sum_{i,k}^{(2)} > 0 \right\} \right)$$

$$\boxed{b^{(1)}} \frac{\partial}{\partial b^{(1)}} [C(\theta)] = \frac{-1}{K} \sum_{k=0}^{K-1} \sum_{m=0}^q \left( \frac{t_{m,k}}{\hat{p}_{m,k}} \frac{\partial}{\partial b^{(1)}} [\hat{p}_{m,k}] \right)$$

$$\text{Now, } \frac{\partial}{\partial b^{(1)}} [\exp \{ p_{m,k}^{(3)} \}] = \exp \{ p_{m,k}^{(3)} \} \frac{\partial}{\partial b^{(1)}} [\hat{p}_{m,k}]$$

$$= \exp \{ p_{m,k}^{(3)} \} \prod \left\{ \sum_{i,k}^{(3)} > 0 \right\} \frac{\partial}{\partial b^{(1)}} \left[ b^{(3)} + \sum_{i=0}^{63} w_{i,m}^{(3)} p_{i,k}^{(2)} \right]$$

$$= \exp \{ p_{m,k}^{(3)} \} \prod \left\{ \sum_{i,k}^{(3)} > 0 \right\} \left( \sum_{i=0}^{63} w_{i,m}^{(3)} \frac{\partial}{\partial b^{(1)}} [\hat{p}_{i,k}^{(2)}] \right)$$

Continued next page...

$$\text{So, } \frac{\partial}{\partial b^{(1)}} [C(\underline{\theta})] = \frac{-1}{K} \sum_{k=0}^{K-1} \sum_{m=0}^q t_{m,k} \left( \sum_{i=0}^{63} \sum_{j=0}^{127} I_{m,k}^{(3)} I_{i,k}^{(3)} I_{j,i,k}^{(2)} - \sum_{z=0}^q \sum_{l=0}^{63} \sum_{s=0}^{127} \hat{P}_{z,k}^{(3)} I_{z,k}^{(3)} I_{s,z,k}^{(2)} I_{s,k}^{(1)} \right)$$

$$W_{\alpha,\beta}^{(3)}$$

$$\frac{\partial}{\partial w_{\alpha,\beta}^{(3)}} [C(\underline{\theta})] = \frac{-1}{K} \sum_{k=0}^{K-1} \sum_{m=0}^q t_{m,k} \frac{\partial}{\partial w_{\alpha,\beta}^{(3)}} [\hat{P}_{m,k}]. \quad \text{Now,}$$

$$\begin{aligned} \frac{\partial}{\partial w_{\alpha,\beta}^{(3)}} \left[ \exp \left\{ \hat{P}_{m,k}^{(3)} \right\} \right] &= \exp \left\{ \hat{P}_{m,k}^{(3)} \right\} \frac{\partial}{\partial w_{\alpha,\beta}^{(3)}} \left[ \hat{P}_{m,k}^{(3)} \right] \\ &= \exp \left\{ \hat{P}_{m,k}^{(3)} \right\} \left[ \frac{\partial}{\partial w_{\alpha,\beta}^{(3)}} \left[ b^{(3)} + \sum_{i=0}^{63} W_{i,m}^{(3)} P_{i,k}^{(2)} \right] \right] \\ &= \begin{cases} 0 & \text{if } m \neq \beta \\ \exp \left\{ \hat{P}_{\beta,k}^{(3)} \right\} - \hat{P}_{\beta,k}^{(3)} P_{\alpha,k}^{(2)} & \text{otherwise} \end{cases} \end{aligned}$$

So use quotient rule on the following:

$$\hat{P}_{m,k} = \frac{U}{V} = \frac{\exp \left\{ \hat{P}_{m,k}^{(3)} \right\}}{\sum_{z=0}^q \exp \left\{ \hat{P}_{z,k}^{(3)} \right\}}$$

$U'$  is here:

$$U' = \exp \left\{ \hat{P}_{\beta,k}^{(3)} \right\} - \hat{P}_{\beta,k}^{(3)} P_{\alpha,k}^{(2)}$$

$$V' = \exp \left\{ \hat{P}_{\beta,k}^{(3)} \right\} \cdot \hat{P}_{\beta,k}^{(3)} P_{\alpha,k}^{(2)}$$

$$\text{So, } \frac{\partial}{\partial w_{\alpha,\beta}^{(3)}} \left[ \hat{P}_{m,k} \right] = \frac{V U' - U V'}{V^2}$$

$$= \frac{1}{\left( \sum_{z=0}^q \exp \left\{ \hat{P}_{z,k}^{(3)} \right\} \right)^2} \left( \left( \sum_{z=0}^q \exp \left\{ \hat{P}_{z,k}^{(3)} \right\} \right) \exp \left\{ \hat{P}_{\beta,k}^{(3)} \right\} - \hat{P}_{\beta,k}^{(3)} P_{\alpha,k}^{(2)} \right) \left( \sum_{z=0}^q \exp \left\{ \hat{P}_{z,k}^{(3)} \right\} \left( \exp \left\{ \hat{P}_{\beta,k}^{(3)} \right\} - \hat{P}_{\beta,k}^{(3)} P_{\alpha,k}^{(2)} \right) \right)$$

$$= \hat{P}_{m,k} \left( \sum_{z=0}^q \exp \left\{ \hat{P}_{z,k}^{(3)} \right\} - \hat{P}_{\beta,k}^{(3)} P_{\alpha,k}^{(2)} \right)$$

can take:  
 $\left( \sum_{z=0}^q \exp \left\{ \hat{P}_{z,k}^{(3)} \right\} - \hat{P}_{\beta,k}^{(3)} P_{\alpha,k}^{(2)} \right)$

Therefore we arrive at:

$$\frac{\partial}{\partial w_{\alpha,\beta}^{(3)}} [C(\underline{\theta})] = \frac{-1}{K} \sum_{k=0}^{K-1} \sum_{m=0}^q t_{m,k} \left( \sum_{z=0}^q \exp \left\{ \hat{P}_{z,k}^{(3)} \right\} - \hat{P}_{\beta,k}^{(3)} P_{\alpha,k}^{(2)} \right) \hat{P}_{\beta,k}^{(3)} P_{\alpha,k}^{(2)}$$

Still Guangzhou  $\rightarrow$  London

05/02/2025

$$W_{\alpha,\beta}^{(2)}$$

Start with  $\frac{\partial}{\partial W_{\alpha,\beta}^{(2)}} \left[ \exp \left\{ p_{m,k}^{(3)} \right\} \right] = \dots$

$$\dots = \exp \left\{ p_{m,k}^{(3)} \right\} \left[ m_k \frac{\partial}{\partial W_{\alpha,\beta}^{(2)}} \left[ b^{(3)} + \sum_{i=0}^{63} w_{i,m}^{(3)} p_{i,k}^{(2)} \right] \right]$$

$$= \sum_{i=0}^{63} w_{i,m}^{(3)} \frac{\partial}{\partial W_{\alpha,\beta}^{(2)}} \left[ p_{i,k}^{(2)} \right]$$

Now,  $\frac{\partial}{\partial W_{\alpha,\beta}^{(2)}} \left[ p_{i,k}^{(2)} \right] = I_{i,k}^{(2)} \frac{\partial}{\partial W_{\alpha,\beta}^{(2)}} \left[ b^{(2)} + \sum_{j=0}^{127} w_{j,i}^{(2)} p_{j,k}^{(1)} \right]$

$$= I_{i,k}^{(2)} \prod_{\{i=\beta\}} p_{\alpha,k}^{(1)}$$

$$\left( = I_{B,K}^{(2)} \prod_{\{i=\beta\}} p_{\alpha,k}^{(1)} \right)$$

$$\Rightarrow \frac{\partial}{\partial W_{\alpha,\beta}^{(2)}} \left[ \exp \left\{ p_{m,k}^{(3)} \right\} \right] = \exp \left\{ p_{m,k}^{(3)} \right\} I_{m,k}^{(3)} \sum_{i=0}^{63} w_{i,m}^{(3)} I_{B,K}^{(2)} p_{\alpha,k}^{(1)} \prod_{\{i=\beta\}}$$

Interestingly, there is no index indicator here unlike

$$= \exp \left\{ p_{m,k}^{(3)} \right\} I_{m,k}^{(3)} w_{B,m}^{(3)} I_{B,K}^{(2)} p_{\alpha,k}^{(1)}$$

the  $\frac{\partial}{\partial W_{\alpha,\beta}^{(2)}}$  derivative. I think this is because the final layer ~~as~~ weights do not contribute to ALL 9 OF THE FINAL NEURONS, whereas each 2nd (and 1st) layer weight does contribute to EVERY final neuron!

QUOTIENT RULE  $\hat{P}_{m,k}^{(3)} = \frac{\exp \left\{ p_{m,k}^{(3)} \right\}}{\sum_{z=0}^9 \exp \left\{ p_{z,k}^{(3)} \right\}} = \frac{U}{V}$   $U'$  is above.

$$V' = \sum_{z=0}^9 \exp \left\{ p_{z,k}^{(3)} \right\} I_{z,k}^{(3)} w_{B,z}^{(3)} I_{B,K}^{(2)} p_{\alpha,k}^{(1)}$$

$$\frac{\partial}{\partial W_{\alpha,\beta}^{(2)}} \left[ \hat{P}_{m,k}^{(3)} \right] = \frac{VU' - UV'}{V^2}$$

$$= \frac{1}{\left( \sum_{z=0}^9 \exp \left\{ p_{z,k}^{(3)} \right\} \right)^2} \left( \sum_{z=0}^9 \exp \left\{ p_{z,k}^{(3)} \right\} \exp \left\{ p_{m,k}^{(3)} \right\} I_{m,k}^{(3)} w_{B,m}^{(3)} I_{B,K}^{(2)} p_{\alpha,k}^{(1)} \right) = \dots$$

$$= \left( \sum_{z=0}^9 \exp \left\{ p_{z,k}^{(3)} \right\} \right)^2 \left( -\exp \left\{ p_{m,k}^{(3)} \right\} \sum_{z=0}^9 \exp \left\{ p_{z,k}^{(3)} \right\} I_{z,k}^{(3)} w_{B,z}^{(3)} I_{B,K}^{(2)} p_{\alpha,k}^{(1)} \right) = \dots$$

$$= \hat{P}_{m,k} \left( I_{m,k}^{(3)} w_{B,m}^{(3)} + I_{B,k}^{(2)} p_{\alpha,k}^{(1)} - \sum_{z=0}^q \hat{P}_{z,k} w_{B,z}^{(3)} I_{z,k}^{(3)} I_{B,k}^{(2)} p_{\alpha,k}^{(1)} \right)$$

$$\Rightarrow \frac{\partial}{\partial w_{\alpha,B}^{(1)}} [C(\theta)] = \frac{-1}{K} \sum_{k=0}^{K-1} \sum_{m=0}^q I_{m,k}^{(3)} I_{B,k}^{(2)} p_{\alpha,k}^{(1)} \left( I_{m,k}^{(3)} w_{B,m}^{(3)} - \sum_{z=0}^q \hat{P}_{z,k} w_{B,z}^{(3)} I_{z,k}^{(3)} I_{B,k}^{(2)} p_{\alpha,k}^{(1)} \right)$$

$$w_{\alpha,B}^{(1)} \quad \text{Start with } \frac{\partial}{\partial w_{\alpha,B}^{(1)}} [\exp \{ \cdot \cdot \cdot p_{m,k}^{(3)} \}] = \dots$$

$$\dots = \exp \{ p_{m,k}^{(3)} \} \underbrace{I_{m,k}^{(3)} \sum_{i=0}^{63} w_{i,m}}_{(1)} \frac{\partial}{\partial w_{\alpha,B}^{(1)}} [p_{i,k}^{(2)}]$$

$$\text{Now } \frac{\partial}{\partial w_{\alpha,B}^{(1)}} [p_{i,k}^{(2)}] = I_{i,k}^{(2)} \frac{\partial}{\partial w_{\alpha,B}^{(1)}} [b^{(2)}] + \sum_{j=0}^{127} w_{j,i}^{(2)} p_{j,k}^{(1)}$$

$$= I_{i,k}^{(2)} \sum_{j=0}^{127} w_{j,i}^{(2)} \frac{\partial}{\partial w_{\alpha,B}^{(1)}} [p_{j,k}^{(1)}]$$

$$= I_{i,k}^{(2)} \sum_{j=0}^{127} w_{j,i}^{(2)} \underbrace{I_{j,k}^{(1)} \frac{\partial}{\partial w_{\alpha,B}^{(1)}} [b^{(1)}]}_{(1)} + \sum_{n=0}^{783} w_{n,i}^{(1)} p_{n,k}^{(0)}$$

$$= I_{i,k}^{(2)} \sum_{j=0}^{127} w_{j,i}^{(2)} I_{j,k}^{(1)} \underbrace{\sum_{j=0}^{127} w_{j,j}^{(1)}}_{(1)} p_{\alpha,k}^{(0)}$$

$$= I_{i,k}^{(2)} w_{B,i}^{(2)} I_{B,k}^{(1)} p_{\alpha,k}^{(0)} \quad \begin{aligned} & \text{(since summing over } j \text{)} \\ & \text{(zeroes all except } j=B \text{)} \end{aligned}$$

So, substituting back in:

$$\frac{\partial}{\partial w_{\alpha,B}^{(1)}} [\exp \{ p_{m,k}^{(3)} \}] = \exp \{ p_{m,k}^{(3)} \} I_{m,k}^{(3)} \sum_{i=0}^{63} w_{i,m}^{(3)} I_{i,k}^{(2)} w_{B,i}^{(2)} I_{B,k}^{(1)} p_{\alpha,k}^{(0)}$$

$$\text{(QUOTIENT RULE)} \quad \hat{P}_{m,k} = \frac{u}{v} = \frac{\exp \{ p_{m,k}^{(3)} \}}{\sum_{z=0}^q \exp \{ p_{z,k}^{(3)} \}}. \quad u = \sum_{z=0}^q \exp \{ p_{z,k}^{(3)} \}, \quad v = \sum_{i=0}^{127} w_{i,m}^{(3)} I_{i,k}^{(2)} w_{B,i}^{(2)} I_{B,k}^{(1)} p_{\alpha,k}^{(0)}$$

$$\Rightarrow \frac{\partial}{\partial w_{\alpha,B}^{(1)}} [\hat{P}_{m,k}] = \frac{vu' - uv'}{v^2} = \hat{P}_{m,k} \left( I_{m,k}^{(3)} \sum_{i=0}^{63} w_{i,m}^{(3)} I_{i,k}^{(2)} w_{B,i}^{(2)} I_{B,k}^{(1)} p_{\alpha,k}^{(0)} \right) - \sum_{z=0}^q \hat{P}_{z,k} I_{z,k}^{(3)} \sum_{i=0}^{63} w_{i,z}^{(3)} I_{i,k}^{(2)} w_{B,i}^{(2)} I_{B,k}^{(1)} p_{\alpha,k}^{(0)}$$

$$\Rightarrow \frac{\partial}{\partial w_{\alpha,B}^{(1)}} [C(\theta)] = \frac{-1}{K} \sum_{k=0}^{K-1} \sum_{m=0}^q I_{m,k}^{(3)} I_{B,k}^{(2)} p_{\alpha,k}^{(1)} \left( I_{m,k}^{(3)} \sum_{i=0}^{63} w_{i,m}^{(3)} I_{i,k}^{(2)} w_{B,i}^{(2)} I_{B,k}^{(1)} p_{\alpha,k}^{(0)} - \sum_{z=0}^q \hat{P}_{z,k} I_{z,k}^{(3)} \sum_{i=0}^{63} w_{i,z}^{(3)} I_{i,k}^{(2)} w_{B,i}^{(2)} I_{B,k}^{(1)} p_{\alpha,k}^{(0)} \right)$$