

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

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

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

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

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

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

- Step $\theta_j^{[n+1]} := \theta_j^{[n]} - \nu \frac{\partial}{\partial \theta_j} [C(\underline{\theta})]$ 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:
(minus 1, due to 0 indexing)

- $L_0 = 783$ (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 ~~$i+1$~~ has a weight parameter associated
(i.e. for layer i , $W^{(i)} = \{w_{i,j}^{(i)} : i \in [0, L_i] \cap \mathbb{Z}, j \in [0, L_{i+1}] \cap \mathbb{Z}\}$). ~~$i+1$~~

- 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:

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

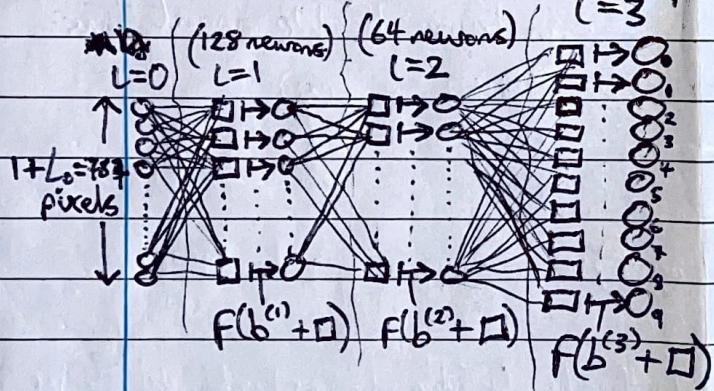
Before differentiating $C(\Theta)$ by all these new parameters:

$$\underline{\Theta} = \{b^{(1)}, b^{(2)}, b^{(3)}\} \quad \text{and} \quad U \nabla J^{(1)} \quad U \nabla J^{(2)} \quad U \nabla J^{(3)}$$

let us write clearly how neuron values in $C(\Theta)$ 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, $\hat{p}_{n,k} = \frac{\exp(p_{n,k}^{(3)})}{\sum_{z=0}^8 \exp(p_{z,k}^{(3)})}$

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 $\frac{\partial}{\partial \theta_j} [C(\Theta)]$ 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 part in each case to kick things off:

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

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

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

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

using our old friend defined earlier $\sum_{j=0}^{i-1} t_{j,k}$:= sum of j^{th} bias + weighted sum going into j^{th} neuron in layer l'

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

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

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

$$= \exp\{\hat{p}_{m,k}\} \prod_{z=0}^q \underbrace{\left\{ \sum_{m,k}^{(3)} > 0 \right\}}_{=1}$$

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

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

(derived earlier so copied here!)

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

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

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

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

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

the $\prod_{m,k} \left\{ \sum_{m,k}^{(3)} > 0 \right\}$ term comes

from differentiating ReLU.

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

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

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

So, use quotient rule on the following:

$$\hat{p}_{m,k} = \frac{u}{v} = \frac{\exp \{ p_{m,k}^{(3)} \}}{\sum_{z=0}^9 \exp \{ p_{z,k}^{(3)} \}}$$

u is given here

$$v' = \sum_{z=0}^9 \exp \{ p_{z,k}^{(3)} \} \prod \left\{ \sum_{z,k}^{(3)} > 0 \right\} \left(\sum_{i=0}^{63} w_{i,z} \prod \left\{ \sum_{i,k}^{(2)} > 0 \right\} \right)$$

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

$$= \frac{1}{\left(\sum_{z=0}^9 \exp \{ p_{z,k}^{(3)} \} \right)^2} \left(\left(\sum_{z=0}^9 \exp \{ p_{z,k}^{(3)} \} \right) \exp \{ p_{m,k}^{(3)} \} \prod \left\{ \sum_{m,k}^{(3)} > 0 \right\} \left(\sum_{i=0}^{63} w_{i,m} \prod \left\{ \sum_{i,k}^{(2)} > 0 \right\} \right) \right) \\ - \left(\exp \{ p_{m,k}^{(3)} \} \sum_{z=0}^9 \exp \{ p_{z,k}^{(3)} \} \prod \left\{ \sum_{z,k}^{(3)} > 0 \right\} \left(\sum_{i=0}^{63} w_{i,z} \prod \left\{ \sum_{i,k}^{(2)} > 0 \right\} \right) \right)$$

$$= \hat{p}_{m,k} \left(\sum_{i=0}^{63} w_{i,m} \prod \left\{ \sum_{m,k}^{(3)} > 0, \sum_{i,k}^{(2)} > 0 \right\} - \sum_{z=0}^9 \hat{p}_{z,k} \sum_{i=0}^{63} w_{i,z} \prod \left\{ \sum_{z,k}^{(3)} > 0, \sum_{i,k}^{(2)} > 0 \right\} \right)$$

$$\Rightarrow \frac{\partial}{\partial b} [C(\theta)] = \frac{-1}{K} \sum_{k=0}^{K-1} \sum_{m=0}^9 t_{m,R_k} \left(\sum_{i=0}^{63} w_{i,m} \prod \left\{ \sum_{m,k}^{(3)} > 0, \sum_{i,k}^{(2)} > 0 \right\} - \sum_{z=0}^9 \hat{p}_{z,k} \sum_{i=0}^{63} w_{i,z} \prod \left\{ \sum_{z,k}^{(3)} > 0, \sum_{i,k}^{(2)} > 0 \right\} \right)$$

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

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

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

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

Continued next page...

$b^{(1)}$ cont...

Keep unpacking the derivative from prev. page...

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

$$= \prod_{\{ \sum_{i,k}^{(2)} > 0 \}} \left(\sum_{j=0}^{127} w_{j,i}^{(2)} \frac{\partial}{\partial b^{(1)}} [p_{j,k}^{(1)}] \right) \rightarrow y = \sum_{i,k}^{(1)}$$

$$= \prod_{\{ \sum_{i,k}^{(2)} > 0 \}} \left(\sum_{j=0}^{127} w_{j,i}^{(2)} \frac{\partial}{\partial y} [f(y)] \frac{\partial}{\partial b^{(1)}} [b^{(1)} + \sum_{n=0}^{783} w_{n,j}^{(1)} p_{n,k}^{(0)}] \right)$$

$$= \prod_{\{ \sum_{i,k}^{(2)} > 0 \}} \left(\sum_{j=0}^{127} w_{j,i}^{(2)} \prod_{\{ \sum_{j,k}^{(1)} > 0 \}} \right) = 1$$

So, $\frac{\partial}{\partial b^{(1)}} [\exp\{p_{m,k}^{(3)}\}] = \exp\{p_{m,k}^{(3)}\} \prod_{\{ \sum_{z,k}^{(3)} > 0 \}} \left(\sum_{i=0}^{63} w_{i,m}^{(3)} \prod_{\{ \sum_{i,k}^{(2)} > 0 \}} \left(\sum_{j=0}^{127} w_{j,i}^{(2)} \prod_{\{ \sum_{j,k}^{(1)} > 0 \}} \right) \right)$

So use quotient rule on following:

$$\hat{p}_{m,k} = \frac{u}{v} = \frac{\exp\{p_{m,k}^{(3)}\}}{\sum_{z=0}^9 \exp\{p_{z,k}^{(3)}\}}$$

u' is given here

~~$\cancel{u} = \exp\{p_{m,k}^{(3)}\}$~~

$$v' = \sum_{z=0}^9 \exp\{p_{z,k}^{(3)}\} \prod_{\{ \sum_{z,k}^{(3)} > 0 \}} \left(\sum_{i=0}^{63} w_{i,z}^{(3)} \prod_{\{ \sum_{i,k}^{(2)} > 0 \}} \left(\sum_{j=0}^{127} w_{j,i}^{(2)} \prod_{\{ \sum_{j,k}^{(1)} > 0 \}} \right) \right)$$

So, $\frac{\partial}{\partial b^{(1)}} [\hat{p}_{m,k}] = \frac{vu' - uv'}{v^2}$

$$= \frac{1}{\left(\sum_{z=0}^9 \exp\{p_{z,k}^{(3)}\} \right)^2} \left(\left(\sum_{z=0}^9 \exp\{p_{z,k}^{(3)}\} \right) \exp\{p_{m,k}^{(3)}\} \prod_{\{ \sum_{m,k}^{(3)} > 0 \}} \left(\sum_{i=0}^{63} w_{i,m}^{(3)} \prod_{\{ \sum_{i,k}^{(2)} > 0 \}} \left(\sum_{j=0}^{127} w_{j,i}^{(2)} \prod_{\{ \sum_{j,k}^{(1)} > 0 \}} \right) \right) \right)$$

$$= \hat{p}_{m,k} \left(\sum_{i=0}^{63} \sum_{j=0}^{127} \prod_{\{ \sum_{m,k}^{(3)} > 0 \}} \left(\sum_{z=0}^9 \sum_{j=0}^{127} w_{j,i}^{(2)} \prod_{\{ \sum_{z,k}^{(3)} > 0 \}} \left(\sum_{i=0}^{63} w_{i,z}^{(3)} \prod_{\{ \sum_{i,k}^{(2)} > 0 \}} \left(\sum_{j=0}^{127} w_{j,i}^{(2)} \prod_{\{ \sum_{j,k}^{(1)} > 0 \}} \right) \right) \right) \right)$$

Define shorthand $I_{i,k}^{(1)} := \prod_{\{ \sum_{j,k}^{(1)} > 0 \}}$

$$= \hat{p}_{m,k} \left(\sum_{i=0}^{63} \sum_{j=0}^{127} I_{i,k}^{(3)} m_{j,k} w_{i,m}^{(3)} I_{i,k}^{(2)} w_{i,j}^{(2)} I_{i,k}^{(1)} - \sum_{i=0}^{63} \sum_{j=0}^{127} \hat{p}_{z,k} I_{z,k}^{(3)} w_{i,z}^{(3)} I_{z,k}^{(2)} w_{i,j}^{(2)} I_{i,k}^{(1)} \right)$$

$$= \hat{p}_{m,k} \left(\sum_{i=0}^{63} \sum_{j=0}^{127} I_{i,k}^{(3)} m_{j,k} w_{i,m}^{(3)} I_{i,k}^{(2)} w_{i,j}^{(2)} I_{i,k}^{(1)} - \sum_{i=0}^{63} \sum_{j=0}^{127} \sum_{z=0}^9 \hat{p}_{z,k} I_{z,k}^{(3)} w_{i,z}^{(3)} I_{z,k}^{(2)} w_{i,j}^{(2)} I_{i,k}^{(1)} \right)$$

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

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

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

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

So we 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:

$$U' = \exp \left\{ \hat{P}_{B,k}^{(3)} \right\} I_{B,k}^{(3)} P_{\alpha,k}^{(2)} \left[\sum_{m=B}^q \hat{P}_{m,k}^{(3)} \right]$$

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

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

$$= \frac{1}{\left(\sum_{z=0}^q \exp \left\{ \hat{P}_{z,k}^{(3)} \right\} \right)^2} \left(\left(\sum_{z=0}^q \exp \left\{ \hat{P}_{z,k}^{(3)} \right\} \right) \exp \left\{ \hat{P}_{B,k}^{(3)} \right\} I_{B,k}^{(3)} P_{\alpha,k}^{(2)} \left[\sum_{m=B}^q \hat{P}_{m,k}^{(3)} \right] \right.$$

$$\left. - \exp \left\{ \hat{P}_{m,k}^{(3)} \right\} \left(\exp \left\{ \hat{P}_{B,k}^{(3)} \right\} I_{B,k}^{(3)} P_{\alpha,k}^{(2)} \right) \right)$$

$$= \hat{P}_{m,k} \left(I_{\sum_{m=B}^q \hat{P}_{m,k}^{(3)}} I_{B,k}^{(3)} P_{\alpha,k}^{(2)} - \hat{P}_{B,k}^{(3)} I_{B,k}^{(3)} P_{\alpha,k}^{(2)} \right) \quad \text{can take:} \\ \left(I_{\sum_{m=B}^q \hat{P}_{m,k}^{(3)}} - \hat{P}_{B,k}^{(3)} \right) I_{B,k}^{(3)} P_{\alpha,k}^{(2)}$$

Therefore we arrive at:

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

Still Guangzhou \rightarrow London ii

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

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

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

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

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

$$= I_{i,k}^{(2)} \prod_{\{i=B_3^2\}} p_{\alpha,k}^{(1)}$$

$$(= I_{B,k}^{(2)} \prod_{\{i=B_3^2\}} p_{\alpha,k}^{(1)})$$

$$\Rightarrow \frac{\partial}{\partial W_{\alpha, \beta}^{(2)}} [\exp \{ p_{m,k}^{(3)} \}] = \exp \{ p_{m,k}^{(3)} \} \underbrace{I_{m,k} \sum_{i=0}^{63} w_{i,m}^{(3)} I_{B,k}^{(2)} \prod_{\{i=B_3^2\}} p_{\alpha,k}^{(1)}}$$

Interestingly, there is no index indicator here unlike

$$= \exp \{ p_{m,k}^{(3)} \} I_{m,k} w_{B,m} I_{B,k} p_{\alpha,k}$$

the $W_{\alpha, \beta}^{(3)}$ derivative. I think this is because the final layer 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 $\hat{P}_{m,k} = \frac{\exp \{ p_{m,k}^{(3)} \}}{\sum_{z=0}^9 \exp \{ p_{z,k}^{(3)} \}} = \frac{U}{V}$ U' is above.

$$V' = \sum_{z=0}^9 \exp \{ p_{z,k}^{(3)} \} I_{z,k}^{(3)} \underbrace{w_{B,z} I_{B,k}^{(2)} \prod_{\{i=B_3^2\}} p_{\alpha,k}^{(1)}}$$

$$\frac{\partial}{\partial W_{\alpha, \beta}^{(2)}} [\hat{P}_{m,k}] = \frac{vU' - uV'}{\sqrt{2}}$$

$$= \frac{1}{(\sum_{z=0}^9 \exp \{ p_{z,k}^{(3)} \})^2} \left(\begin{array}{l} \left(\sum_{z=0}^9 \exp \{ p_{z,k}^{(3)} \} \right) \exp \{ p_{m,k}^{(3)} \} I_{m,k}^{(3)} w_{B,m} I_{B,k}^{(2)} \prod_{\{i=B_3^2\}} p_{\alpha,k}^{(1)} \\ - \exp \{ p_{m,k}^{(3)} \} \sum_{z=0}^9 \exp \{ p_{z,k}^{(3)} \} I_{z,k}^{(3)} w_{B,z} I_{B,k}^{(2)} \prod_{\{i=B_3^2\}} p_{\alpha,k}^{(1)} \end{array} \right) = \dots$$

$$= \hat{P}_{m,k} \left(I_{m,k}^{(3)} W_{B,m}^{(3)} I_{B,k}^{(2)} P_{d,k}^{(1)} - \sum_{z=0}^q \hat{P}_{z,k} W_{B,z}^{(3)} I_{z,k}^{(3)} I_{B,k}^{(2)} P_{d,k}^{(1)} \right)$$

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

$$W_{\alpha,\beta}^{(1)}$$

$$\text{Start with } \frac{\partial}{\partial W_{\alpha,\beta}^{(1)}} [\exp \{ \hat{S}^{(3)} \hat{P}_{m,k} \}] = \dots$$

$$\dots = \exp \{ \hat{P}_{m,k}^{(3)} \} I_{m,k}^{(3)} \sum_{i=0}^{63} W_{i,m} \frac{\partial}{\partial W_{\alpha,\beta}^{(1)}} [P_{i,k}]$$

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

$$= I_{i,k}^{(2)} \sum_{j=0}^{127} W_{j,i} \frac{\partial}{\partial W_{\alpha,\beta}^{(1)}} [P_{j,k}^{(1)}]$$

$$= I_{i,k}^{(2)} \sum_{j=0}^{127} W_{j,i} I_{j,k}^{(1)} \frac{\partial}{\partial W_{\alpha,\beta}^{(1)}} [b^{(1)} + \sum_{n=0}^{783} W_{n,j} P_{n,k}^{(0)}]$$

$$= I_{i,k}^{(2)} \sum_{j=0}^{127} W_{j,i} I_{j,k}^{(1)} \underbrace{\sum_{\{j=\beta\}}}_{\text{since } j=\beta} P_{\alpha,k}^{(0)}$$

$$= I_{i,k}^{(2)} W_{B,i} I_{B,k}^{(1)} P_{d,k}^{(0)}$$

(since summing over j zeroes all except $j=\beta$)

So, subbing back in:

$$\frac{\partial}{\partial W_{\alpha,\beta}^{(1)}} [\exp \{ \hat{S}^{(3)} \hat{P}_{m,k} \}] = \exp \{ \hat{P}_{m,k}^{(3)} \} I_{m,k}^{(3)} \sum_{i=0}^{63} W_{i,m} I_{i,k}^{(2)} W_{B,i} I_{B,k}^{(1)} P_{d,k}^{(0)}$$

$$\text{(QUOTIENT RULE)} \quad \hat{P}_{m,k} = \frac{u}{v} = \frac{\exp \{ \hat{P}_{m,k}^{(3)} \}}{\sum_{z=0}^q \exp \{ \hat{P}_{z,k}^{(3)} \}}. \quad (1) \quad v' = \sum_{z=0}^q \exp \{ \hat{P}_{z,k}^{(3)} \} \sum_{i=0}^{63} W_{i,z} I_{i,k}^{(2)} W_{B,i} I_{B,k}^{(1)} P_{d,k}^{(0)}$$

$$\frac{\partial}{\partial W_{\alpha,\beta}^{(1)}} [\hat{P}_{m,k}] = \frac{v u' - u v'}{v^2} = \hat{P}_{m,k} \left(I_{m,k}^{(3)} \sum_{i=0}^{63} W_{i,m} I_{i,k}^{(2)} W_{B,i} I_{B,k}^{(1)} P_{d,k}^{(0)} \right.$$

$$\left. - \sum_{z=0}^q \hat{P}_{z,k} I_{z,k}^{(3)} \sum_{i=0}^{63} W_{i,z} I_{i,k}^{(2)} W_{B,i} I_{B,k}^{(1)} P_{d,k}^{(0)} \right)$$

$$\Rightarrow \frac{\partial}{\partial W_{\alpha,\beta}^{(1)}} [C(\underline{\theta})] = \frac{-1}{K} \sum_{k=0}^{K-1} \sum_{m=0}^q t_{m,k}^{(1)} I_{B,k}^{(0)} \left(I_{m,k}^{(3)} \sum_{i=0}^{63} W_{i,m} I_{i,k}^{(2)} W_{B,i}^{(2)} \right.$$

$$\left. - \sum_{z=0}^q \hat{P}_{z,k} I_{z,k}^{(3)} \sum_{i=0}^{63} W_{i,z} I_{i,k}^{(2)} W_{B,i}^{(2)} \right).$$