

Backprop on Fully Connected Feedforward Network 1 Hidden Layer

By the end of a forward pass:

$$Z_1 = XW_1 + b_1 \rightarrow A_1 = \text{ReLU}(Z_1)$$

$$Z_2 = A_1W_2 + b_2 \rightarrow A_2 = \text{Softmax}(Z_2)$$

$$\text{loss} = -\frac{1}{N} \sum_i \sum_j y_{ij} \log(A_{2ij})$$

How does changing the weights and biases affect the loss?

We want gradient of loss wrt weights and biases, $\frac{\partial L}{\partial W}, \frac{\partial L}{\partial b}$

We reach through chain rule of calculus -

$$\text{Loss per sample: } L = - \sum_j y_j \log(A_{2j})$$

remember \log amplifies differences near zero and cuts on the probability like results of softmax - multiplying by y_j uses one-hot encoding to select the probability assigned to the 'right' class and zero out the others. We then negate b/c $\log(x)$, $x \in (0, 1)$ will be negative.

We do the derivatives backwards based on dependencies

loss depends on A_2

A_2 depends on $Z_2 = A_1W_2 + b_2$

Z_2 depends on A_1

A_1 depends on $Z_1 = XW_1 + b_1$

Z_1 depends on W_1, b_1 & input x

It's easiest to see how the output A_2 changes the loss, and then we propagate that change backward thru each layer using chain rule - avoids recomputing derivatives repeatedly

$$\text{First, do } \frac{\partial L}{\partial A_{2j}} \text{ for one sample: } \frac{\partial L}{\partial A_{2j}} = \frac{\partial}{\partial A_{2j}} (-y_j \log(A_{2j}))$$

$$\frac{\partial L}{\partial A_{2j}} = -\frac{y_j}{A_{2j}} \quad \text{note: } \log \text{ here is } \ln, \frac{\partial}{\partial x}(\ln x) = \frac{1}{x}$$

Next we know $\frac{\partial L}{\partial Z_2} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial Z_2}$ from the chain rule

$$Z_2 = A_1W_2 + b_2, \quad L = -y_j \log(A_{2j})$$

Softmax: $A_{2j} = \frac{e^{Z_{2j}}}{\sum_k e^{Z_{2k}}}$ where the denominator is normalizing

$$\frac{\partial A_{2j}}{\partial Z_{2k}} = \begin{cases} A_{2j}(1-A_{2j}) & j=k \\ A_{2j}A_{2k} & j \neq k \end{cases} \quad \text{how changing one } Z \text{ affects all probabilities}$$

$$\rightarrow \frac{\partial L}{\partial Z_{2k}} = \sum_{j=1}^C \frac{\partial L}{\partial A_{2j}} \frac{\partial A_{2j}}{\partial Z_{2k}} \quad \text{the chain rule allows sum across all paths}$$

If you do this algebra, all the cross terms cancel ($i \neq j$) and you get

$$\frac{\partial L}{\partial Z_{2k}} = A_{2k} - y_{2k}$$

and that's why cross entropy loss + softmax is a popular combo.

Mental Model:

- $Z_2 \rightarrow \text{softmax} \rightarrow A_2 \rightarrow \text{loss}$ as flowchart
- each knob Z_{2k} affects all A outputs, each A output affects L
- arrows from all Z_{2k} to all A_j , each weighted by $\frac{\partial A_j}{\partial Z_{2k}}$
- arrows from all A_j to L , weighted by $\frac{\partial L}{\partial A_j}$
- chain rule = multiply along each arrow, sum contributions into $\frac{\partial L}{\partial Z_{2k}}$

Next in line is $\frac{\partial L}{\partial W_2}$ & $\frac{\partial L}{\partial b_2}$:

$$\text{by chain rule: } \frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial Z_2} \cdot \frac{\partial Z_2}{\partial W_2} \quad Z_2 = \sum_i A_{1i} W_{2i}^{(k,j)} + b_2^{(j)}$$

$$= (A_{2k} - y_{2k}) A_1 \quad \frac{\partial L}{\partial W_2} = A_1^T \cdot \frac{\partial L}{\partial Z_2}$$

and to make that multiplication work in the matrix world, we transpose A_1 , so the inner dims of A_1 & $\frac{\partial L}{\partial Z_2}$ match

$$\text{by chain rule: } \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial Z_2} \cdot \frac{\partial Z_2}{\partial b_2} \quad \& \quad \frac{\partial Z_2}{\partial b_2} = 1$$

$$\rightarrow \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial Z_2} \quad \text{or} \quad \frac{\partial L}{\partial b_2} = \sum_{\text{samples}} \frac{\partial L}{\partial Z_2}$$

$$\text{Next we want } \frac{\partial L}{\partial Z_1} = \frac{\partial L}{\partial A_1} \cdot \frac{\partial A_1}{\partial Z_1} = \left(\frac{\partial L}{\partial Z_2} \cdot \frac{\partial Z_2}{\partial A_1} \right) \frac{\partial A_1}{\partial Z_1}$$

where $\frac{\partial Z_2}{\partial A_1} = W_2$, how output pre-activ depends on hidden layer activation

$\frac{\partial A_1}{\partial Z_1}$ = derivative of activation function

$$\text{thus } \frac{\partial L}{\partial Z_1} = \frac{\partial L}{\partial Z_2} \cdot W_2^T \odot \sigma'(Z_1) \quad \text{where } \odot \text{ is element wise mult \& } \sigma'(Z_1) \text{ is deriv of ReLU}$$

$$\text{then by chain rule } \frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial Z_1} \cdot \frac{\partial Z_1}{\partial W_1} \quad \text{but } Z_1 = XW_1 + b_1, \quad \frac{\partial Z_1}{\partial W_1} = X$$

$$\text{so } \frac{\partial L}{\partial W_1} = X^T \cdot \frac{\partial L}{\partial Z_1}$$

$$\text{and } \frac{\partial L}{\partial b_1} = \sum_{\text{samples}} \frac{\partial L}{\partial Z_{1k}}$$