

Grad Problem Solving

Submitted By: Nafis Neehal

661990881

Problem 1:

- Image size = $3 \times N \times N$ (RGB)
- L layers with h node/hidden layer, n nodes in the output layer [typically, input layers are not considered as layers as they don't have any learnable parameters to optimize like weights or bias. But, I'm considering input layer as a separate layer here for mathematical convenience].
- Input \rightarrow Hidden = $[(3 \times N \times N + 1) \times h]$ number of learnable parameters, $3 \times N \times N \times h$ weights and h biases
- For (L-2) number of Hidden Layers = $[(L-2) \times \{(h \times h) + h\}]$ number of total learnable parameters, where each hidden layer has $h \times h$ number of weights as input and h number of biases.
- Hidden^[L-1] \rightarrow Output = $[(h \times n) + n]$ number of total learnable parameters.
- Final Formula:
 $[(3 \times N \times N + 1) \times h] + [(L-2) \times \{(h \times h) + h\}] + [(h \times n) + n]$

Problem 2:

- $3 \times N \times N$ at input
- Conv filters, $k \times k$ size, number of channels = number of filters = d
- Padding = 1, Stride = 1 will keep the image dimension same:
 - **Proof:**
Let image size $n \times n = 6 \times 6$, filter size = 3×3 , padding
Output image dimension = $\left(\frac{n+2p-f}{s} + 1\right) \times \left(\frac{n+2p-f}{s} + 1\right) = 6 \times 6$
- L #of conv layers, F #of FC layers, $n^{(o)}$ #of nodes in output layer
- Theory \rightarrow Conv layers will only change depth if number of filter increases from d, image size will remain $N \times N$ as padding=1 is set.
- Theory \rightarrow Max-pool layers will pull down the image size (height and width only, not depth) into half of the input (as filter size 2×2 , so stride will be 2 as non-overlapping regions).
- Each conv is followed by max-pool! So L number of Convs and L number of Max-pools.
- Filters in Conv layers can only be considered as learnable parameters, Max-Pool doesn't have any of that, it just pulls out the max value of the selected image region.
- As, in each layer image size will reduce by a factor of 2,
 - At i^{th} layer, image dimension will be $\frac{dxN \times N}{2^{2i}}$
 - Number of learnable parameters in each CONV layer will be $[dx3 \times 3 + d] = [d \times \{(3 \times 3) + 1\}]$, when number of bias parameters is d (1 bias for each filter).
- So, to generalize, if the CONV filter size is $f \times f$, number of filters is d, and there are L #of CONV filters, then total learnable all CONV parameters are =

$$L \times [d \times \{(f \times f) + 1\}]$$

- For FC layers, number of learnable paramers =
Number of activations in previous layer x number of activations in current layer + 1
- So if we have F number of FC layers, and h nodes in each FC layers and $n^{(o)}$ number of nodes in the output layer, then:
 - FC->FC = $(F-2) \times h$
 - Conv-last-layer->FC-1st-layer = $[d \times \{(fxf)+1\}] \times h$
 - Total for all FC layers: $[(F-2) \times h] + [[d \times \{(fxf)+1\}] \times h]$
- Total number of learnable parameters in output layer:
 - FC-Last-Layer->Output = $n^{(o)} \times h$

Problem 3(a):

The proof is a little bit unorthodox and uses inference rather than a solid mathematical proof. But the logic is perfect.

Let's say-

Output at final layer = [1, 3.5, 757, -5]

#this output vector covers normal small number, decimal, large number and negative number

For softmax, for each output node Z_i ,

Final output will be = $\text{Softmax}(Z_i) = \frac{e^{Z_i}}{\sum_{k=1}^k e^{Z_k}}$

Now, $e^1 = 2.718$, $e^{3.5} = 33.1545$, $e^{757} = 5.76$, $e^{-5} = 0.007$

So, while applying Softmax to each of the output node, denominator always will be =

$2.718 + 33.1545 + 5.76 + 0.001 = 41.639$

And for each softmax, one of the 4 exponential values will be as numerator.

So, $\frac{e^{Z_i}}{\sum_{k=1}^k e^{Z_k}}$ can never be greater than 1 and can be at **max very close to 1** if one of the four values is actually equal to denominator.

Again, as it is an exponential function, whatever the input is, the output will always have to be **greater than 0**. So, the value of softmax is always lying in between 0 and 1 which represents probability distribution.

Problem 3(b):

$$\text{Let } \sum_{i=1}^k e^{Z_k} = \text{SUM}$$

$$P_j = \frac{e^{Z_j}}{\sum_{i=1}^k e^{Z_k}} = \frac{e^{Z_j}}{\text{SUM}}$$

$$\frac{\partial P_i}{\partial Z_k} = \frac{\partial}{\partial Z_k} \left(\frac{e^{Z_i}}{\text{SUM}} \right)$$

$$\text{By applying, } f'(U/V) = \frac{V f'(U) - U f'(V)}{V^2},$$

$$\frac{\partial P_i}{\partial Z_k} = \frac{e^{Z_i} \cdot \text{SUM} - e^{Z_i} e^{Z_k}}{\text{SUM}^2} \quad [\text{for } i=j \text{ case}]$$

$$= \frac{e^{Z_i} (\text{SUM} - e^{Z_k})}{\text{SUM} \cdot \text{SUM}} = P_i * (1 - P_i)$$

$$\frac{\partial P_i}{\partial Z_k} = \frac{e^{Z_i} \cdot \text{SUM} - e^{Z_i} e^{Z_k}}{\text{SUM}^2} \quad [\text{for } i \neq j \text{ case}]$$

$$= \frac{0 - e^{Z_i} e^{Z_k}}{\text{SUM}^2} \quad [\text{for } i \neq j \text{ case}]$$

$$= - P_i * P_k$$

$$\text{We know, Kronecker delta function, } \delta_{ij} = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}$$

Combining both ($i=j$) and ($i \neq j$) cases,

$$\frac{\partial P_i}{\partial Z_k} = \begin{cases} P_i * (1 - P_i), & i = j \\ - P_i * P_k, & i \neq j \end{cases}$$

$$= P_i (\delta_{ik} - P_k)$$

Problem 3(C)

Cross Entropy Loss FN: $L(p,y) = -\sum_i Y_i * \log(P_i) = -Y_1 * \log(P_1) - Y_2 * \log(P_2) - \dots - Y_k * \log(P_k)$

Net Output at output layer = Z_i

$P_i = \text{Softmax}(Z_i)$

$$\frac{\partial L}{\partial Z_i} = \frac{\partial P_i}{\partial Z_i} X \frac{\partial L}{\partial P_i} \text{ [Applied chain rule]}$$

$$\frac{\partial L}{\partial P_i} = -\frac{Y_i}{P_i} \text{ for } \forall i = 1 \dots k$$

$$\frac{\partial P_i}{\partial Z_k} = \begin{cases} P_i * (1 - P_i), & i = j \\ -P_i * P_k, & i \neq j \end{cases} \text{ [from 3(b)]}$$

Putting it together,

$$\begin{aligned} \frac{\partial P_i}{\partial Z_i} X \frac{\partial L}{\partial P_i} &= \begin{bmatrix} P_1(1 - P_1) & \dots & -P_1 P_k \\ \vdots & \ddots & \vdots \\ -P_1 P_k & \dots & P_k(1 - P_k) \end{bmatrix} \begin{bmatrix} -\frac{Y_1}{P_1} \\ \dots \\ -\frac{Y_k}{P_k} \end{bmatrix} \\ &= \begin{bmatrix} -Y_1 + Y_1.P_1 + \sum_{i \neq 1}^k Y_i P_1 \\ \dots \\ -Y_k + Y_k.P_k + \sum_{i \neq 1}^k Y_i.P_k \end{bmatrix} \\ &= \begin{bmatrix} -Y_1 + P_1 \sum_{i=1}^k Y_i \\ \dots \\ -Y_k + P_k \sum_{i=1}^k Y_i \end{bmatrix} \\ &= \begin{bmatrix} -Y_1 + P_1 \\ \dots \\ -Y_k + P_k \end{bmatrix}, \text{ since } \sum_{i=1}^k Y_i = 1 \text{ [softmax is probability distribution]} \\ &= P_i - Y_i \text{ for } \forall i = 1 \dots k \end{aligned}$$