

## Assignment-03

K. Surya Prakash  
E-E18 BTech 11026

Q1)  
=

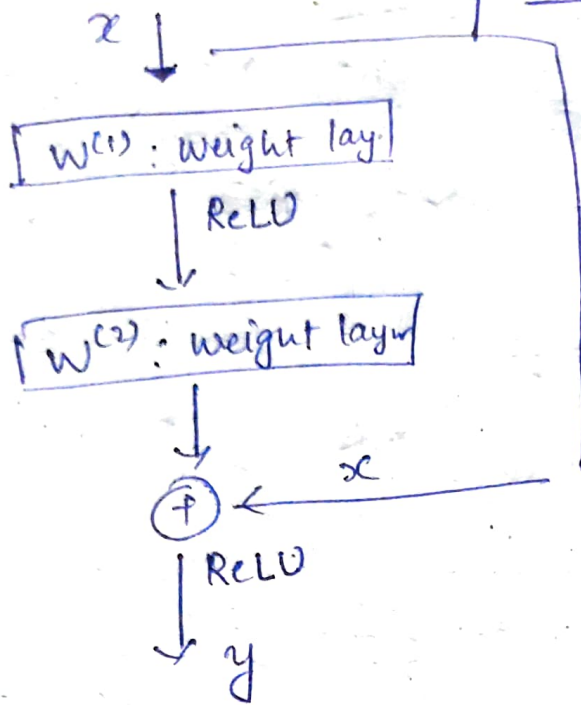
$h^{(1)}$

$z^{(1)}$

$h^{(2)}$

$h^{(3)}$

$y$



Forward pass:

o/p 1st blk:

$$h^{(1)} = x * w^{(1)}$$

$$z^{(1)} = \text{ReLU}(h^{(1)})$$

o/p 2nd blk:

$$h^{(2)} = z^{(1)} * w^{(2)}$$

$$h^{(3)} = h^{(2)} + x$$

$$y = \text{ReLU}(h^{(3)})$$

Note:  $A \odot B \Rightarrow$  Element wise product  
which result in same dim

$A * B \Rightarrow$  Conv. to give same dim.  
output

Assumptions:

1)  $x, h^{(1)}, h^{(2)}, h^{(3)}$  all have the same dimensions

$\Rightarrow$  The convolutions are done st, the sufficient padding is to be done.

$\Rightarrow$  Assuming  $\frac{\partial L}{\partial y}$  is already known

\* Backpropagation:

Need to find  $\frac{\partial L}{\partial w^{(1)}}, \frac{\partial L}{\partial w^{(2)}}, \frac{\partial L}{\partial x}$

$$1) \frac{\partial L}{\partial w^{(2)}} = \frac{\partial L}{\partial y} (h^{(3)} > 0) \underbrace{(1)}_1 \cdot \frac{\partial h^{(3)}}{\partial h^{(2)}} \cdot \frac{\partial h^{(2)}}{\partial w^{(2)}}$$

$$\frac{\partial L}{\partial w^{(2)}} = \left( \frac{\partial L}{\partial y} (h^{(3)} > 0) \right) \otimes 1 * (z^{(1)})$$

$$\boxed{\frac{\partial L}{\partial w^{(2)}} = \left( \frac{\partial L}{\partial y} (h^{(3)} > 0) \right) * z^{(1)}}$$

( $\otimes$ ) element wise product

$$\frac{\partial L}{\partial w^{(1)}} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial h^{(2)}} \cdot \frac{\partial h^{(2)}}{\partial z^{(1)}} \cdot (h^{(2)} > 0) \cdot \frac{\partial h^{(1)}}{\partial w^{(1)}}$$

$$\frac{\partial L}{\partial h^{(1)}} \neq X \rightarrow (1)$$

$$\frac{\partial L}{\partial w^{(1)}} = \left( \frac{\partial L}{\partial y} (h^{(2)} > 0) \neq (w^{(2)}_{\text{flipped}}) (0) \right)$$

$$(h^{(1)} > 0) \neq X$$

$$\frac{\partial L}{\partial w^{(1)}} \left\{ \left( \left( \frac{\partial L}{\partial y} (h^{(2)} > 0) \neq w^{(2)}_{\text{flip}} \right) (0) (h^{(1)} > 0) \right) \right\}$$

~~X~~

$$\frac{\partial L}{\partial w^{(1)}} = \frac{\partial L}{\partial h^{(1)}} \neq X$$

} Short form

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y}^{(1)} \frac{\partial y}{\partial h^{(3)}}^{(1)} \frac{\partial h^{(3)}}{\partial x}$$

$$= \frac{\partial L}{\partial y}^{(1)} (h^{(3)} > 0)^{(1)} \frac{\partial h^{(3)}}{\partial x} \rightarrow (ii)$$

\* Computing  $\frac{\partial h^{(3)}}{\partial x} = \frac{\partial h^{(2)}}{\partial x} + 1$  represent  
all ones  
matrix

$$\Rightarrow \frac{\partial h^{(2)}}{\partial x} = \left( * w_{flip}^{(2)} \right)^{(1)} (h^{(1)} > 0) \left( * w^{(1)} \right)$$

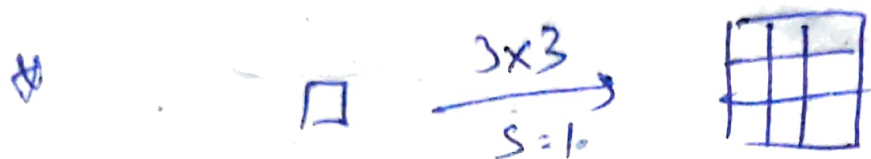
$$\Rightarrow \frac{\partial h^{(2)}}{\partial x} = \left( \left( * w_{flip}^{(2)} \right)^{(1)} (h^{(1)} > 0) \right) * w^{(1)}$$

Putting all together:

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial h^3} \left( 1 + \frac{\partial h^{(2)}}{\partial x} \right)$$

$$= \frac{\partial L}{\partial h^3} + \left( \frac{\partial L}{\partial h^3} * w_{flip}^{(2)} \right)^{(1)} (h^{(1)} > 0) * \underline{\underline{w^{(1)}}}$$

Q2) 4 consecutive  $3 \times 3$  conv layers  
→ stride: 1



For every layer each dimension of the support increases by '2'.

$$\Rightarrow 1 \times 1 \rightarrow 3 \times 3 \rightarrow 5 \times 5 \rightarrow 7 \times 7 \rightarrow 9 \times 9$$

$C_1 \qquad C_2 \qquad C_3 \qquad C_4$

∴ A single pixel in the image contributes

to a  $9 \times 9$  patch in the 4th  
non-image layer

$$\Rightarrow \text{Total effected pixels} = 81 //$$



Q3)

→ If no. of hidden units  $\rightarrow \uparrow$

⇒ Model can learn more complex

relations b/w O/P and I/P

⇒ But with more complexity

model has: less bias (Soln is closer to optimal solution)

: high variance

As model now learns complex relations, noise in the dataset is also captured, this results in fluctuations in the result

Model will be more sensitive

more hidden units }

less bias

high variance

$$Q4) y_k(x/w) = \sum_j w_{kj}^{(2)} \sigma \left( \sum_i w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)}$$

~~But~~ we will be using tanh activation instead of sigmoid. & we desire to get same output :  $y_k(x/w)$

\* Definition:

$$w_{kj}^{(2)} ; w_{ji}^{(1)} ; w_{k0}^{(2)} ; w_{j0}^{(1)}$$

Weights <sup>with</sup> of sigmoid model

$$t_{kj}^{(2)} ; t_{ji}^{(1)} ; t_{k0}^{(2)} ; t_{j0}^{(1)}$$

weights with tanh model.

\* Relation b/w tanh & sigmoid

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \frac{2}{1 + e^{-2z}} - 1$$

$$= 2\sigma(2z) - 1$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Let  $z_j^0 = \sum_i w_{ji}^{(1)} x_i^0 + w_{j0}^{(1)}$  } with Sigmoid weights.

$a_j^0 = \sum_i t_{ji}^{(1)} x_i^0 + t_{j0}^{(1)}$  } with tanh weights.

\* For Sigmoid model:

$$y_k(x, w) = \sum w_{kj}^{(2)} \sigma(z_j^0) + w_{k0}^{(2)} \rightarrow (1)$$

$$y_k(x, t) = \sum t_{kj}^{(2)} \tanh(a_j^0) + t_{j0}^{(2)}$$

$$= \sum (2t_{kj}^{(2)} \cdot \sigma(2a_j^0) - 1) + t_{j0}^{(2)}$$

$$= \sum_j 2t_{kj}^{(2)} \sigma(2a_j^0) + t_{j0}^{(2)} - \sum_j 2t_{kj}^{(2)} \rightarrow (2)$$



Comparing (1) & (2)

$$2a_j^0 = z_j^0 \longrightarrow (1)$$

$$\Rightarrow 2 \sum t_{ji}^{(1)} x_i^0 + 2t_{j0}^{(1)} = 2 \sum w_{ji}^{(1)} x_i^0 + w_{j0}^{(1)}$$

By comparing:

$$t_{ji}^{(1)} = \frac{w_{ji}^{(1)}}{2}$$

$$t_{j0}^{(1)} = \frac{w_{j0}^{(1)}}{2}$$

Also:

$$t_{kj}^{(2)} = \frac{w_{kj}^{(2)}}{2}$$

$$t_{k0}^{(2)} - 2 \sum t_{kj}^{(2)} = w_{k0}^{(2)}$$

Linear transformations

∴ By making the above transformations  
to the existing weights we can  
obtain the same output even by  
using tanh activation. Hence she  
is correct.

Q5) Quadratic error:

$$E(w) \approx E(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*)$$

$H$ : Hessian evaluated @  $w^*$

\* We know that, if  $w^*$  is the optimal soln

$H$  is a PSD

$\Rightarrow$  Eigen values are:

$$H u_i^0 = \lambda_i^0 u_i^0 \quad (\lambda_i^0 \geq 0 \because H: \text{PSD})$$

where,

$u_i^0$ : Eigen vectors of  $H$ , st

$u_i^0$  are orthonormal

$$\Rightarrow \boxed{u_i^{0T} u_j^0 = \delta_{ij}} \rightarrow (1)$$

\* Let us express  $w - w^*$  as linear combination of eigen vectors

$$\boxed{w - w^* = \sum_i \alpha_i u_i^0} \rightarrow (2)$$

$$\therefore \frac{1}{2} (w - w^*)^T H (w - w^*) = \left(\frac{1}{2}\right) \left(\sum_i \alpha_i u_i^0\right)^T H \left(\sum_j \alpha_j u_j^0\right)$$

$$= \frac{1}{2} \left( \sum_i \alpha_i u_i \right)^T \cdot \left( \sum_j \lambda_j u_j \right)$$

$$= \frac{1}{2} \left( \sum_i \alpha_i u_i \right)^T \left( \sum_j \alpha_j \lambda_j u_j \right)$$

$$= \frac{1}{2} \sum_i \sum_j (\lambda_j) (\alpha_i \alpha_j u_i^T u_j)$$

\* We know that from (1)

$$u_i^T u_j = \delta_{ij} = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases}$$

$$= \frac{1}{2} \sum_i \lambda_i \alpha_i^2$$

$$\therefore \underbrace{E(w) - E(w^*)}_{\geq 0} \approx \frac{1}{2} \sum_i \lambda_i \alpha_i^2$$

$\Rightarrow$  Since 'H' is a PSD  $\Rightarrow E(w) \geq E(w^*)$

$$v^T H v \geq 0$$

$$\Rightarrow E(w) - E(w^*) \geq 0$$

$$\Rightarrow \boxed{\frac{1}{2} \sum_i \lambda_i \alpha_i^2 \geq 0}$$

⇒ This can be written in the form;

$$\sum \lambda_i d_i^2 = c^2 : c^2 = 2(E(w) - E(w^*))$$

\* for a constant error :  $E(w) = \text{constant}$

$$\Rightarrow c^2 = \text{constant}$$

∴  $w^*$  : fixed vector

$$\boxed{\sum \lambda_i d_i^2 = c^2} : (\text{positive constant})$$

↳ This represents elliptical contours

where axes are  $u_i$  (eigen vectors of  $H$ )

$$\Rightarrow w = w^* + \sum d_i u_i$$

$$\hookrightarrow w - w^* = \sum d_i u_i$$

~~A trans~~

Translation of axes from origin as center to  $w^*$  as center.

\* The length of axis:  $\alpha_j^0$  is computed by setting  $\alpha_i^0 |_{i \neq j} = 0$

$$\Rightarrow \boxed{\alpha_j^0 = \sqrt{\frac{c^2}{\lambda_j^0}}}$$

We know that  $\lambda_j^0$ : eigen value of  $H \geq 0$

$$\therefore \left( \alpha_j^0 \propto \frac{1}{\sqrt{\lambda_j^0}} \right) \quad \checkmark$$

\* length of axis is inversely proportional to square root of corresponding eigen value ( $\lambda_j^0$ )



Q6) We can use transfer learning?

Available dataset:

20 images of 200 classes

⇒ Comment: Small dataset

Available deployed model:

→ Trained on large dataset.

~~Dataset~~  
Trained data is similar to available data.

→ We can borrow the deployed model specifications (architecture, pre-trained weights) and retrain by making some tweaks to the model.

## General understanding

- 1) Freeze the 'Source task generic layers' & 'Source task special layers'.

⇒ These layers are responsible for feature extraction. Since both data are alike the features are likely to be same.

- 2) Discard the existing classification layer (which classifies 1000 classes).

- 3) Attach a new classification layer (which classifies 200 classes), with random initialise weights.

- 4) Train the network, with the available data, make sure weight updation should be done in the newly attached 'classification layer'

2 The final layer should have 200 <sup>neurons</sup> ~~classes~~ to classify 200 classes.