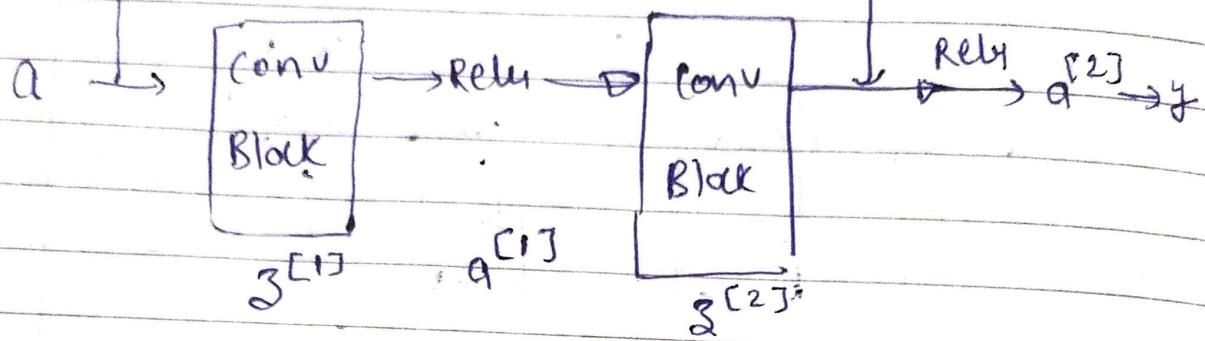


Ans-1

Identity



$a \rightarrow$  Last Layer activation as input to Residual Block.

$$z^{[1]} = w^{[1]} a + b^{[1]}$$

$$a^{[1]} = g(z^{[1]}) = \text{ReLU}(z^{[1]})$$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$$y = a^{[2]} = g(z^{[2]} + a) = \text{ReLU}(z^{[2]} + a)$$

$$\frac{\partial E}{\partial y} = \frac{\partial E}{\partial a^{[2]}} = \delta^{[2]} \rightarrow \text{derivative of Loss wrt Resblock output}$$

$$\frac{\partial E}{\partial z^{[2]}} = \frac{\partial E}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} = \delta^{[2]} \times g'(z^{[2]} + a)$$

Similarly on solving

$$\frac{\partial E}{\partial z^{[1]}} = \delta^{[2]} \cdot \delta^{[1]} \cdot g'(z^{[1]}) \left[ \delta^{[1]} = \frac{\partial E}{\partial a^{[1]}} \right]$$

Ans-2 →

Given

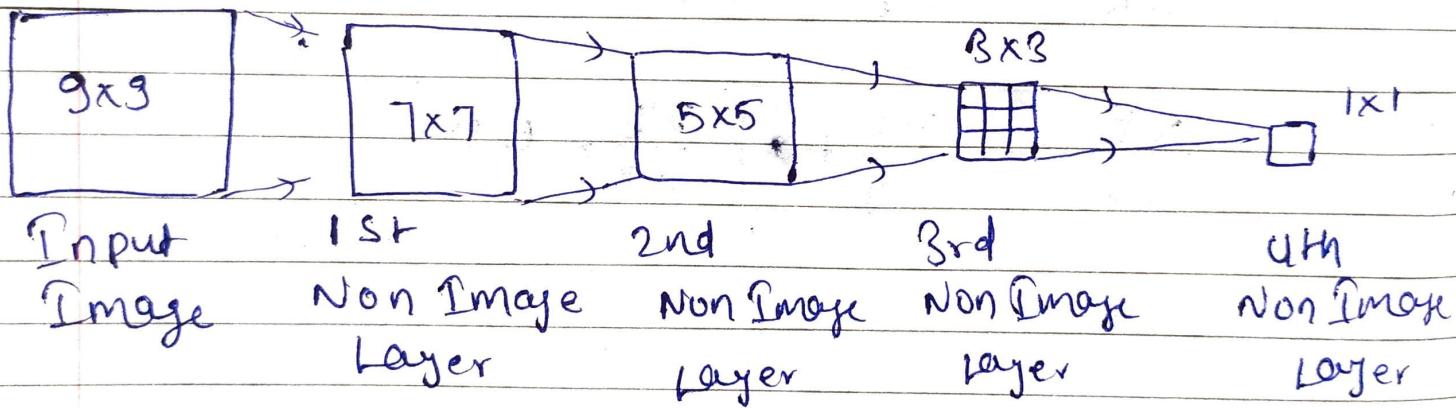
Filter size =  $3 \times 3$

Number of convolution layers = 4

stride = 1

Pooling = No

For every convolution layer the size of image layer reduces from  $n$  to  $n-2$ .



Hence Input image dimension of  $9 \times 9$  will contribute to one activation neuron in the 4th non image layer.

So support size for one activation  
 $9 \times 9 = 81$  pixels.

ANS-3 →

If the number of hidden units is increased  
Network will be able to learn more complex  
function representation.

And as we know more complicated models  
result in Lower bias & high variance.  
models with high complexity tends to  
overfit as they try to follow training data too  
closely & does not generalize well on testing  
Data.

Ans-4 → Given a two layer neural network

$$y_k(n, w) =$$

$$h \circ \left( \sum_{j=1}^m w_{kj}^{[2]} \sigma \left( \sum_{i=1}^n w_{ji}^{[1]} n_i + w_{j0} \right) + w_{k0}^{[2]} \right)$$

$h$  is a linear Activation function.

$\sigma(a)$  is sigmoid Activation function.

first let's prove that tanh is rescaled & shifted version of the sigmoid function.

Sigmoid Property : →

$$\left[ 1 - \frac{1}{1+e^{-n}} = \frac{1}{1+e^n} \Rightarrow 1 - \sigma(n) = \sigma(-n) \right]$$

$$\rightarrow \tanh(n) = \frac{e^n - e^{-n}}{e^n + e^{-n}}$$

$$= \frac{e^n + e^{-n} - 2e^{-n}}{e^n + e^{-n}}$$

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

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

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

$$= 1 - 2 + 2\sigma(2n) \quad [\text{Sigmoid Property}]$$

$$\tanh(n) = 2\sigma(2n) - 1 \quad \text{eq ①}$$

which prove that tanh is rescaled & shifted version of sigmoid function.

This shift of -1 will be taken care of by bias learning of hidden layer.

The inner  $\sigma$  in  $2\sigma(2n)-1$ , which scale input by a factor of 2, will be handled by weights present in first hidden layer.

& the outer scaling of 2 will be handled by weights in 2nd layer which contain linear activation function.

$$y_k(n, w) = \sum_{j=1}^M \frac{w_{kj}^{[2]}}{2} \cdot \sigma \left( \sum_{i=1}^D 2 \left[ \frac{w_{ji}^{[1]} n_i}{2} + \frac{w_{j0}^{[1]}}{2} \right] \right) + w_{k0}^{[2]}$$

$$y_k(n, w) = \sum_{j=1}^M \frac{w_{kj}^{[2]}}{2} \cdot \sigma \left( \sum_{i=1}^D 2 \left[ \frac{w_{ji}^{[1]} \cdot x_i}{2} + \frac{w_{j0}^{[1]}}{2} \right] \right) + w_{k0}^{[2]} - \frac{w_{kj}^{[2]} w_{j0}^{[1]}}{2}$$

$$y_k(n, w) =$$

$$\sum_{j=1}^M \frac{w_{kj}^{[2]}}{2} \left\{ \sigma \left( \sum_{i=1}^n \left[ \frac{w_{ji}^{[1]} \cdot x_i}{2} + \frac{w_{jo}^{[1]}}{2} \right] \right) - 1 \right\} + w_{ko}^{[2]} + \frac{w_{kj}^{[2]}}{2}$$

which is same as.

$$\rightarrow \sum_{j=1}^M \frac{w_{kj}^{[2]}}{2} \tanh \left[ \frac{w_{ji}^{[1]} \cdot x_i}{2} + \frac{w_{jo}^{[1]}}{2} \right] + w_{ko}^{[2]} + \frac{w_{kj}^{[2]}}{2}$$

→ By eq ①

$$\rightarrow \sum_{j=1}^M \tilde{w}_{kj}^{[2]} \tanh \left[ \tilde{w}_{ji} \cdot x_i + \tilde{w}_{jo} \right] + \tilde{w}_{ko}^{[2]}$$

so only weights will change for learning  
same function.

Ans-5→

## Quadratic Error

$$E(w) \approx E(w^*) + \frac{1}{2} (w-w^*)^T H(w-w^*) - 0$$

H= Hessian matrix evaluated at  $w^*$

Eigen value equation

$$\rightarrow H u_i = d_i u_i \quad -②$$

We have to show that contours of constant error are ellipses whose axes are aligned with eigen vectors  $u_i$  with the lengths that are inversely proportional to the square root of the corresponding eigen values  $d_i$ .

$$\text{Our hint} \rightarrow w-w^* = \sum_i \alpha_i u_i \quad -③$$

→ Eigen vectors for orthonormal set

$$u_i^T u_j = \delta_{ij} \quad -④$$

$$\text{putting } w-w^* = \sum_i \alpha_i u_i$$

into Quadratic Error we get →

$$E(w) = E[w^*] + \frac{1}{2} \sum_i \alpha_i^* \mu_i \cdot H(\{\alpha_i^* \mu_i\})$$

using eq 2.44 we get

$$E(w) = E[w^*] + \frac{1}{2} \sum_i \alpha_i^* \alpha_i^2$$

for finding constant Error contour putting the value of  $E(w)$  to  $c$  we get

$$E[w^*] + \frac{1}{2} \sum_i \alpha_i^* \alpha_i^2 = c$$

which is

$$\sum_i \alpha_i^* \alpha_i^2 = 2c - 2E[w^*] = \bar{c} \quad [\text{Another constant}]$$

This is the equation for an ellipse -  
the axes of this ellipse are aligned with the coordinates described by the variables  $\alpha_i$ .

for finding length of 1 axis  $j$ , we set  $\alpha_i = 0$  for all  $i \neq j$  & solving  $\alpha_j$  gives

$$\alpha_j = \left( \frac{\bar{c}}{d_j} \right)^{1/2}$$

which is inversely proportional to square root of eigen value  $d_j$ .

□ Hence prove

ANS-6 →

Number of species = 200

Images per species = 20

With small amount of data like this it is impossible to train a decent size neural network from scratch.

So the only option for deploying deep learning model is to use transfer learning.

→ Base model being trained on 1 million images from 1000 classes learn basic features such as edges & blobs in initial layers so we can freeze some initial layers while retraining based model on target dataset.

→ Or we can set very low Learning rate for initial layer, medium Learning rate for middle layers and High Learning rate for Last layer

- We change Last Layer Prediction from 1000 class to 20 class as target dataset only contain 20 species.
- We devide the available images into training & validation & test set if needed.  
or just train & test set so we are well aware that our model is not overfitting or underfitting.
- With this small number of images it may still be difficult to transfer base model on target dataset so we can use image augmentation like Random Crop & Flip, Shear, Colour Augmentation & so on to increase number of images.  
& by this way our model will learn more robust representation.
- Early stopping can be performed by simultaneously checking the model output on test set.
- If base model is Large (Parameter wise) some of the Last layers can be removed & then softmax with 20 output can be added.