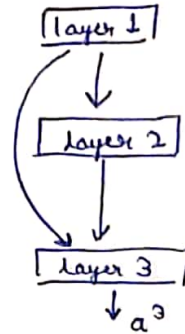


1. Consider Resnet Model :

Let $W^{2,3}$ be weight matrix for layer 2 to layer 3.

Let $W^{1,3}$ be weight matrix for connecting layer 1 to 3.



\therefore Forward propagation is given by -

$$a^3 = f(z^3 + W^{1,3} \cdot a^1) \quad \text{where } a^{(i)} \text{ is activations of neurons in layer } i.$$

f is activation function of layer 3.

$$z^3 = W^{2,3} a^2 + b^3$$

Now coming to back propagation :

~~$$\Delta W^{2,3} \rightarrow \Delta_{W^{2,3}} L(W; x, y) = \delta^3 (a^{(2)})^T$$~~

$$\Delta_{W^{1,3}} L(W; x, y) = \delta^3 (a^{(1)})^T$$

and for updating,

$$W^{1,3} = W^{1,3} - \eta \Delta_{W^{1,3}} L(W; x, y)$$

$$W^{2,3} = W^{2,3} - \eta \Delta_{W^{2,3}} L(W; x, y)$$

where η is learning rate.

&

$\delta^{(k)}$ is error term in k^{th} layer.

2. With each layer support increases from $n \times n$ to $(n+2) \times (n+2)$.

In first layer it is 3×3 . In second layer it will be

5×5 . In third layer it will be 7×7 . In fourth layer,

it will be 9×9 .

\therefore Support of a neuron in 4^{th} non-image layer = $9 \times 9 = 81$ pixels

3. We know that more complex models have lower bias and higher variance.

And Increasing the number of hidden units will lead to more complex model. Hence it will lead to lower bias and higher variance.

$$\begin{aligned}
 4. \quad \tanh(x) &= \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{1 - e^{-2x}}{1 + e^{-2x}} \quad \left(\begin{array}{l} \text{Divide by } e^x \\ \text{in numerator and} \\ \text{denominator} \end{array} \right) \\
 &= \frac{1}{1 + e^{-2x}} - \frac{(1 + e^{-2x} - 1)}{1 + e^{-2x}} = \sigma(2x) - 1 + \frac{1}{1 + e^{-2x}} \\
 &= 2\sigma(2x) - 1
 \end{aligned}$$

\therefore She is correct as there is a relation between $\sigma(x)$ and $\tanh(x)$.

$$\begin{aligned}
 \therefore \text{ we have } \quad \tanh\left(\frac{x}{2}\right) &= 2\sigma(x) - 1 \\
 \Rightarrow \sigma(x) &= \frac{1 + \tanh(x/2)}{2}
 \end{aligned}$$

So we will substitute value of σ in y :

$$\begin{aligned}
 y_k &= \sum_{j=1}^M w_{kj}^{(2)} \left(\frac{1}{2} + \frac{1}{2} \tanh\left(\frac{\sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)}}{2} \right) \right) + w_{k0}^{(2)} \\
 &= \sum_{j=1}^M \frac{w_{kj}^{(2)}}{2} \tanh\left(\frac{\sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)}}{2} \right) + w_{k0}^{(2)} + \sum_{j=1}^M \frac{w_{kj}^{(2)}}{2}
 \end{aligned}$$

Hence we can write as -

$$y_k = \sum_{j=1}^M w'_{kj}^{(2)} \tanh\left(\frac{\sum_{i=1}^D w'_{ji}^{(1)} x_i + w'_{j0}^{(1)}}{2} \right) + w'_{k0}^{(2)}$$

$$\text{where } w'_{kj}^{(2)} = \frac{w_{kj}^{(2)}}{2} \quad \text{for } j=1 \text{ to } M$$

$$w'_{ji}^{(1)} = \frac{w_{ji}^{(1)}}{2} \quad \text{for } j=1 \text{ to } M \text{ and } i=0 \text{ to } D.$$

$$w'_{k0}^{(2)} = w_{k0}^{(2)} + \sum_{j=1}^M \frac{w_{kj}^{(2)}}{2}$$

Hence,
Parameters of two models differ only by linear transformation. }

5. It is given that $H u_i = \lambda_i u_i$ — (i) where $\{u_i\}$ forms orthonormal set. i.e. $\langle u_i, u_j \rangle = 1$ if $i = j$
 $= 0$ otherwise — (ii)

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

Putting $w - w^* = \sum_i \alpha_i u_i$

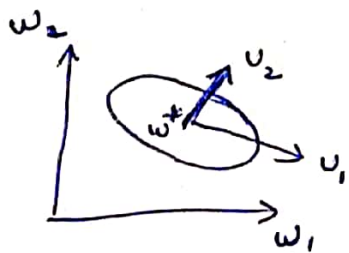
we have $E(w) \approx E(w^*) + \frac{1}{2} \left(\sum_i \alpha_i u_i \right)^T H \left(\sum_i \alpha_i u_i \right)$

$$\therefore E(w) \approx E(w^*) + \frac{1}{2} \left(\sum_i \alpha_i u_i \right)^T \left(\sum_i \alpha_i \lambda_i u_i \right) \text{ \{from (i)\}}$$

$$\approx E(w^*) + \frac{1}{2} \sum_i \alpha_i^2 \lambda_i \text{ \{from (ii)\}}$$

\therefore Contours of constant error are ellipses whose axes are aligned

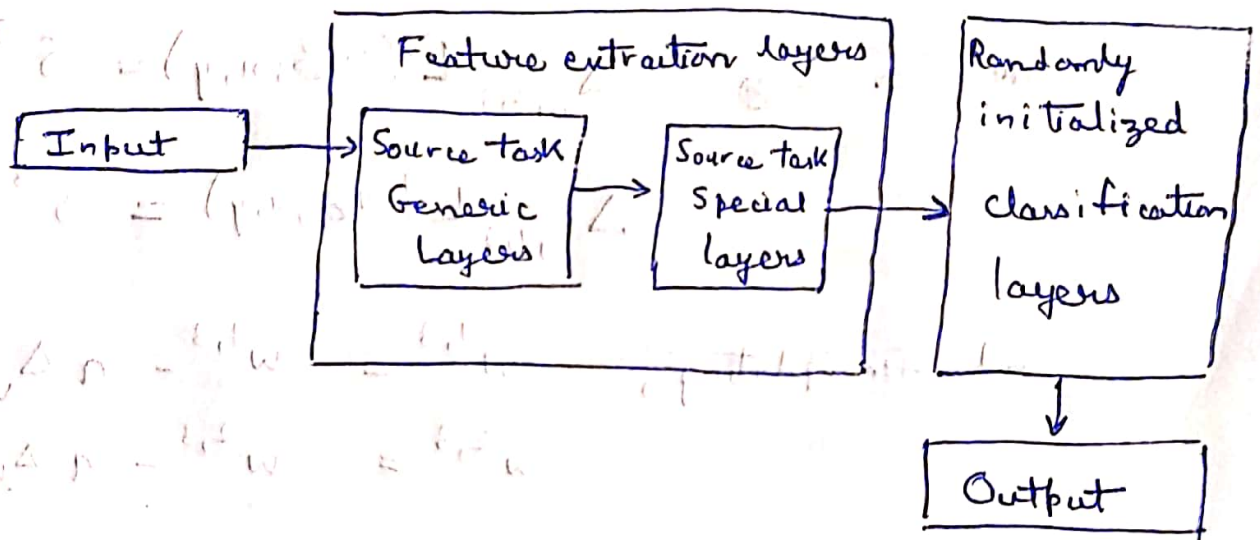
with u_1, u_2 and length of axes are $\frac{1}{\sqrt{\lambda_1}}$ and $\frac{1}{\sqrt{\lambda_2}}$.



Ellipse is centered at w^* .

6. We know that target and source datasets are similar because ~~one is of~~ & both includes images of animals from two different national parks.

So also for Kaziranga National Park, size of dataset is small (only 20 images of all 200 species as compared to 1 million images). So we will use following model:



In this model we use specialized features extractor of model deployed at Olympic National Park as they are similar for source and target datasets. We randomly initialize parameters of classification layer and train them.

Rest of the network remains frozen / unchanged in order to avoid overfitting.