## Assignment-3
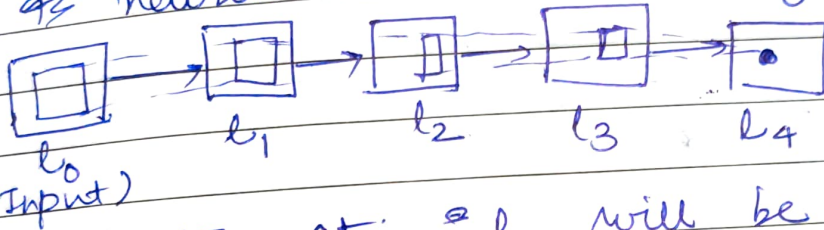
Puneet - Mangla    CS17BTECH11029

Q.2  When a 3×3 filter with stride 1 is applied on an n×n input, the dimensions get reduced to $(n-2) \times (n-2)$

For $9^{th}$ neuron in $4^{th}$ non-image layer



$l_0$ (Input)    $l_1$    $l_2$    $l_3$    $l_4$

the support at $l_3$ will be  3×3
"       "      at $l_2$  "    "   5×5
"       "      at $l_1$  "    "   7×7
"       "      at $l_0$  "    "   9×9

Thus, the support is at $l_0$ or input layer is 81

Q.3  Adding an extra hidden layer will decrease the bias and increase the variance.
Adding a hidden layer will increase the representation power and thus will lead to lower bias on given run.

However, Adding layer will increase the complexity and, making loss surfaces complex and difficult to converge at a same point, leading to diff variance in performance.

Q.5

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

and H has eigenvalues $\lambda_i$ corresponding to eigenvector $u_i$

By linear independence of eigenvectors $u_i$

$$(w - w^*) = \sum_{i \neq} \alpha_i u_i = U\alpha$$

where $U = [u_1, u_2, u_3 -- u_n]$ and $\alpha = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{bmatrix}$

H has eigenvalue decomposition ie

$$H = U \Lambda U^T \text{ where } \Lambda = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & 0 \\ & & \ddots & \\ 0 & & & \lambda_n \end{bmatrix}_{n \times n}$$

and $U U^T = I$

now $E(w) = E(w^*) + \frac{1}{2}(U\alpha)^T U \Lambda U^T (U\alpha)$

$$= E(w^*) + \frac{1}{2} \alpha^T U^T U \Lambda U^T U \alpha \qquad \begin{cases} U^T U = I \\ \downarrow \end{cases}$$
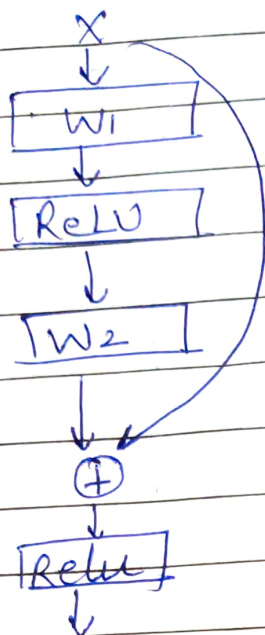
$$= E(w^*) + \frac{1}{2} \alpha^T \Lambda \alpha$$

$$= E(w^*) + \frac{1}{2}\left[ \alpha_1^2 \lambda_1 + \lambda_2 \alpha_2^2 -- \lambda_n \alpha_n^2 \right]$$

$$= E(w^*) + \frac{1}{2} \sum_i \lambda_i \alpha_i^2 \qquad \begin{cases} \frac{x_1^2}{a_1^2} + \frac{x_2^2}{a_2^2} -- = 1 \end{cases}$$

This represents an ellipse. Since we represented vectors as linear combination of $v_i$'s, the axis will be $v_i$'s

now the length $a_i = \dfrac{\sqrt{2(E(w) - E(w^*))}}{\sqrt{\lambda_i}}$ Hence $a_i \propto \dfrac{1}{\sqrt{\lambda_i}}$

**Q.1**



Diagram (top to bottom): x → [W₁] → [ReLU] → [W₂] → ⊕ → [ReLU] →, with an "identity" skip connection.

Assumption: we assume a simple 2 layer fully connected residual block with relu activation. Let $\sigma$ = Relu

Bias is false.

$$y = \sigma\!\!\!\underset{\downarrow}{\overset{R}{\phantom{=}}}\!\!\!\big(w_2 \cdot \sigma(w_1 x) + x\big)$$

feed forward equation

$$\frac{\partial y}{\partial W_2} = \big(w_2 \sigma(w_1 x) + x\big)\; \sigma(w)$$

$$\frac{\partial y}{\partial w_2^{ij}} = \sigma'\big(w_2\,\sigma(w_1 x) + x\big) \cdot \big(\sigma(w_1 x)_i\big)$$

Let $z = w_1 x$

$$\frac{\partial y}{\partial z_i} = \sigma'\big(w_2\,\sigma(w_1 x) + x\big) \cdot \sum_{\forall j} w_2^{ji} \cdot \sigma'(z_i)$$

$$\frac{\partial y}{\partial w_1^{ij}} = \frac{\partial y}{\partial z_j} \cdot \frac{\partial z_j}{\partial w_2^{ij}} = \frac{\partial y}{\partial z_j} \cdot x_i$$

$$\frac{\partial y}{\partial x_j} = \sigma'\big(w_2\,\sigma(w_1 x) + x\big) + \left[\sum_{\forall i} \frac{\partial y}{\partial z_i} \cdot \frac{\partial z_i}{\partial x_j}\right]$$

$$= \sigma'\big(w_2\,\sigma(w_1 x) + x\big) + \left[\sum_{\forall i} \frac{\partial y}{\partial z_i} \cdot \sigma'(x_j \cdot w_{ji})\right]$$

**Q.4** $\sigma(a) = \dfrac{1}{1+e^{-a}}$  also  $\sigma(-a) = 1 - \sigma(a)$

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

$$= \sigma(2a) - \sigma(-2a) \quad \{ \because \sigma(-a) = 1 - \sigma(a) \}$$

$$= 2\cdot\sigma(2a) - 1 \qquad \text{or} \qquad \sigma(a) = \frac{\tanh(a/2) + 1}{2}$$

$$y_k(x,w) = \sum_{j=1}^{M} w_{kj}^{(2)} \, \sigma\left( \sum_{i=1}^{D} w_{ji}^{(1)} x_i + w_{jo}^{(1)} \right) + w_{ko}^{(2)} \qquad (1)$$

$$y_k(x,\tilde{w}) = \sum_{j=1}^{M} \tilde{w}_{kj}^{(2)} \, \tanh\left( \sum_{i=1}^{D} \tilde{w}_{ji}^{(1)} x_i + \tilde{w}_{jo}^{(1)} \right) + \tilde{w}_{ko}^{(2)} \qquad (2)$$

$$y_k(x,w) = \sum_{j=1}^{M} w_{kj}^{(2)} \left[ \frac{\tanh\left( \dfrac{1}{2}\,\sum_{i=1}^{D} w_{ji}^{(1)} x_i + w_{jo}^{(1)} \right) + 1}{2} \right]$$

$$+ w_{ko}^{(2)}$$

$$y_k(x,w) = \sum_{j=1}^{M} w_{kj}^{(2)} \, \frac{1}{2} \tanh\left( \sum_{i=1}^{D} w_{ji}^{(1)} x_i + w_{jo}^{(1)} \right)$$

$$+ \sum_{j=1}^{M} w_{kj}^{(2)} \frac{1}{2} + w_{ko}^{(2)} \qquad (3)$$

Comparing (2) and (3)

$$\tilde{w}_{kj}^{(2)} = \frac{w_{kj}^{(2)}}{2} \;,\quad \tilde{w}_{ji}^{(1)} = \frac{w_{ji}^{(1)}}{2} \;,\quad \tilde{w}_{jo}^{(1)} = \frac{w_{jo}^{(1)}}{2}$$

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

**Q.6** Since the custom data set is small and both target and source datasets are similar.

- We well use transfer learning, where ~~speciali~~ specialized feactures and generic layers both will be fixed.

- Parameters of final classification will be initialized again and trained.