

Name - Pratham Singhal 2021082

ASSIGNMENT - I

Answer 1) (a) MSE does not take account of the probabilistic nature b/w dependent and independent variables because it assumes the relationships between two ~~to be~~ as linear. Also in classification based learning, it leads to slow convergence.

(b) BCE =  $[-y \log \hat{y} + (1-y) \log (1-\hat{y})]$

(c) BLE =  $-[y \log \hat{y} + (1-y) \log (1-\hat{y})]$   
 $\Rightarrow y = 0 \quad \hat{y} = 0.9$

$$BLE = -1 \log (0.1)$$

$$BCE = 2.303$$

(d) Avg Loss =  $\frac{1}{N} \sum_{i=1}^N (y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i))$   
=  $-\frac{1}{3} [1 \log 0.1 + 1 (\log 0.8) + 1 \log 0.3]$   
=  $-\frac{1}{3} [\log (2.303 + 0.223 + 1.204)]$

$$\text{Avg Loss} = 1.243$$

(e) with L2 regularization.

$$L_{\text{BCE}} = \frac{1}{m} \sum -y \log \hat{y} + -(1-y) \log(1-\hat{y}) \\ + \frac{\lambda}{2m} \underbrace{\sum_{i=1}^n w_i^2}_{\text{Regularization Term}}$$

$$L_{\text{BCE}} = \frac{1}{m} \sum -y \log \sigma(z) - (1-y) \log(1-\sigma(z)) \\ + \frac{\lambda}{2m} \sum w_i^2$$

$$\frac{\partial L}{\partial w} = \frac{-1}{m} \left( \sum y \cdot \frac{1}{\sigma(z)} \cdot \sigma'(z) [1 - \sigma(z)] \cdot \frac{\partial w^T x}{\partial w} \right. \\ \left. + (1-y) \frac{\partial w x}{\partial w} \cdot \frac{1}{1 - \sigma(z)} [-\sigma(z)(1 - \sigma(z))] \right) \\ + \frac{\lambda}{2m} \cdot 2 \sum_{i=1}^n w_i^2$$

$$\frac{\partial L}{\partial w} = \frac{1}{m} \left[ - \left( \sum (y - \hat{y}) x_i \right) + \lambda \sum_{i=1}^n w_i^2 \right]$$

$$W_{\text{new}} = W_{\text{old}} - \alpha \left[ \frac{1}{m} \left[ E \left( \sum_{j=1}^{i=m} (y_j \hat{y}_j) X_j \right) + \lambda \sum_{i=0}^n \right] \right]$$

f)

KL divergence is a measure of the average distribution diverges from a second, expected probability distribution.

Cross entropy  $H(P|Q)$  measure of average no. of bits needed to identify an event from a set of possibilities if a coding scheme is used based on a given prob of distribution  $Q$ , rather than the "true" distribution  $P$ .

KL - divergence

$$D_{KL}(P||Q) = \sum P(x) * \log \left( \frac{P(x)}{Q(x)} \right)$$

i.e.

$$\begin{aligned} D_{KL}(P||Q) &= \sum P(x) (\log P(x) - \log Q(x)) \\ &= \sum P(x) \log P(x) - \sum P(x) \log Q(x) \\ &\rightarrow \sum P(x) \log P(x) + D_{KL}(P||Q) = -\sum P(x) \log Q(x) \end{aligned}$$

$$\text{ie } \exp \log^P + D_{KL}(P||Q) = \exp P \log Q$$

$$H^P + D_{KL}(P||Q) = H(P, Q)$$

Aus 2) a)  $W_2 = k \times Da$   
 $b_2 = k \times 1$   
 $(Da \times m)$

b)  $\frac{\partial \hat{y}_k}{\partial z_k^{(2)}} = \hat{y}_k(1 - \hat{y}_k)$

c)  $\frac{\partial \hat{y}_k}{\partial z_i^{(2)}} = -\hat{y}_k \hat{y}_i$

d)  $L = \sum -y_i \log \hat{y}_i$  where,  $\hat{y}_i = \text{softmax}(z_i^{(2)})$

$$\frac{\partial L}{\partial z_i} = \frac{\partial \sum -y_i \log \hat{y}_i}{\partial z_i} \quad \hat{y}_i = \text{softmax}(z_i)$$

$$= \frac{\partial}{\partial z_i} \sum -y_i \log (\text{softmax}(z_i))$$

$$= -y_i \sum \frac{\partial}{\partial z_i} \log (\text{softmax}(z_i))$$

$$-\frac{y_i}{\text{softmax}(z_i(w))} \cdot \frac{\partial}{\partial z_i} \text{softmax}(z_i(w))$$

$$= -\frac{y_i}{\text{softmax}(z_i(w))} \cdot \frac{\partial}{\partial z_i} \frac{e^{z_i(w)}}{\sum e^{z_j(w)}}$$

case 1

$$-\frac{y_i}{\hat{y}_i} \cdot \hat{y}_i(1 - \hat{y}_i) \quad (i \neq k)$$

$$\frac{\partial L}{\partial z_i} = -y_i(1 - \hat{y}_i)$$

case 2  $i = k$

$$\frac{\partial L}{\partial z_i} = -\frac{y_i}{\hat{y}_i} - \sum_{j \neq i} \frac{y_j}{\hat{y}_j} \hat{y}_j \quad [i = k]$$

$$\frac{\partial L}{\partial z_i} = \frac{y_i}{\hat{y}_i} - \sum_{j \neq i} \frac{y_j}{\hat{y}_j} \hat{y}_j$$

$$\frac{\partial L}{\partial z_i} = y_i \hat{y}_i$$

⑥ Numerical stability in softmax is due to computation of exponents of large +ve/-ve values causing an overflow/underflow error.

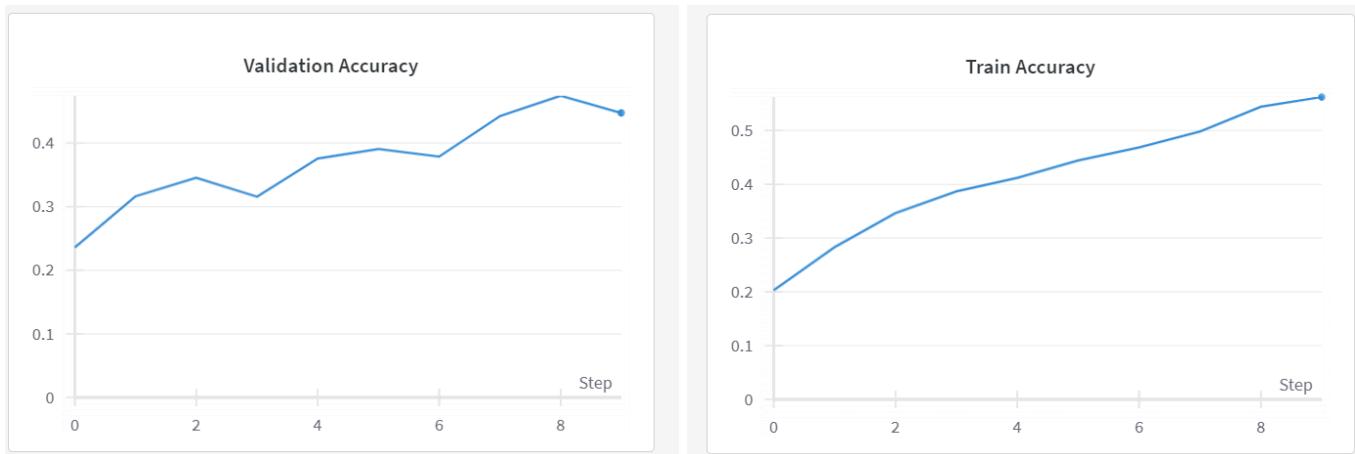
Overflow can lead to  $\infty$  & underflow  
can lead to 0.

This can issue a possibly  
meaningful gradient to the function

$$\hat{y}_k = \frac{e^{z_k^{(2)} - \max(z^{(2)})}}{\sum_{j=1}^K e^{z_j^{(2)} - \max(z^{(2)})}}$$

Answer2>

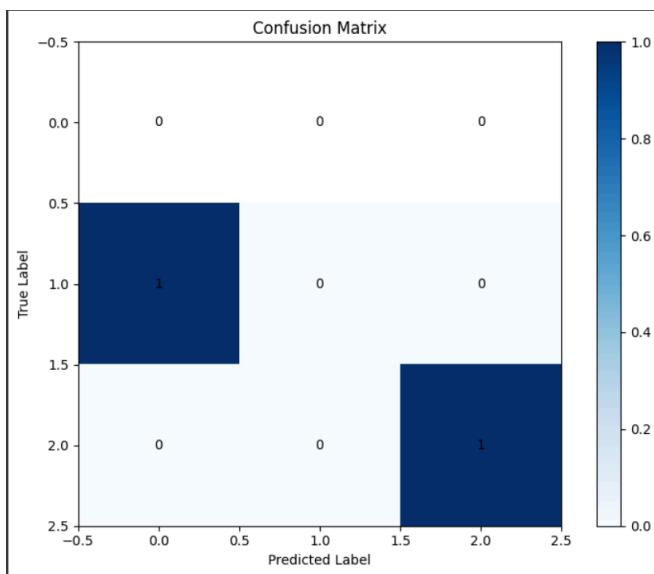
Ans 2.2.c>



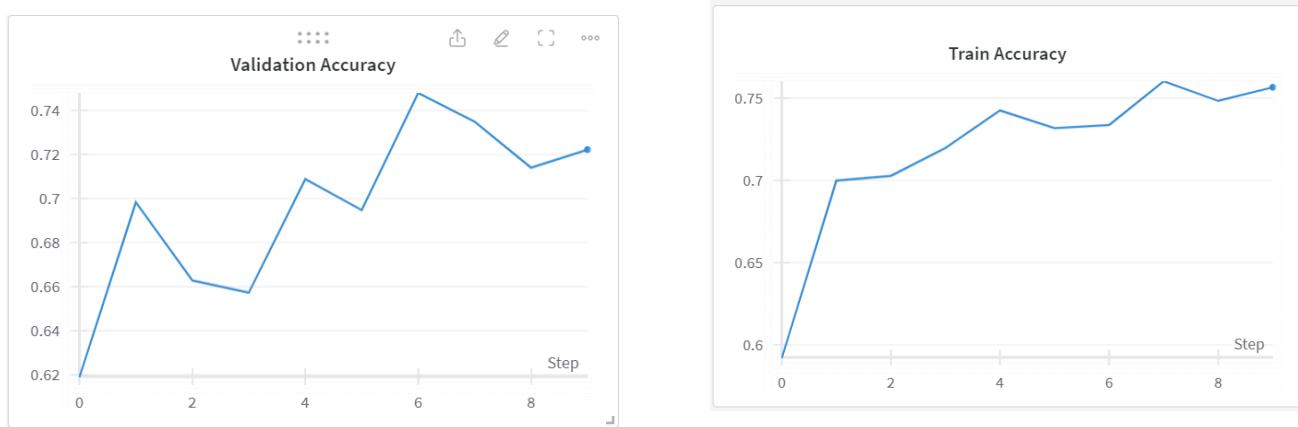
The trained CNN model is overfitting to a slight extent, since the training accuracy is 56% and validation accuracy is 46%.

ans2.2.d>

Test accuracy is 58% and F1 score = 0.5

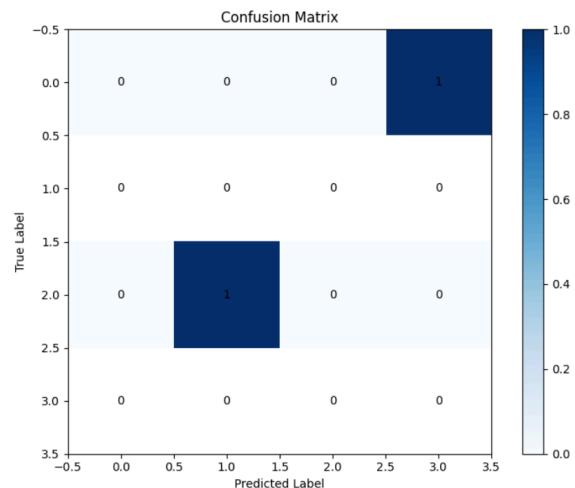


Ans 2.3.b>

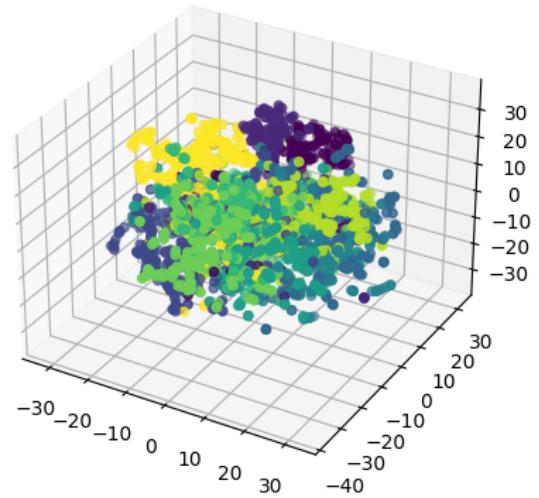
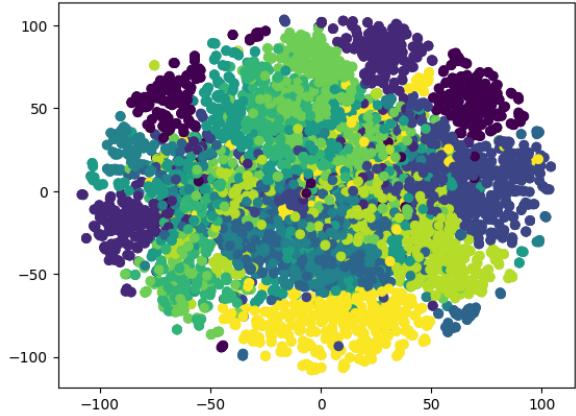


The model is not overfitting, since the training accuracy is 75% and validation accuracy is 72%. The gap is low between them.

Ans2.3.c> Test Accuracy = 71% and f score = 0.68



Ans 2.3.d>



Ans2.4.c>

The test accuracy = 71 % and F1 score = .69

Ans2.6> the cnn model performed a little less accurately compared to all the other models. It was overfitting to a bit. The resnet model without data augmentation performed better than the other all models. This was the best accurate model, but was overfitting a little. The resnet model with data augmentation was not overfitting but performed a little less accurately. This was due to the different types of datasets.