Theory Report     Abhishek Sushil    2021441

Q1.

a) MSE → simply minimizes dist b/w real & predicted values
   Based on squared $L_2$ norm of actual & predict$^n$
   ∴ No basis for false +ves or false -ve as
   it takes squared dist
   Eg: Here labelling a sweet papaya is not sweet
   doesn't affect customer but the opposite does

b) Bin Cross Entropy $(\hat{y}, \hat{y}) = - \left[ y \log_2 \hat{y} + (1-y) \log_2 (1-\hat{y}) \right]$
   real ← → predicted probability

c) $\hat{y}$ for class 0 $= 0.9$ ∴ $\hat{y}$ for class 1 $= 0.1$

   $$\mathcal{L}(y, \hat{y}) = - \left[ 1 \log_2 0.1 + 0 \cdot \log_2 0.9 \right]$$

   $$= - \log_2 0.1 \approx 3.32$$

d) $S_1 \rightarrow \mathcal{L}(1, 0.1) = 3.32$    as above
   $S_2 \rightarrow \mathcal{L}(0, 0.2) = - \left[ 0 \log 0.2 + 1 \log 0.8 \right]$
   $$= 0.32$$
   $S_3 \rightarrow \mathcal{L}(0, 0.7) = - \left[ 0 \log 0.7 + 1 \log 0.3 \right]$
   $$= 1.73$$

   ∴ Average $= 1.79$

e) $$W = W - \alpha \cdot \frac{\partial L_{BCE}}{\partial W} + 2\lambda W$$

More penalty as $2\lambda W$ term present $\therefore$ a decrease in weights for model A that uses $L_2$ regularizer.
Advantage for model A as $L_2$ prevents overfitting

f) KL Divergence
→ Used to see how far apart 2 probability distribution are based on how well one distribution is likely to generate samples from the other

$$D_{KL}(P(x) | Q(x)) = \sum_{x \in X} P(x) \left( \frac{\log(P(x))}{\log(Q(x))} \right)$$

Cross Entropy
→ Diff. b/w actual & predicted probability distribution

$$H(P(x), Q(x)) = - \sum_{x \in X} P(x) \cdot \log(P(x))$$

$$D_{KL} = \sum_{x \in X} P(x) \left( \frac{\log(P(x))}{\log(Q(x))} \right)$$

$$= \sum_{x \in X} P(x) \left[ \log(P(x)) - \log(Q(x)) \right]$$

$$= \sum_{x \in X} P(x) \log P(x) - \sum_{x \in X} P(x) \log Q(x)$$

$$D_{KL} = -H(Px) + H(P(x), Q(x))$$

$$\Rightarrow D_{KL}(P(x) \| Q(x)) + H(P(x)) = H(P(x), Q(x))$$

**Q2.**

**a)**

$\dim(W^{[2]}) = k \times Da$

$\dim(b^{[2]}) = k \times 1$

$\dim$ of hidden layer ay 'm' samples - $Da \times m$

**b)**

$$\hat{y}_k = \frac{e^{z_k^{[2]}}}{\sum_{i=0}^{k} e^{z_i^{[2]}}}$$

$$\therefore \quad \frac{\partial \hat{y}_k}{\partial z_k^{[2]}} = \partial \left[ \frac{e^{z_k^{[2]}}}{\sum e^{z_i^{[2]}}} \right] \Big/ \partial z_k^{[2]}$$

$$= \frac{1}{\left( \sum_i e^{z_i^{[2]}} \right)^2} \left[ \frac{\partial e^{z_k^{[2]}}}{\partial z_k^{[2]}} \cdot \sum_{i=0}^{k} e^{z_i^{[2]}} - \frac{\partial \sum_{i=0}^{k} e^{z_i^{[2]}}}{\partial z_k^{[2]}} \cdot e^{z_k^{[2]}} \right] \quad \text{--- Division Rule}$$

$$= \frac{1}{\left( \sum_{i=0}^{k} e^{z_i^{[2]}} \right)^2} \left[ e^{z_k^{[2]}} \cdot \sum_{i=0}^{k} e^{z_i^{[2]}} - e^{z_k^{[2]}} \cdot e^{z_k^{[2]}} \right]$$

$$= \frac{e^{z_k^{[2]}}}{\left( \sum_{i=0}^{k} e^{z_i^{[2]}} \right)} - \left[ \frac{e^{z_k^{[2]}}}{\sum_{i=0}^{k} e^{z_i^{[2]}}} \right]^2$$

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

c) $i \neq k$

$$\hat{y}_k = \frac{e^{z_k^{[2]}}}{\sum_{i=0}^{k} e^{z_i^{[2]}}}$$

$$\frac{\partial \hat{y}_k}{\partial z_i^{[2]}} = \frac{\partial \left[ \frac{e^{z_k^{[2]}}}{\sum_{i}^{k} e^{z_i^{[2]}}} \right]}{\partial z_i^{[2]}} \longrightarrow \frac{\partial e^{z_k}}{\partial z_i} = 0$$

$$= \frac{1}{\left( \sum_{i=0}^{k} e^{z_i^{[2]}} \right)^2} \left[ 0 - e^{z_k^{[2]}} \cdot e^{z_i^{[2]}} \right]$$

$$= - \frac{e^{z_k^{[2]}}}{\sum_{i=0}^{k} e^{z_i^{[2]}}} \circ \frac{e^{z_i^{[2]}}}{\sum_{i=0}^{k} e^{z_i^{[2]}}}$$

$$\frac{\partial \hat{y}_k}{\partial z_i^{[2]}} = - \hat{y}_k \circ \hat{y}_i$$

d) i) $q = k$

$$\frac{\partial L}{\partial z_k^{[2]}} = \frac{\partial L}{\partial \hat{y}_k} \circ \frac{\partial \hat{y}_k}{\partial z_k^{[2]}} \qquad \text{chain rule}$$

$$= \frac{\partial \left[ - \sum_{i=0}^{k} y_i \log \hat{y}_i \right]}{\partial \hat{y}_k} \circ \left( \hat{y}_k (1 - \hat{y}_k) \right) \hat{y}_i \quad \text{from b}$$

$$= \frac{y_k \log \hat{y}_k}{\partial \hat{y}_k} \left( \hat{y}_k (1 - \hat{y}_k) \right)$$

$$= \frac{y_k (\hat{y}_k)(1 - \hat{y}_k)}{y_k}$$

$$\frac{\partial L}{\partial z_k^{[2]}} = \frac{y_k}{\searrow 1} (1 - \hat{y}_k) = (1 - \hat{y}_k)$$

i) B $i \neq k$

$$\frac{\partial L}{\partial z_i^{[2]}} = \frac{\partial L}{\partial \hat{y}_k} \cdot \frac{\partial \hat{y}_k}{\partial z_i^{[2]}}$$

$$= \frac{\partial \left[ -\sum_{i=0}^{K} y_i \log \hat{y}_i \right]}{\partial \hat{y}_k} \cdot \left( -\hat{y}_i \hat{y}_k \right) \quad \text{--- from } \circled{c}$$

$$= \frac{y_k}{\hat{y}_k} \cdot -1 \times \hat{y}_i \cdot \hat{y}_k$$

$$= -1 \cdot y_k \cdot \hat{y}_i$$

$$\frac{\partial L}{\partial z_i^{[2]}} = -\hat{y}_i$$

e) Numerical Instability can be encountered when dealing w/ very large or very small values. Those can turn when taken in exponents very large values lead to numerical overflow & very small nos. to underflow

Assume the final layer has a vector w/ values in a similar range we can normalise the vector before soft -maxing process

$[a, \quad b, \quad c] \rightarrow$ last layer $\qquad \bar{x} = \text{mean} = \frac{a+b+c}{3}$

$\oplus$ $\qquad\qquad\qquad\qquad\qquad \sigma = \text{std dev}(a,b,c)$

$\left[ \frac{a-\bar{x}}{\sigma}, \quad \frac{b-\bar{x}}{\sigma}, \quad \frac{c-\bar{x}}{\sigma} \right] \rightarrow$ modified last lay

Now apply soft max on the modified layer