

## Theory

Answer 1

- a) We have a classification problem where we have to label papaya

$$p(x_i) \begin{cases} \text{sweet} = 1 & n_i = 1 \\ \text{not} = 0 & n_i = 0 \end{cases}$$

We want to maximize the likelihood function of  $p(x, \theta)$ : function to correctly predict the label

$$L(x_1, x_2, \dots, x_n) = P(x_1=x_1, x_2=x_2, \dots, x_n=x_n | y_1=y_1, y_2=y_2, \dots, \theta)$$

$$= \prod_{y_i} p(x_i) \cdot \prod_{1-y_i} 1-p(x_i)$$

$$L(\theta) = \prod p(n_i)^{y_i} (1-p(n_i))^{1-y_i}$$

we need to find  $\theta$  parameter that maximizes likelihood function which is BINARY CROSS ENTROPY function

MSE is NOT a good option because:

- 1) MSE assumes error to be a Normal distribution, But in the classification problem error can take only 2 values {0, 1} and is a Bernoulli R.V. So, use of MSE doesn't make sense.

- 2) Non Linear Relationship: MSE is a linear function while use of sigmoid function make the model's output non linear
- 3) Poor Indication of Misclassification: MSE doesn't provide clear indication of how many instances were misclassified, it focuses on magnitude of error.
- 4) Gradient Vanishing Problem

(b) BCE : Binary Cross Entropy Loss

$$L(y, \hat{y}) = - \sum_{c=1}^n y_{0,c} \log(\hat{y}_{0,c})$$

For binary cross entropy :

$$L(y, \hat{y}) = -(y_i \log \hat{y}_i + (1-y_i) \log(1-\hat{y}_i))$$

c)  $\hat{y} = 0.9$

$$\begin{aligned} L(y, \hat{y}) &= -(0 \log 0.9 + 1 \log 0.1) \\ &= -\log 0.1 = +2.303 \end{aligned}$$

d)  $y = [1 \ 0 \ 0]$   
 $\hat{y} = [0.1 \ 0.2 \ 0.7]$

$$\begin{aligned} L(y, \hat{y}) &= - \left( \sum_{i=0}^n y_i \log \hat{y}_i \right) & L = BCE(y_1, \hat{y}) + \frac{BCE(y_2, \hat{y}) + BCE(y_3, \hat{y})}{3} \\ &= - (1 \log 0.1 + 0 \log 0.8 + 0 \log 0.3) + 1 \log 0.2 \\ &= -\log 0.1 + 1 \log 0.2 \\ &= +2.303 + 0.2231 + 1.2041 \approx 1.24 \end{aligned}$$

Date: \_\_\_ / \_\_\_ / \_\_\_

a) contd.

1<sup>st</sup> sample :  $BCE(y_1, \hat{y}_1) \approx 2.3026$   
2<sup>nd</sup> sample :  $BCE(y_2, \hat{y}) \approx 0.2231$   
3<sup>rd</sup> sample :  $BCE(y_3, \hat{y}) = 1.2041$

$$L = BCE(y_1, \hat{y}_1) + \frac{BCE(y_2, \hat{y})}{3} + BCE(y_3, \hat{y})$$

$$L \approx 2.3026 + \frac{0.2231}{3} + 1.2041 \approx 1.24$$

e)  $w_{\text{new}} \leftarrow w - \alpha \left( \frac{\partial L_{\text{BCE}}}{\partial w} + \lambda w \right)$

We penalize big weights so as to prevent overfitting.  
In  $L_2$  we penalize it with  $\frac{\partial L_{\text{BCE}}}{\partial w} + \lambda w$

If  $\mathcal{S}$  have 2 model's A & B

A uses  $L_2$  regularization  
B does NOT

any weight  $w$  that is larger than the approximate central tendency would be penalized in A

Hence weights in A would [generally be smaller] than B

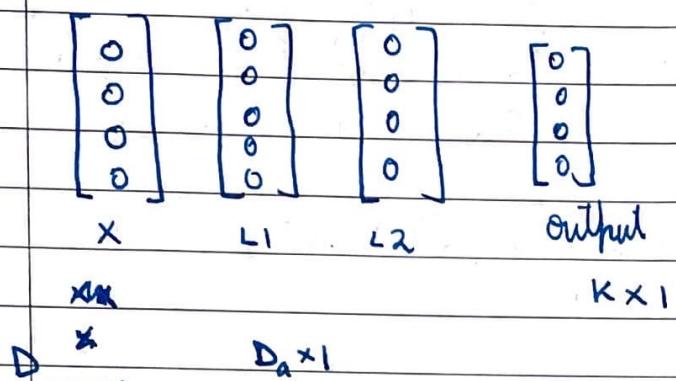
exception: if  $w$  is -ve

then  $\lambda w + \frac{\partial L_{\text{BCE}}}{\partial w}$  term is may also

be negative and in that case  $w$  is pushed towards 0

Hence, in the case when  $w$  is -ve,  
 $w_B \rightarrow w_A$   $w_A < w_B$

Ans 2



$$z^1 = W^1 x + b^1$$

$$a^1 = \text{Leaky ReLU}(z^1, \alpha=0.01)$$

$$z^2 = W^2 a^1 + b^2$$

$$\hat{y} = \text{softmax}(z^2)$$

$$L = -\sum y_i \log \hat{y}_i$$

a)  $z^1 = W^1 x + b^1$   
 $D_a \times 1 \quad \downarrow \quad D_n \times 1 + D_a \times 1$   
 $D_a \times D_x$

$$\text{shape of } W^1: D_a \times D_x$$

$$b^1 = D_a \times 1$$

$$a^1 = D_a \times 1$$

$$z^2 = W^2 a^1 + b^2$$

$K \times 1 \quad D_a \times 1 \quad K \times 1$

$w^2 : K \times K$ if hidden layer if batch size  
is m $w^2 : K \times D_a$  $b^2 : K \times 1$  $z^{[2]} : D_a \times m$ (b) We have got  $\hat{y}$ : which is a  $K \times 1$  matrix $y$ : one hot encoding vector  
~~of~~  $K \times 1$  matrix

We calculate loss as

$$L = - \sum y_i \log(\hat{y}_i)$$

and backpropagate loss

$$\hat{y}_i = \frac{e^{z_i^{[2]}}}{\sum_{j=1}^K e^{z_j^{[2]}}}$$

$$\hat{y}_k = \frac{e^{z_k^{[2]}}}{\sum_{i=1}^K e^{z_i^{[2]}}} =$$

$$-(1 - y_k)$$

$$\frac{\partial \hat{y}_k}{\partial z_k^{[2]}}$$

$$= \sigma(z_k^{[2]}) (1 - \sigma(z_k^{[2]})) =$$

$$y_k(1 - y_k)$$

 $\sigma(x)$ : signifies sigmoid fn of on x

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

$$\Theta = \boxed{\frac{y_i/K - \hat{y}_i}{y_i(1 - y_i)}} \quad \boxed{-y_i y_k}$$

(d)  $L = - \sum y_i \log \hat{y}_i$

$$\frac{\partial L}{\partial y} = \left[ \frac{\partial L}{\partial y_1}, \frac{\partial L}{\partial y_2}, \dots, \frac{\partial L}{\partial y_n} \right]$$

$$\frac{\partial L}{\partial z^{(2)}} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial z^{(2)}}$$

$$\sum_{j=1}^n \frac{\partial L}{\partial y_j} \cdot \frac{\partial y_j}{\partial z^{(2)}}$$

$\therefore$  only

$$\frac{\partial L}{\partial z^{(2)}} = \frac{\partial L}{\partial y_k} \cdot \frac{\partial y_k}{\partial z^{(2)}}$$

(rest all terms  
are 0 because of  
1 hot encoding vector)

if  $i == k$ :  $\frac{\partial L}{\partial z^{(2)}} = -\frac{1}{y_k} \cdot y_k (1-y_k) = \boxed{\frac{y_k - 1}{y_k}}$   
 $(k = \neq i)$

else if  $i != k$ :  $\frac{\partial L}{\partial z^{(2)}} = -\frac{1}{y_k} \cdot y_k \cdot y_i = \boxed{-y_i}$

e) softmax function :  $\frac{e^{x_i}}{\sum e^{x_i}}$

since softmax exponentiates  $x_i$ , for large  $x$   
 $e^{x_i}$  will be a really big value and hard to  
compute

Date: \_\_\_ / \_\_\_ / \_\_\_

Hence softmax function is numerically **UNSTABLE**

A modification in softmax function is done

$$\frac{e^{x_i - \beta}}{\sum e^{x_i - \beta}}$$

:

where  $\beta$  ref is the max( $x_i$ )  $\forall i$

this transforms all  $x$  into values between [0, 1] and normalizes and provides softmax output