

#1

(a) The answer is number 1.

$$L = \|E(A) - E(P)\|^2 - \|E(A) - E(N)\|^2$$

(b) The answer is number 2.

$$L = \|C(M) - C(G)\|^2 + \|S(M) - S(G)\|^2$$

#2

(a) $\hat{y} = \sigma(wx+b)$ where σ is sigmoid function.

(b)

$$(b) \text{Cross-Entropy Loss} = -(y \log(\hat{y}) + (1-y) \log(1-\hat{y}))$$

$$(c) L = -\frac{1}{m} \sum_{i=1}^m (y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i))$$

$$\frac{\partial L}{\partial \hat{y}_i} = \frac{1}{m} \left(-\frac{y_i}{\hat{y}_i} + \frac{1-y_i}{1-\hat{y}_i} \right)$$

$$\frac{\partial \hat{y}}{\partial w} = \sigma(wx+b) (1-\sigma(wx+b)) \cdot x$$

$$\frac{\partial \hat{y}}{\partial b} = \sigma(wx+b) (1-\sigma(wx+b)) \cdot 1$$

$$\begin{aligned} dw &= \frac{1}{m} \sum_{i=1}^m \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w} \\ db &= \frac{1}{m} \sum_{i=1}^m \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial b} \end{aligned}$$

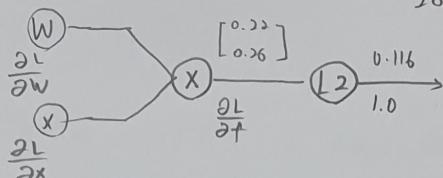
$$\therefore W_{\text{new}} = W_{\text{old}} - \alpha \cdot dw \quad \text{where } \alpha \text{ is learning rate, } dw \text{ and } db \text{ is above.}$$

$$b_{\text{new}} = b_{\text{old}} - \alpha \cdot db$$

#3

$$L_2 \text{ loss} = 0.1$$

$$W = \begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} \quad x = \begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$$



$$Wx = \begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} \begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} = \begin{bmatrix} 0.02 + 0.20 \\ -0.06 + 0.32 \end{bmatrix} = \begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$\frac{\partial L}{\partial W} = \sum_{i=1}^2 2 \cdot (Wx)_i \cdot x_i^T = 2(Wx)x^T, \quad \frac{\partial L}{\partial x} = \sum_{i=1}^2 2 \cdot (Wx)_i \cdot W = 2W^T(Wx)$$

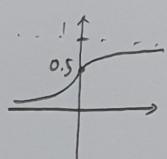
$$\therefore dW = 2 \cdot \begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix} \begin{bmatrix} 0.2 & 0.4 \end{bmatrix}, \quad dx = 2 \cdot \begin{bmatrix} 0.1 & -0.3 \\ 0.5 & 0.8 \end{bmatrix} \begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$= \begin{bmatrix} 0.088 & 0.196 \\ 0.104 & 0.208 \end{bmatrix} \quad = \begin{bmatrix} -1.112 \\ 0.636 \end{bmatrix}$$

#4

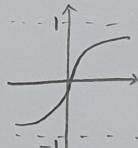
(a)

$$(i) \text{ Sigmoid} : a(z) = \frac{1}{1+e^{-z}}$$



Sigmoid function exists between (0,1). It is specially used for models where we have to predict the probability. It is differentiable, but when z goes big or small, it has gradient vanishing.

$$(ii) \tanh : a(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



Tanh function is also like sigmoid but better. It exists between (-1,1). The negative inputs will be mapped negative. It is differentiable, but it has gradient vanishing problem.

$$(iii) \text{ ReLU} : a(z) = \max(0, z)$$



ReLU function is the most widely used in the world and it helps to solve the gradient vanishing problem. But, any negative inputs turns the value into zero and it has not proper mapping problem.

$$(iv) \text{ Leaky ReLU} : a(z) = \max(0.01z, z)$$



Leaky ReLU function solves gradient vanishing problem and also has non-zero value for negative inputs. It helps to increase the range of the ReLU function.

(b)

$$(i) \text{ sigmoid} : a'(z) = a(z)(1-a(z))$$

$$(ii) \tanh : a'(z) = 1 - \{a(z)\}^2$$

$$(iii) \text{ ReLU} : a'(z) = \begin{cases} 0 & \text{for } z < 0 \\ 1 & \text{for } z \geq 0 \end{cases}$$

$$(iv) \text{ Leaky} : a'(z) = \begin{cases} 0.01 & \text{for } z < 0 \\ 1 & \text{for } z \geq 0 \end{cases}$$

- (a) The purpose of regularization is to prevent overfitting. It means that we have to reduce high variance.
- (b) The effect of regularization on weights W reduces the magnitude of weight ($= |\text{weight}|$) and prevents weights value from getting bigger. By doing this, the degree of decision boundary represented by weights decreases.
- (c) The regularization penalizes to make complex decision boundary which is the result of focusing on minimizing the training loss only. The regularization term in the loss function, it reduces the magnitude of weights which in turns reduces fluctuation in the regression. So, the regularization prevents overfitting and high variance.