

# Neural Network And Deep Learning

## ★ Week 2.

### → Binary Classification

- Logistic regression is an algorithm for binary classification.
- In binary classification our goal is to learn a classifier that can input an image represented by a feature vector  $x$  and predict the corresponding label  $y$  is 1 or 0. For example, if the image is of a cat,  $y$  will be 1 if it's a cat or else it will be 0.

### → Notation (Logistic regression):

- (i) A training example is denoted by a pair  $(x, y)$  where,  $x \in \mathbb{R}^{n_x}$  ( $n_x$ -dimensional feature vector)  $y \in \{0, 1\}$  (label)

- (ii)  $n$  (training examples),  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})\}$

$n_{\text{train}}$ : to emphasise the number of training examples

$n_{\text{test}}$ : testing examples.

$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(n)} \\ | & | & & | \end{bmatrix} \quad \begin{matrix} \uparrow \\ n_x \\ \downarrow \end{matrix}$$

$X \in \mathbb{R}^{n_x \times n}$   
 $X_{\text{shape}} = (n_x, n)$

$\xleftarrow{\quad n \quad}$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(n)} \end{bmatrix} \quad \begin{matrix} \uparrow \\ 1 \times n \\ \downarrow \end{matrix}$$

$Y \in \mathbb{R}$   
 $Y_{\text{shape}} = (1, n)$

### → Logistic Regression:

Given  $x$ , want  $\hat{y} = P(y=1|x)$

$\hat{y} \Rightarrow$  estimate of  $y$ , probability of the chance that  $y$  is equal to one given the input features  $x$ .

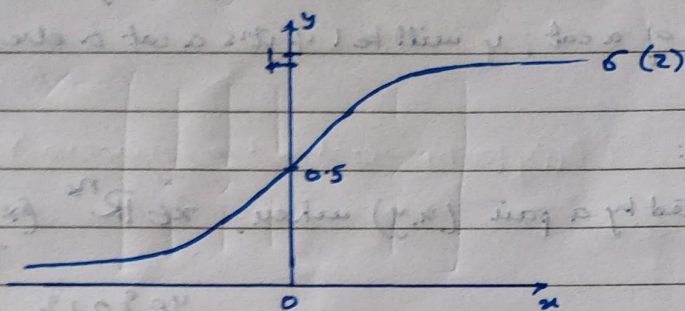


$$x \in \mathbb{R}^n$$

Parameters:  $w \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$

Output:  $\hat{y} = \sigma(w^T x + b)$  ( $w^T$  = transpose of vector  $w$ )

Sigmoid  
 $\sigma^a$



$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

if  $z$  large (+ve)

$$\sigma(z) \approx \frac{1}{1+0} = 1$$

if  $z$  large (-ve)

$$\sigma(z) \approx \frac{1}{1 + \text{Big number}} \approx 0$$

→ To train the parameters, we need to define a cost function.

$$\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b), \text{ where } \sigma(z^{(i)}) = \frac{1}{1 + e^{-z^{(i)}}}$$

Criterion  $\{ (x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}) \}$ , want  $\hat{y}^{(i)} \approx y^{(i)}$  (training set values)  
↓  
output

Loss (cost) function:  $L(\hat{y}, y)$

$L \Rightarrow$  how good the output  $\hat{y}$  is when the true label is  $y$ .  
 function used to measure how accurate the output is.

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

If  $y=1$ :  $L(\hat{y}, y) = -\log \hat{y}$  ← want  $\log \hat{y}$  large, want  $\hat{y}$  large (max 1)

If  $y=0$ :  $L(\hat{y}, y) = -\log(1-\hat{y})$  ← want  $\log(1-\hat{y})$  large, want  $\hat{y}$  small (min 0).

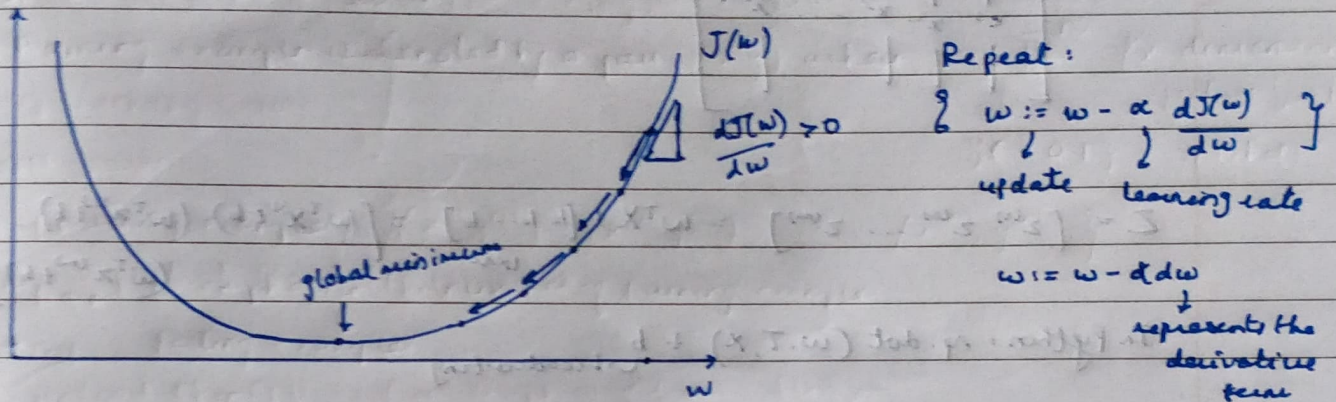


Loss function is defined for a single training example, cost function is defined for a set of training examples (cost of the parameters).

$$\text{Cost Function: } J(w, b) = \frac{1}{n} \sum_{i=1}^n L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{n} \sum_{i=1}^n [y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})]$$

→ Since we want the cost function to be minimum, we find  $w, b$  which minimize the cost function.

For this we use the concept of gradient descent.



$$J(w, b) \Rightarrow \left. \begin{array}{l} w := w - \alpha \frac{\partial J(w, b)}{\partial w} \\ b := b - \alpha \frac{\partial J(w, b)}{\partial b} \end{array} \right\} \text{for parameters } w \text{ and } b \text{ is the cost } J$$

→ Computation graph

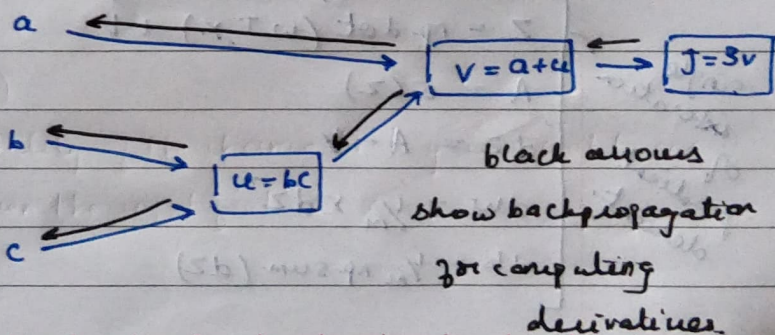
Example:

$$J(a, b, c) = 3(a + bc)$$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$





→ **Vectorization**: it is the art of removing explicit for loops from your code. In deep learning code, if there are more for loops, it slows down the program hence decreasing efficiency, so vectorization is used.

In Python: `np.dot(w, b)` → vectorization.

→ **Vectorising logistic regression**: it is a machine learning algorithm that uses logistic regression.

$$X = \begin{bmatrix} | & | & | & \dots & | \\ x^1 & x^2 & x^3 & \dots & x^{(n)} \\ | & | & | & \dots & | \end{bmatrix}$$

$$Z = [z^{(1)} \ z^{(2)} \ \dots \ z^{(n)}] = W^T X + [b \ b \ \dots \ b] = [(w^T x^{(1)} + b) \ (w^T x^{(2)} + b) \ \dots \ (w^T x^{(n)} + b)]$$

In Python: `np.dot(w.T, X) + b` → broadcasting

→ **Vectorising logistic Regression's gradient descent**.

$$dZ = [dz^{(1)} \ dz^{(2)} \ \dots \ dz^{(n)}]$$

$$db = \frac{1}{n} \sum_{i=1}^n dz^{(i)}$$

In Python: `db = 1/n np.sum(dZ)`

$$dW = \frac{1}{n} X dZ^T$$

→ **Implementing logistic Regression**:

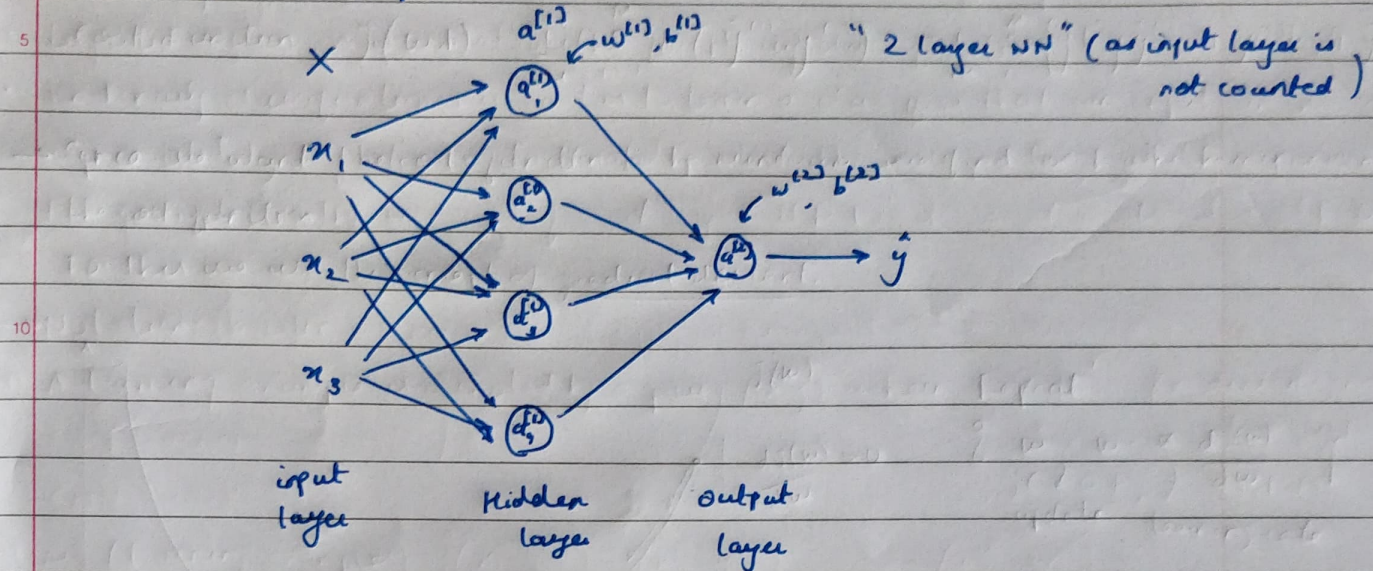
single iteration of gradient descent.

$$\left[ \begin{array}{l} Z = np.dot(w.T, X) + b \\ A = \sigma(Z) \\ dZ = A - Y \\ dW = \frac{1}{n} X dZ^T \\ db = \frac{1}{n} np.sum(dZ) \end{array} \right. \quad \left[ \begin{array}{l} w := w - \alpha dW \\ b := b - \alpha db \end{array} \right]$$



### \* Week 3

#### → Neural Network Representation:



- Considering the first node of the hidden layer:

$$z_1^{(1)} = w_1^{(1)T} X + b_1^{(1)}$$

$$\rightarrow a_1^{(1)} = \sigma(z_1^{(1)})$$

$$\left( \begin{matrix} a_i^{(l)} \rightarrow \text{layer} \\ a_i \rightarrow \text{node in layer} \end{matrix} \right)$$

#### → Formulas for computing derivatives

- Forward Propagation

$$z^{(1)} = w^{(1)} X + b^{(1)}$$

$$A^{(1)} = g^{(1)}(z^{(1)})$$

$$z^{(2)} = w^{(2)} A^{(1)} + b^{(2)}$$

$$A^{(2)} = g^{(2)}(z^{(2)}) = \sigma(z^{(2)})$$

- Back Propagation

$$dz^{(2)} = A^{(2)} - Y \quad Y = [y^{(1)} y^{(2)} \dots y^{(n)}]$$

$$dw^{(2)} = \frac{1}{n} dz^{(2)} A^{(1)T}$$

$$db^{(2)} = \frac{1}{n} \text{np.sum}(dz^{(2)}, \text{axis}=1, \text{keepdims}=\text{True})$$

$$dz^{(1)} = \frac{w^{(2)T} dz^{(2)}}{(n^{(1)}, m)} + \underbrace{g^{(1)'}(z^{(1)})}_{(n^{(1)}, m)} \text{element wise}$$

$$dw^{(1)} = \frac{1}{n} dz^{(1)} X^T$$

$$db^{(1)} = \frac{1}{n} \text{np.sum}(dz^{(1)}, \text{axis}=1, \text{keepdims}=\text{True})$$