

1.1

Introduction:

The aim of this project is to implement the artificial neural network algorithm to solve the multiclass (10 class) classification problem for a given dataset.

This assignment teaches us how the algorithms written in the sklearn library actually works and in depth mathematical intuition behind these algorithms.

The dataset used in this project consists of six independent features i.e. X1, X2 ,X3 ,X4 ,X5 ,X6 and a dependent feature Y. Our goal is to use these independent features to predict the values of dependent feature

1.2 Code Design, Methodology and Derivations:

The code can be divided into following subsections for clear understanding:

A. Importing Required Libraries:

You need to install and import numpy, pandas and matplotlib libraries to successfully run this project

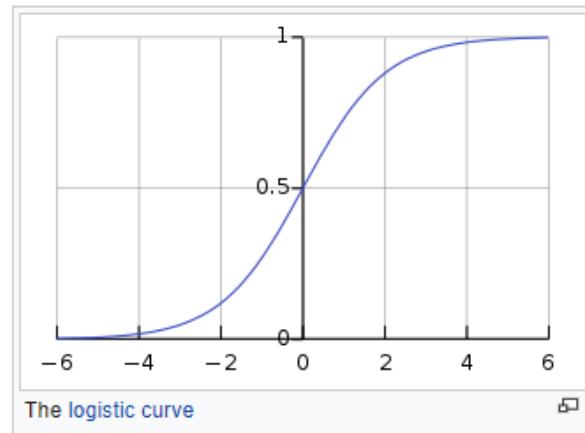
B. Data Preprocessing:

1. First we are reading the csv file using the pandas library
2. Then we are normalizing the data set using
 $X \rightarrow (X - \text{mean})/(\text{standard deviation})$
3. Shuffling the data
4. Performing test train split so that we use 70% of the data for training and 30% of the data for testing.

C. Data Learning With Two Layer Artificial Neural Network

Firstly we have defined the **sigmoid function**. Mathematical formula for sigmoid function is

$$S(x) = \frac{1}{1 + e^{-x}}$$



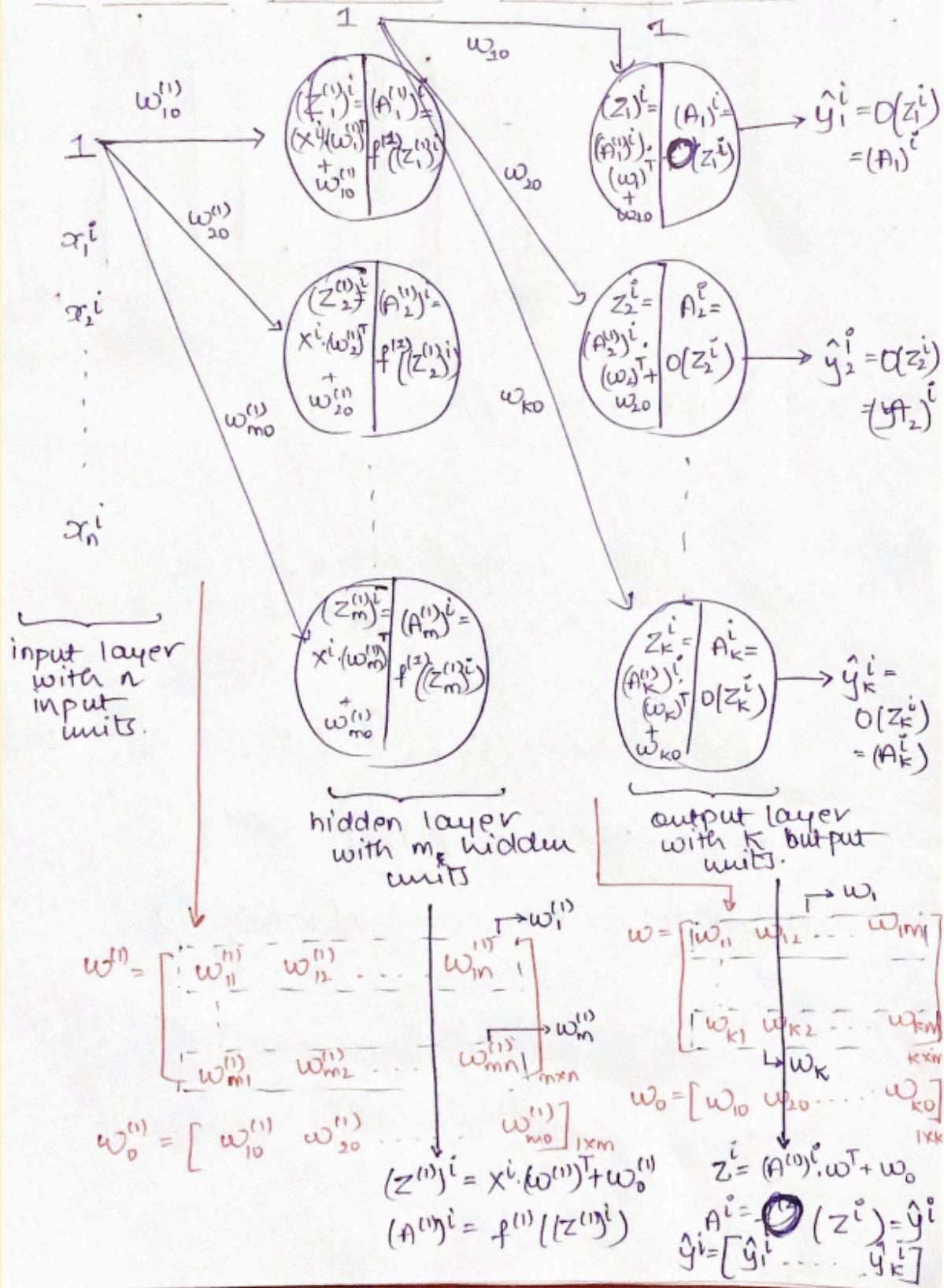
Next we defined the **derivative of the sigmoid function**. Mathematical formula for first derivative of sigmoid function is given by

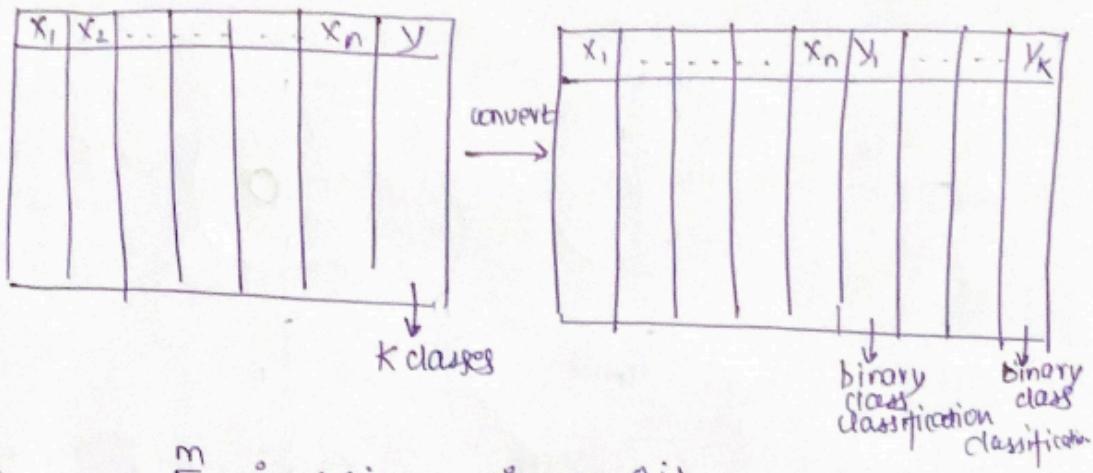
$$S'(x) = S(x) * (1 - S(x))$$

Now a function called **stochastic gradient descent** is defined which is the core of the entire algorithm. It takes the data, learning rate, no of iterations and no of hidden units as inputs. In this we have used stochastic gradient descent algorithm to optimize the unconstrained optimization problem formed in ANN. Since we are using stochastic gradient descent algorithm the cost function due to just one example.

This function can be further divided into following parts , namely: forward propagation , backward propagation and weight updation. The formulas used in all these can be understood by understanding the complete math intuition of the algorithm provided in notes below and the coding is done using similar notations for easy understanding.

2 LAYER NETWORK - ARCHITECTURE





$$\text{Error} = \sum_{i=1}^m -y_i^i \ln(\hat{y}_i^i) - (1-y_i^i) \ln(1-\hat{y}_i^i) + \\ \dots + \sum_{i=1}^m -y_k^i \ln(\hat{y}_k^i) - (1-y_k^i) \ln(1-\hat{y}_k^i)$$

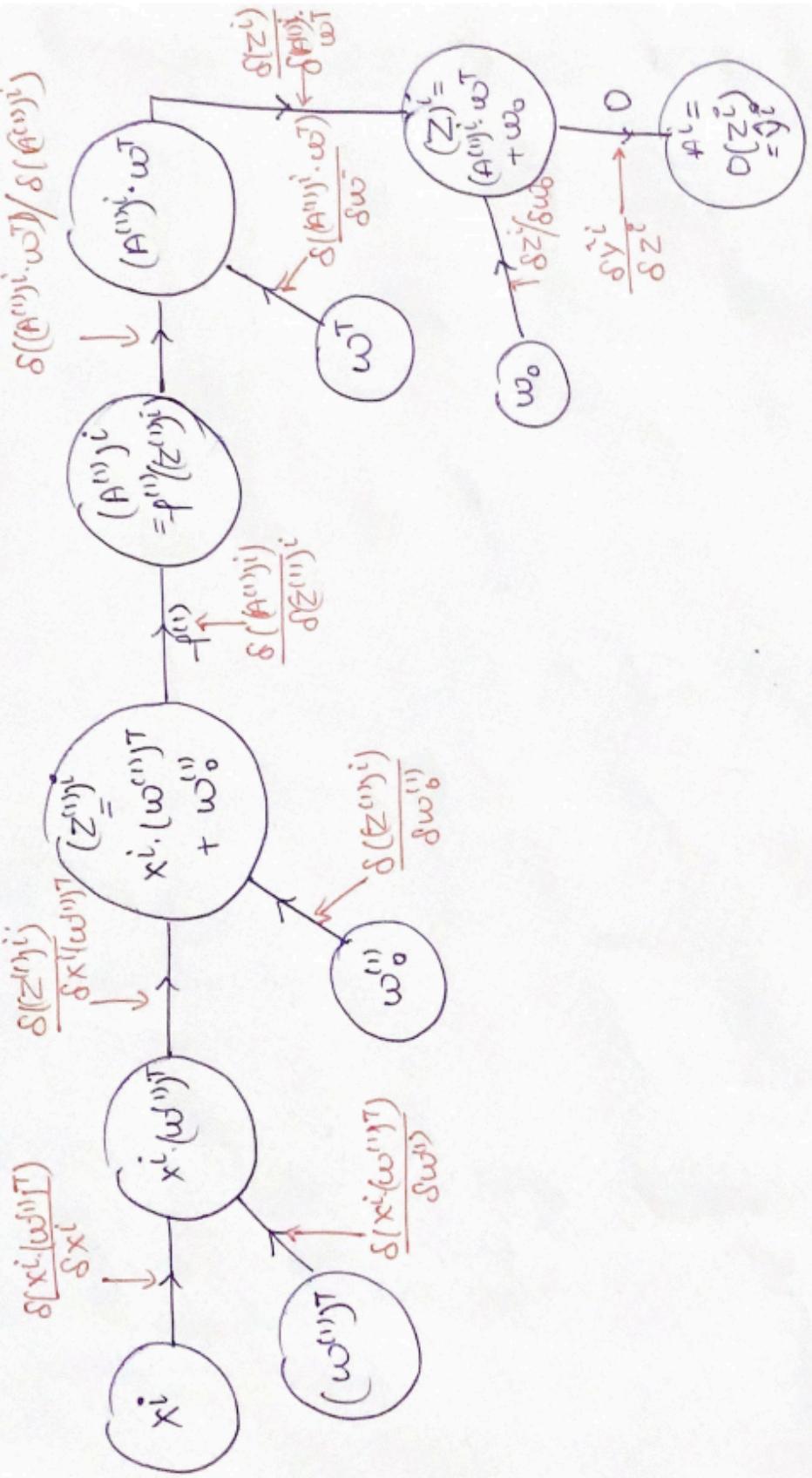
Stochastic Error = $-y_1^i \ln(\hat{y}_1^i) - (1-y_1^i) \ln(1-\hat{y}_1^i) + \dots + -y_k^i \ln(\hat{y}_k^i) - (1-y_k^i) \ln(1-\hat{y}_k^i)$

using matrix form
 $y^i = [y_1^i \ y_2^i \ \dots \ y_k^i], \hat{y}^i = [\hat{y}_1^i \ \hat{y}_2^i \ \dots \ \hat{y}_k^i]$

$$= -y^i \ln((\hat{y}^i)^T) - (1-y^i) \ln(1-(\hat{y}^i)^T)$$

$$= -y^i \ln(O(z^i)^T) - (1-y^i) \ln(1-O(z^i)^T)$$

COMPUTATION GRAPH FOR ERROR FUNCTION.



CALCULATING PARTIAL DERIVATIVES

$$\begin{aligned} \boxed{\frac{\delta \text{Error}}{\delta w_0}} &= \frac{\delta \text{Error}}{\delta y^i} \cdot \frac{\delta y^i}{\delta z^i} \cdot \frac{\delta z^i}{\delta w_0} \\ &= \cancel{\frac{\delta y^i}{\delta z^i}} \cdot \frac{\delta \text{Error}}{\delta z^i} \cdot \frac{\delta z^i}{\delta w_0} \\ &= \left(\frac{-y^i}{O(z^i)} \cdot O'(z^i) - \frac{(1-y^i)}{1-O(z^i)} \cdot O'(1-O(z^i)) \right) \cdot (1) \end{aligned}$$

$$\begin{aligned} \boxed{\frac{\delta \text{Error}}{\delta w}} &= \frac{\delta \text{Error}}{\delta z^i} \cdot \frac{\delta z^i}{\delta ((A^{(1)})^i \cdot w^T)} \cdot \frac{\delta ((A^{(1)})^i \cdot w^T)}{\delta w} \\ &= \cancel{\frac{\delta \text{Error}}{\delta z^i}} \left| \frac{\delta \text{Error}}{\delta w_0} \cdot (1+0) \cdot ((A^{(1)})^i) \right| \end{aligned}$$

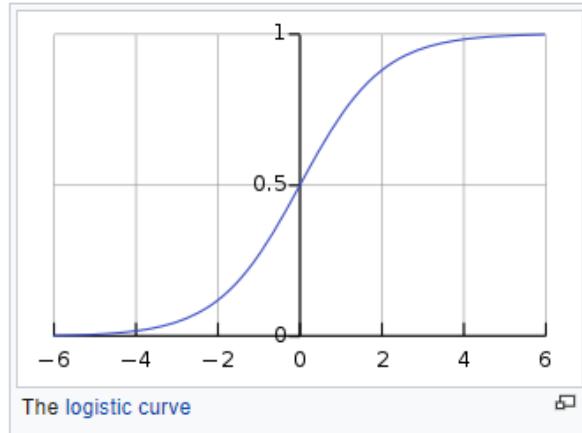
$$\begin{aligned} \boxed{\frac{\delta \text{Error}}{\delta w_0^{(1)}}} &= \frac{\delta \text{Error}}{\delta z^i} \cdot \frac{\delta z^i}{\delta ((A^{(1)})^i \cdot w^T)} \cdot \frac{\delta ((A^{(1)})^i \cdot w^T)}{\delta ((A^{(1)})^i)} \cdot \frac{\delta ((A^{(1)})^i)}{\delta (z^{(1)})^i} \\ &\quad \cdot \frac{\delta (z^{(1)})^i}{\delta (w_0^{(1)})} \\ &= \boxed{\left(\frac{\delta \text{Error}}{\delta w_0} \right) (1+0) (w^T) \cdot (-l^{(1)}(z^{(1)})^i) \cdot (0+1)} \end{aligned}$$

$$\begin{aligned} \boxed{\frac{\delta \text{Error}}{\delta w^{(1)}}} &= \frac{\delta \text{Error}}{\delta z^i} \cdot \frac{\delta z^i}{\delta ((A^{(1)})^i \cdot w^T)} \cdot \frac{\delta ((A^{(1)})^i \cdot w^T)}{\delta ((A^{(1)})^i)} \cdot \frac{\delta ((A^{(1)})^i)}{\delta (z^{(1)})^i} \\ &\quad \cdot \frac{\delta (z^{(1)})^i}{\delta (x^i \cdot w^{(1)})^T} \cdot \frac{\delta (x^i \cdot w^{(1)})^T}{\delta w^{(1)}} \\ &= \boxed{\left(\frac{\delta \text{Error}}{\delta w_0} \right) (1+0) \cdot (x^i)} \end{aligned}$$

D. Data Learning Using Three Layer Artificial Neural Network

Firstly we have defined the **sigmoid function**. Mathematical formula for sigmoid function is

$$S(x) = \frac{1}{1 + e^{-x}}$$



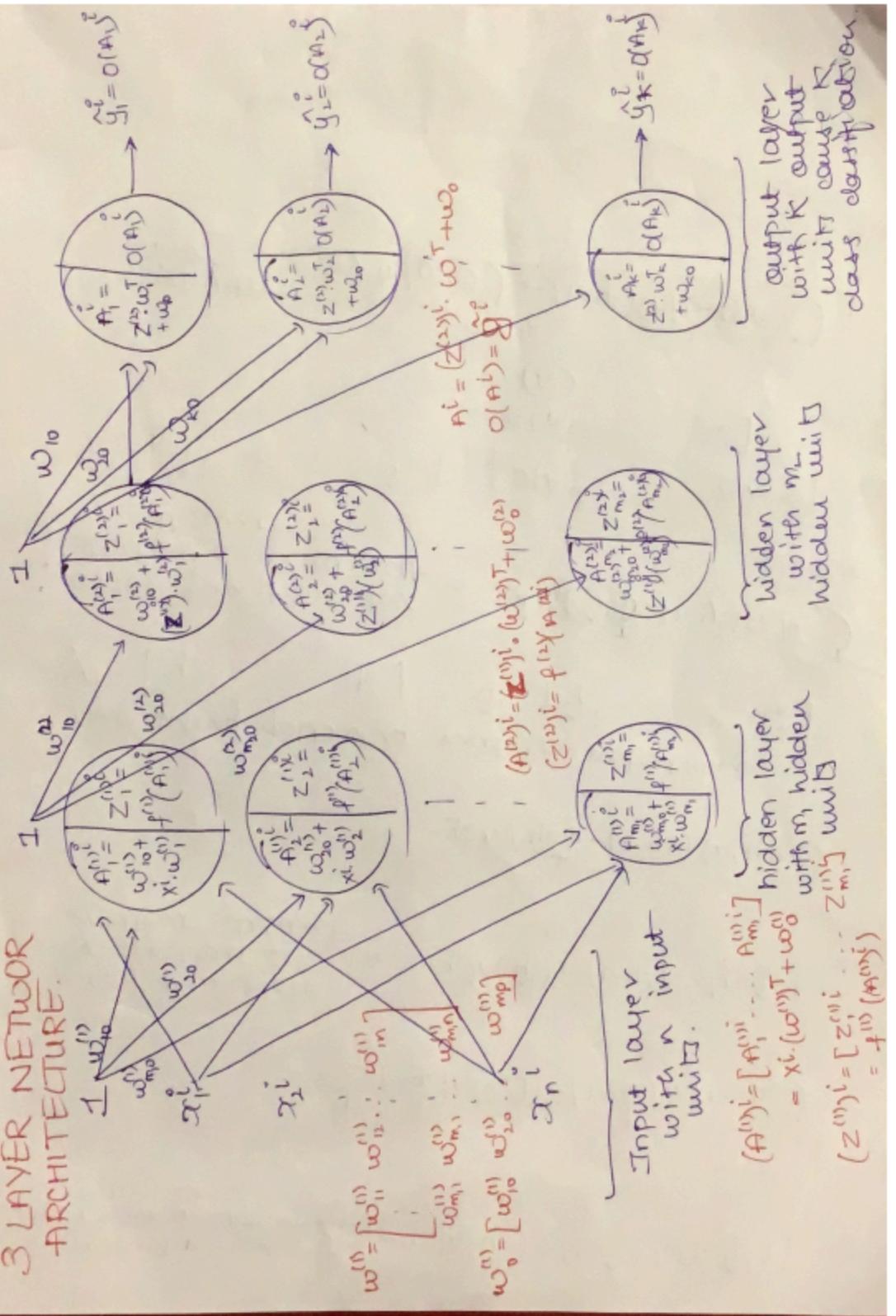
Next we defined the **derivative of the sigmoid function**. Mathematical formula for first derivative of sigmoid function is given by

$$S'(x) = S(x) * (1 - S(x))$$

Now a function called **stochastic gradient descent** is defined which is the core of the entire algorithm. It takes the data, learning rate, no of iterations and no of hidden units in hidden layer 1 and no of hidden units in input layer 2 as inputs. In this we have used stochastic gradient descent algorithm to optimize the unconstrained optimization problem formed in ANN. Since we are using stochastic gradient descent algorithm the cost function due to just one example.

This function can be further divided into following parts , namely: forward propagation , backward propagation and weight updation. The formulas used in all these can be understood by understanding the complete math intuition of the algorithm provided in notes below and the coding is done using similar notations for easy understanding.

3 LAYER NETWERK ARCHITECTURE.



$$\therefore \text{error function for gradient descent} = \sum_{i=1}^n (-y_i^i \ln(\hat{y}_i^i) - (1-\hat{y}_i^i) \ln(1-\hat{y}_i^i))$$

+

$$\sum_{k=1}^m (-y_k^i \ln(\hat{y}_k^i) - (1-\hat{y}_k^i) \ln(1-\hat{y}_k^i))$$

$$\therefore \text{error function for stochastic gradient descent} = -y_i^i \ln(\hat{y}_i^i) - (1-y_i^i) \ln(1-\hat{y}_i^i)$$

+

$$-y_k^i \ln(\hat{y}_k^i) - (1-y_k^i) \ln(1-\hat{y}_k^i)$$

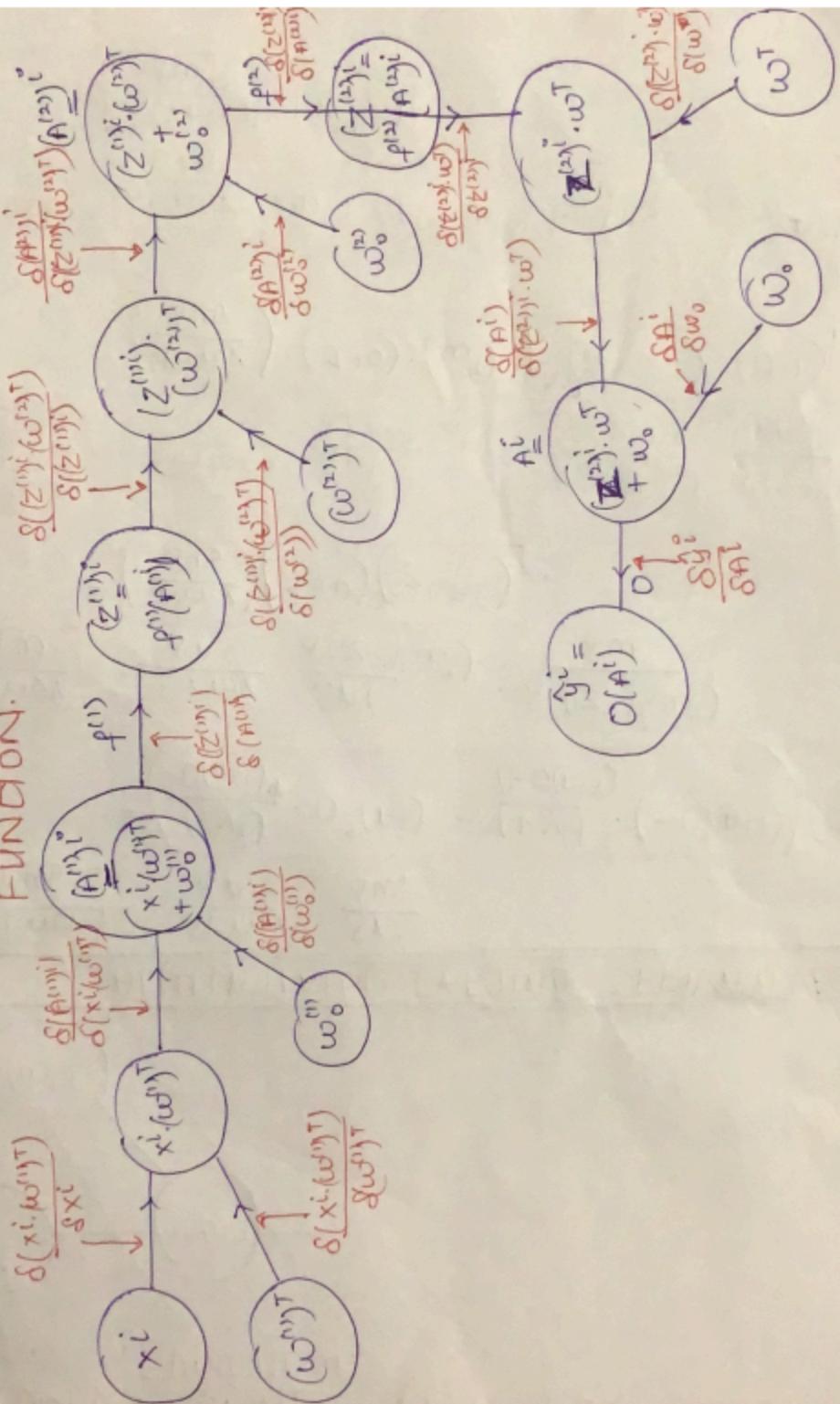
now using following matrix notation
 $y^i = [y_1^i \dots y_K^i]$, $\hat{y}^i = [\hat{y}_1^i \dots \hat{y}_K^i]$

$$\therefore \text{error function for stochastic gradient descent} = -y_i^i \ln(\hat{y}^i)^T - (1-y_i^i) \ln(1-\hat{y}^i)^T$$

$$\begin{aligned} \hat{y}^i &= [y_1^i \ y_2^i \dots \hat{y}_K^i] = [O(A_1)^i \dots O(A_K)^i] \\ &= O[A_1^i \dots A_K^i] \\ &= O(A^i) \end{aligned}$$

$$\boxed{\text{Stochastic Error.} = -y_i^i \ln(O(A^i)^T) - (1-y_i^i) \ln(1-(O(A^i))^T)}$$

COMPUTATIONAL GRAPH FOR ERROR FUNCTION.



CALCULATING PARTIAL DERIVATIVES:

$$\boxed{\frac{\delta \text{Error}}{\delta w_0}} = \frac{\delta \text{Error}}{\delta A^i} \cdot \frac{\delta A^i}{\delta w_0}$$

$$= \left(\frac{(-y^i)}{O(A^i)} \cdot O^9(A^i) - \frac{(1-y^i)}{(1-O(A^i))} \cdot (-O^9(A^i)) \right) (0+1)$$

$$\boxed{\frac{\delta \text{Error}}{\delta w}} = \frac{\delta \text{Error}}{\delta A^i} \cdot \frac{\delta A^i}{\delta (Z^{(2)})^i \cdot w^T} \cdot \frac{\delta ((Z^{(2)})^i \cdot w^T)}{\delta w}$$

$$= \left(\frac{\delta \text{Error}}{\delta w_0} \right) (1+0) ((Z^{(2)})^i)$$

$$\boxed{\frac{\delta \text{Error}}{\delta w_0^{(2)}}} = \frac{\delta \text{Error}}{\delta A^i} \cdot \frac{\delta A^i}{\delta ((Z^{(2)})^i \cdot w^T)} \cdot \frac{\delta ((Z^{(2)})^i \cdot w^T)}{\delta (Z^{(2)})^i} \cdot \frac{\delta (Z^{(2)})^i}{\delta (A^{(2)})^i} \cdot \frac{\delta f^{(2)}}{\delta w_0^{(2)}}$$

$$= \left(\frac{\delta \text{Error}}{\delta w_0} \right) \cdot (1+0) \cdot (w^T) \cdot (f^{(2)}(A^{(2)i}) \cdot 1) \cdot (0+1)$$

$$\boxed{\frac{\delta \text{Error}}{\delta w^{(2)}}} = \left[\frac{\delta \text{Error}}{\delta w_0} \cdot (w^T) \cdot (f^{(2)}(A^{(2)i}) \cdot 1) (1+0) \cdot ((Z^{(2)})^i) \right]$$

$$= \left(\frac{\delta \text{Error}}{\delta w_0^{(2)}} \right) \cdot ((Z^{(2)})^i)$$

$$\boxed{\frac{\delta \text{Error}}{\delta w_0^{(1)}}} = \left[\frac{\delta \text{Error}}{\delta w_0} \cdot (w^T) \cdot (f^{(1)}(A^{(1)i}) \cdot 1) \cdot (1+0) \cdot (w^{(2)})^T \cdot \right.$$

$$\left. (f^{(1)}(A^{(1)i}) \cdot 1) \cdot (0+1) \right]$$

$$\boxed{\frac{\delta \text{Error}}{\delta w^{(1)}}} = \left(\frac{\delta \text{Error}}{\delta w_0^{(1)}} \right) \cdot (1+0) \cdot (x^i)$$

