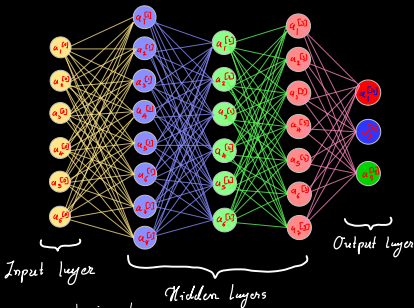


Forward Propagation & Backward Propagation in Neural Networks



where

$a_n^{[m]}$

n : neuron number in a layer

m : layer number of a neuron

$$a_1^{[2]} = af(w_{11} \times a_1^0 + w_{12} \times a_2^0 + w_{13} \times a_3^0 + w_{14} \times a_4^0 + w_{15} \times a_5^0 + w_{16} \times a_6^0 + \beta)$$

Forward Propagation

$$\begin{aligned} Z_1 &= w_1 \times A_0 + \beta_1 & Z_2 &= w_2 \times A_1 + \beta_2 & Z_3 &= w_3 \times A_2 + \beta_3 \\ A_1 &= f(Z_1) & \rightarrow A_2 &= f(Z_2) & \rightarrow A_3 &= f(Z_3) = \text{O/p prediction} \end{aligned}$$

Process

Assign random weight (w) at initially



Model Training



Update the value of weights in such a way so that it can predict more accurately.

Cost function



$$W_{\text{new}} = W_{\text{old}} - \alpha * \frac{\partial \text{cost}}{\partial w}$$

Where

α = learning rate

$\frac{\partial \text{cost}}{\partial w}$ = slope

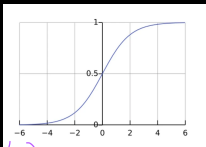
Activation functions

1. Sigmoid function

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

→ Output range 0 to 1 hence we can use it on output layer.

(Binary classification)



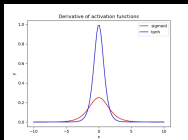
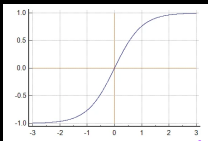
2. Tanh function

$$\tanh = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

→ Also called as hyperbolic tangent function.

→ Range -1 to 1.

→ We need to take derivative of activation functions in order to calculate derivative of cost function.



→ As this image shows the derivatives of both function.

→ Here tanh has more derivative than sigmoid that's why by using tanh in hidden layers of NN, we can accelerate the training process much faster.

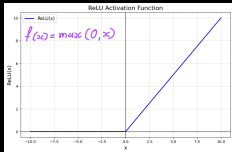
→ In tanh, average of data is close to '0' (data is normalized centre around '0')

→ Hence if we pass normalized data to next layer, it makes training much more easier.

Note: Both sigmoid & tanh has smaller derivative value at both corner side. this will leads to slow learning also called as vanishing gradient problem. this problem can be solved by next activation function

3. ReLU (Rectified Linear Unit)





→ Overcomes the Vanishing Gradient Problem

$$\rightarrow \frac{\partial f(x)}{\partial x} = \begin{cases} 1, & x > 0 \\ 0, & x < 0 \end{cases}$$

→ also called as Piece-wise linear.

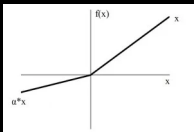
→ Advantages of both linear & non-linear property.

→ linearity of a function can help to overcome vanishing gradient problem and

also makes training faster.

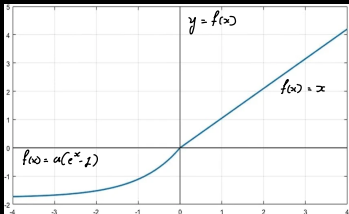
Variations of ReLU

[a] Leaky ReLU



$$f(x) = \max(0.01 * x, x)$$

[b] ELU [Exponential Linear Unit]



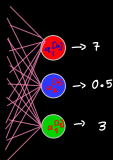
$$f(x) = a(e^x - 1)$$

$$f(x) = \max(a * (e^x - 1), x)$$

4. Softmax function

→ For multi-class classification.

first method...



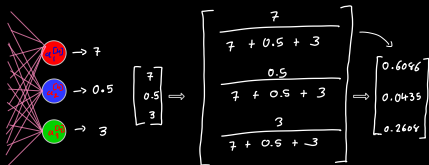
converting into probability

$$\begin{bmatrix} 7 \\ 0.5 \\ 3 \end{bmatrix} \Rightarrow \begin{bmatrix} e^7 \\ e^{0.5} \\ e^3 \end{bmatrix} \Rightarrow$$

$$\begin{bmatrix} \frac{e^7}{e^7 + e^{0.5} + e^3} \\ \frac{e^{0.5}}{e^7 + e^{0.5} + e^3} \\ \frac{e^3}{e^7 + e^{0.5} + e^3} \end{bmatrix} \Rightarrow \begin{bmatrix} 0.9782 \\ 0.00147 \\ 0.0179 \end{bmatrix}$$



Second method.



→ We can also calculate probability as shown in above, but by using exponential method, it will expand the value hence we can get more accuracy.

$$\text{Softmax } f_i(x) = \frac{e^{u_i}}{\sum_k e^{u_k}}$$

Note: If the op of neural network is continuous in this case, it is not suitable to use any kind of non-linear activation function at output neuron.

→ In this situation, we will just take the linear op without activation function.

Cost functions

Regression /
Linear regression

Binary Classification

Multi-class classification.

1. Regression problem

↳ ex. House price prediction

$$\text{cost} = \frac{1}{n} \sum_{i=1}^m |\hat{y}_i - y_i| \quad \text{on} \quad \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 \quad \text{Loss} = |\hat{y}_i - y_i|$$

Mean Absolute Error Mean Squared Error

for i^{th} observation



Created with
Notewise

2. For binary classification

$$\text{error} = -[y_i \cdot \log(u_i) + (1-y_i) \log(1-u_i)] \text{ for 1 observation}$$

$$\text{cost} = -\frac{1}{m} \sum [y_i \cdot \log(u_i) + (1-y_i) \log(1-u_i)]$$

Note if $u_i \rightarrow 0$ then error will be \downarrow

$u_i \rightarrow 1$ then error will be \uparrow

$u_i \rightarrow y_i$ then we will get less error.

! $u_i \rightarrow y_i$ then we will get more error.

→ This method is also called as 'Binary Cross Entropy'.

3. for multi-class classification

Suppose

$$y_i = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, u_i = \begin{bmatrix} 0.1 \\ 0.5 \\ 0.3 \\ 0.2 \end{bmatrix} \left| \begin{array}{l} \text{error} = \\ [y_{i2} \cdot \log(u_{i2}) + y_{i2} \cdot \log(u_{i2}) + y_{i3} \cdot \log(u_{i3}) + y_{i4} \cdot \log(u_{i4})] \\ \vdots \\ -\frac{[y_{i2} \cdot \log(u_{i2})]}{m} = -\frac{\log(0.5)}{m} \end{array} \right.$$

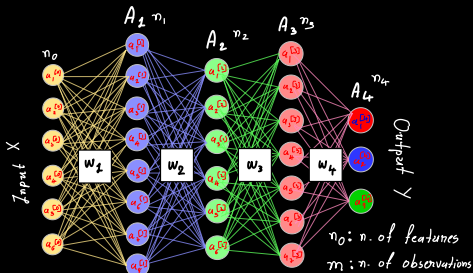
$$\text{Cost} = -\sum_{j=0}^M \sum_{i=0}^N (y_{ij} * \log(u_{ij}))$$

m : Total number of neurons in o/p layer.

n : Total number of observations.



Derivation of back propagation



Shape of $W_1 = (n_1, n_0)$ Shape of $B_1 = (n_1, 1)$ Shape of $X = (n_0, m)$

Shape of $W_2 = (n_2, n_1)$ Shape of $B_2 = (n_2, 1)$ Shape of $A_1 = (n_1, m)$

Shape of $W_3 = (n_3, n_2)$ Shape of $B_3 = (n_3, 1)$ Shape of $A_2 = (n_2, m)$

Shape of $W_4 = (n_4, n_3)$ Shape of $B_4 = (n_4, 1)$ Shape of $A_3 = (n_3, m)$

Shape of $A_4 = (n_4, m)$

Shape of $Y = (n_4, m)$

Activation functions

Form of layer

$i = [1, 2, 3]$ (ReLU / tanh)
 $A_i = f_i(z_i)$

A_4 (Sigmoid / Softmax)

$$Z_k = W_k \cdot A_{k-1} + B_k$$

Process of backward Propagation

$$\frac{\partial L}{\partial W_k}$$

$L = \text{Loss of } i^{\text{th}} \text{ data point}$
 $L \rightarrow a_4 \rightarrow z_4 \rightarrow W_4$

$$\text{loss} = -[y_i \cdot \log(u_i) + (1 - y_i) \log(1 - u_i)]$$

$a_4 = \text{Sigmoid}(z_4) \rightarrow z_4 = W_4 \cdot A_3 + B_4$

$$\frac{\partial L}{\partial W_4} = \frac{\partial L}{\partial a_4} \cdot \frac{\partial a_4}{\partial z_4} \cdot \frac{\partial z_4}{\partial W_4}$$



Steps of NN training

1. Initialize parameters randomly.
2. Repeat the process
 - Forward Propagation
 - Activation function
 - Calculating cost
 - Backward Propagation
 - Updating weights & bias

