

H. MultiLayer perception

multilayer perceptron also called feed forward network

this is a N/w that consists of one or more hidden layers

so that I can perform complex tasks

we can also perform transformation (e.g XOR gate)

(Intention) Existence of hidden layers we can perform complex tasks

exit hand written digits recognition

This one called as feed forward N/w because IP signal is forwarded from next layer then next layer has only 1 direction the data is forwarded.

Hence it is called feed forward Neural N/w.

It has one or more hidden layers.

Here it has one or more hidden layers.

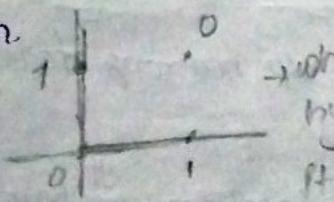
In the sketch we have MLP classifier to create Multilayer Perception.

Advantages (of MLP over single layer perceptron)

The MLP can perform complex tasks.

Perform transformation.

e.g. XOR gate



when it fails this task even though the classification is binary.

it can only classify in 2D

Importance of Hidden layers

Having hidden layers let MLP

① Non-linear transformations so that the data is linearly separable

② can change dimensionality of data

③ Extract low level features

④ Compute gradient (used in error optimization)

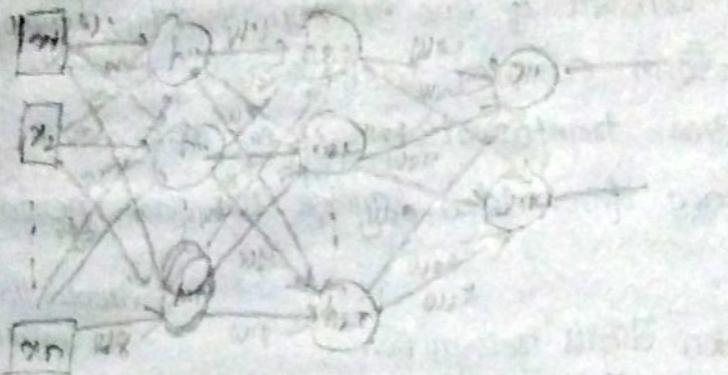
horizontal lines -

1st (even) hidden layer
V. L → 2nd level

curve - L → 3rd levels
of an image

Structure of Multi Layer perception

Input layer hidden layer 1 hidden layer 2 ... hidden layer N OLP Layer



- Every layer is well connected.
- The OLP layer passes data to hidden layers.
- The data from hidden layer 1 is passed to higher hidden layers.
- The hidden layers process the model to get the OLP.
- Each & every OLP is weighted.
- NO of features = NO of OLPs
- If we have 1 OLP neuron - only binary classification
2 OLP " " - Multiclass "
- If we have the no. of OLP neurons depends on the no. of features you have.
- In sklearn the MLP classifier has 100 hidden layer neurons by default.
- We need to compute the error at the OLP layer.

* Conventions



- w_{ij} - weight of the OLP from neuron i to j
- y_j - actual OLP of neuron j
- d_j - target OLP of neuron j (y_i) associated with neuron j

e_j - error value for neuron j
 a_j - Activation funⁿ for j th neuron.
 v_j - induced local field value for neuron j .
 φ_j - Multilayer perceptron

$$y_j = \varphi_j(v_j)$$

$$v_j = \sum_{q=1}^n (w_{jq} + x_q b)$$

$$e_j = d_j - y_j$$

$$\text{cost fun}^n c(w_{j1}) = \frac{1}{N} \sum_{q=1}^n e_j^q$$

$$c_{avg}(w_{j1}) = \frac{1}{N} \sum_{q=1}^n e_j^q$$

N = total no of samples
 n = the n th sample taken

Implementing MLP

Scklearn
 ↳ neural-network → MLP classifier

→ we can use constructor of MLP classifiers.

1st parameter $(100,)$ → 100 neurons

1 hidden layer - sizes = tuple
 $= (10, 5, 3)$ ↳ 70% (8) 90% used
 $\downarrow \quad \downarrow \quad \downarrow$
 $H_1 \quad H_2 \quad H_3$

(hidden layer with 10 neurons)

→ The 1st hidden layer - no. of neurons = 70% of inputs (80) 90% of filp features.

→ 2nd layer hidden ⇒ 2nd parameter

activation ↳ identity
 ↳ logistic
 ↳ tanh
 ↳ relu

3rd parameter - solver
 → we can specify optimization algorithm by using solvers

→ 1 hfgs - Neurons
 sgd - stochastic gradient descent
 adam - variant of sgd.

4th parameter - max_iter

→ it specifies the max no of iterations

→ default 100 has 200 iterations

→ tol → tolerance, tol = 10^{-4}

5th parameter → learning rate 0.01.

(diabetes)

Program

from sklearn import datasets

d = datasets.load_breast_cancer()

d.data # filp features

dfit array ([[5 +, 3 +, 0 -]])

d. target & target feature o/p :- array ([0, 0, 0, 0, ...])

d. feature_names

['Sepal length (cm)',
'Sepal width (cm)',
'Petal length (cm)',
'Petal width (cm)']

Import pandas as pd # to create a dataframe

from sklearn.model_selection import train_test_split

from sklearn.neural_network import MLPClassifier

df = pd.DataFrame(d.data, columns = ['Sepal length', 'S.W', 'P.L', 'P.W'])

df.head()

df.info() # non-null values o/p :- 150 entries with 4 columns no null values

df.columns # names of columns

X = df ['S.L', 'S.W', 'P.L', 'P.W']

y = d.target # target labels

Xtr, Xte, Ytr, Yte = train_test_split(X, y, test_size=0.2, random_state=0)

mc = MLPClassifier(hidden_layer_sizes=(4, 3), max_iter=500, activation = "tanh", solver="sgd")

mc

d.target_names # o/p :- array(['Setosa', 'Versicolor', 'V. Virginica'])

mc.fit(Xtr, Ytr)

logistic - 0.2

mc.predict(Xte) o/p :- array([2, 1, 0, 2, ...])

tanh - 0.56

mc.score(Xte, Yte) # 0.9666 (linearly separable)

* Back Propagation Algorithm

- First, we need to compute the total error value which we call as

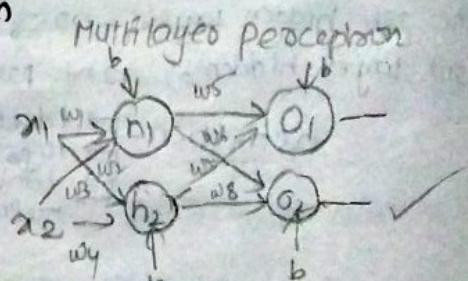
E_{total} .

- The actual o/p is produced at o/p layer.

- We need to minimize the total error value by updating the free parameters (w, b).

- When we initialise the wt & bias with random values we may get the line that not separate the data. So we need to update until error is reduced.

- We need to update all the weight values $0.60 + 0.8 \cdot x_2 + b \rightarrow$
but all is not same possible \leftarrow linear line
~~not separate the~~



- at same time.
- so first we need to update the g/p's associated with o/p layer
- then the hidden layer g/p's will be updated.
- one of the reason for getting error is wrong values
- if we finish the f the hidden layer then we need to move through another hidden layers.
- Back propagation alg make use of gradient descent algorithm.
- It consists of two phase
 - 1) forward pass
 - 2) backward pass.

forward Pass:

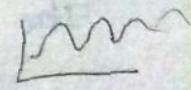
- N/w fed with the g/p
- neurons in (different) layers process the g/p. and forward them to the next layers.

backward pass

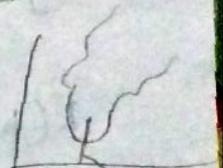
- The total error value is taken computed
- When there is a error
 - uses gradient descent to minimize the error value

↳ it updates the wt, bias by computing the gradient of cost fun?

$$\frac{1}{2} \sum_{i=1}^n e_i^2 \rightarrow \text{convex fun}$$



- when we derivative eqⁿ → we will get slope
- ↳ it will in reverse direction slope
- ↳ it will move to global minimum in ~~reverse~~ opposite direction of slope.



global minimum

$$w(n+1) = w(n) - \eta g(n)$$

↳ Back propagation also make use of this simple rule.

But will make use of different procedures for g/p, e, dr

layer

- Let us see how back propagation alg will produce the error

Computing gradient at

OIP Layer

$$E_{\text{total}} = E_{O1} + E_{O2} + \dots + E_{On}$$

Eqn 1 (E_{O1} over E_{O2})

$$E_{\text{total}} = E_{O1} + E_{O2}$$

$$E_{O1} = \frac{1}{2} \sum_{i=1}^m e_i^2 \quad \text{at filone}$$

$$E_{O2} = \frac{1}{2} \sum_{i=1}^m e_i^2 = \frac{1}{2} \sum_{i=1}^m (d_{O1} - y_{O1})^2 \quad m = \text{size of training set}$$

then

$$\boxed{\text{ad q122}} \quad E_{\text{total}} = \frac{1}{2} \sum_{i=1}^m (d_{O1} - y_{O1})^2 + \frac{1}{2} \sum_{i=1}^m (d_{O2} - y_{O2})^2$$

⇒ when the OIP is not correct then we call that situation is called error situation.

⇒ OIP layer neuron j :

$$e_j = d_j - y_j - ①$$

$$y_j = \psi(v_j) \quad v_j = \text{induced local field value at } j\text{th neuron}$$

$$= \text{Sigmoid}(v_j) = \frac{1}{1+e^{-v_j}} - ②$$

The OIP y depends upon the activation funⁿ we use

$$v_j = \sum_{i=1}^m w_{ji} x_i + b_0 \quad w_{ji} = i\text{th weight for the}$$

$$\hookrightarrow ③ \quad b_0 = b_{\text{bias}} \quad \text{fifth I/P}$$

$$E(w) = \frac{1}{2} \sum_{i=1}^n e_i^2 - ④$$

$$= \frac{1}{2} \sum_{i=1}^n (d_i - y_i)^2$$

⇒ Update Rule

as per gradient descent

$$w_{(n+1)} = w_{(n)} - \eta g(n) \quad (2) \quad w_{(n+1)} = w_{(n)} - \eta g(n)$$

$g(n)$ is gradient of cost funⁿ

Let us assume,
a neural Net with K neurons in output layer

$$E_{\text{total}}(w) = \sum_{j=1}^K E_j(w) = E_1(w) + E_2(w) + \dots + E_K(w) \quad (5)$$

$E_k(w) \rightarrow$ error or value at
10-th neuron

gradient of cost funⁿ

$$g(n) = \frac{\partial E_{\text{total}}}{\partial w_{ji}} \quad (\text{using derivative chain rule})$$

$$= \frac{\partial E}{\partial e} \times \frac{\partial e}{\partial v} \times \frac{\partial v}{\partial w}$$

$$= \frac{\partial E_{\text{total}}}{\partial e_j} \times \frac{\partial e_j}{\partial y_i} \times \frac{\partial y_i}{\partial v_j} \times \frac{\partial v_j}{\partial w_{ji}} \quad (8) \rightarrow \text{term by term}$$

$$\frac{\partial E_{\text{total}}}{\partial e_j} = \frac{\partial}{\partial e_j} (E_1(w) + E_2(w) + \dots + E_K(w)) \quad \text{from (5)}$$

$$= \frac{\partial}{\partial e_j} \left(\frac{1}{2} \sum_{j=1}^m e_j^2 + \frac{1}{2} \sum_{j=1}^m e_j^2 + \dots + \frac{1}{2} \sum_{j=1}^m e_j^2 + \dots + \frac{1}{2} \sum_{j=1}^m e_j^2 \right)$$

$$= \frac{1}{2} \sum_{j=1}^m \frac{\partial}{\partial e_j} (e_1^2 + e_2^2 + \dots + e_K^2)$$

$$= 1/2 \times 2 \times e_j^2$$

$$\frac{\partial E_{\text{total}}}{\partial e_j} = e_j^2$$

$$\frac{\partial e_j}{\partial y_i} = \frac{\partial}{\partial y_i} (d_j - y_i) = -1 \quad \text{ie } \frac{\partial e_j}{\partial y_i} = -1 \quad \text{From (1)}$$

$$\frac{\partial y_i}{\partial v_j} = \frac{\partial}{\partial v_j} \left(\frac{e^{v_j}}{1+e^{v_j}} \right) \quad \text{ie } \frac{\partial y_i}{\partial v_j} = y_j(1-y_j)$$

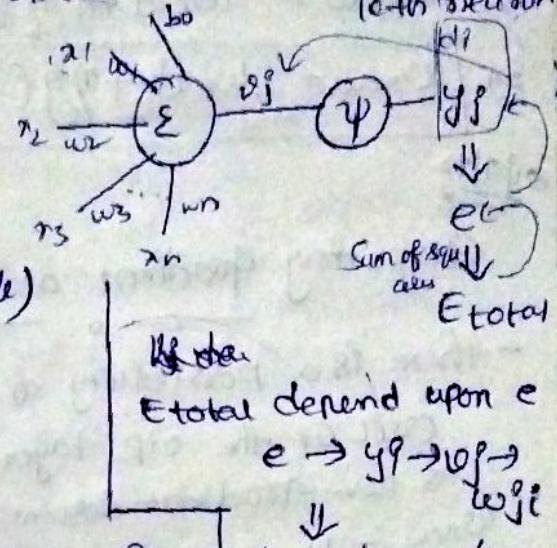
$$\begin{aligned} &= \frac{\partial}{\partial v_j} \left(\frac{e^{v_j}}{1+e^{v_j}} \right) \\ &\text{taking derivation} = \frac{e^{v_j} e^{v_j}}{(1+e^{v_j})(1+e^{v_j})} \\ &\text{or} \quad = y_j(1-y_j) \end{aligned}$$

$$y_j = \frac{1}{1+e^{-v_j}} = \frac{e^{v_j}}{1+e^{v_j}}$$

$$\text{ie } \frac{\partial v_j}{\partial w_{ji}} = \sum_{i=1}^n w_{ji} x_i$$

Consider

$$\frac{\partial v_j}{\partial w_{ji}} = \frac{\partial}{\partial w_{ji}} \left(\sum_{i=1}^n w_{ji} x_i + b_0 \right) = \sum_{i=1}^n \frac{\partial v_j}{\partial w_{ji}}$$



It is clear

E_{total} depend upon e
 $e \rightarrow y_j \rightarrow v_j \rightarrow w_{ji}$

dependency b/w
 E_{total} & weight

$$= \frac{1}{2} \sum_{j=1}^m e_j^2$$

$$g(n) = \frac{\partial E_{\text{total}}}{\partial w_{ji}}$$

$$= e_j^T \times (-1) \times y_f^T \times \Delta_i$$

$$= -y_f(1-y_f) e_j^T \Delta_i$$

$$\Rightarrow w_{\text{new}} = w_{\text{old}} + \eta g(n)$$

$$\Rightarrow \boxed{w_{\text{new}} = w_{\text{old}} + \eta y(1-y) e_j^T x_i} \quad \text{--- (6)}$$

21/9/22

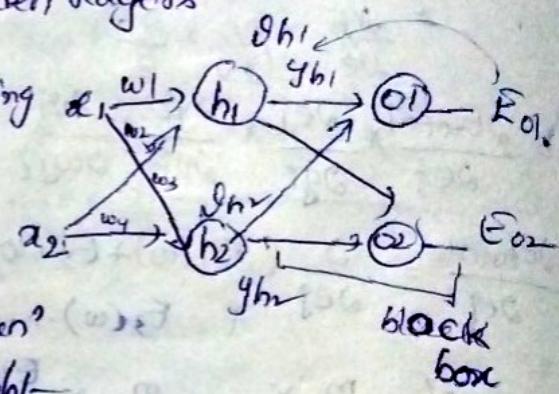
Computing gradient at hidden layers

- there is a possibility of getting error at the OLP layer.

As per gradient descent alg

$$w_{\text{new}} = w_{\text{old}} - \eta g(n)$$

$g(n)$ = Gradient of cost func?
with respect to weight



- $E_{\text{total}} = \text{error caused by hidden layer neurons at OLP layer neuron } j$ [we look back by assuming entire OLP layer is black box]

$$- g(n) = \frac{\partial E_{\text{total}}}{\partial w_{ji}} \quad \text{--- (1)} \quad \begin{matrix} \text{The OLP hidden layer is} \\ \text{black box} \end{matrix}$$

$$\frac{\partial E_{\text{total}}}{\partial w_{ji}} = \frac{\partial E_{\text{total}}}{\partial y_i} \times \frac{\partial y_i}{\partial y_j} \times \frac{\partial y_j}{\partial w_{ji}} \quad \text{--- (2)} \quad \begin{matrix} \text{Cause for the error} \\ \text{at OLP layer} \end{matrix}$$

$$\frac{\partial E_{\text{total}}}{\partial y_i} = \frac{\partial E_{O1}}{\partial y_i} + \frac{\partial E_{O2}}{\partial y_i} \quad \text{--- (3)} \quad \begin{matrix} E_{\text{total}} \rightarrow y_i \rightarrow O_1 \rightarrow O_2 \\ \downarrow \text{dependency} \end{matrix}$$

hidden layer
OLP

$$= \frac{\partial E_{O1}}{\partial y_{h1}} + \frac{\partial E_{O2}}{\partial y_{h1}} \quad \boxed{y_{h1} \in O_1 \leftarrow y_{O1} \in O_2 \in E_{O1}}$$

Consider,

$$\frac{\partial E_{O1}}{\partial y_i} = \frac{\partial E_{O1}}{\partial y_{O1}} \times \frac{\partial y_{O1}}{\partial y_{h1}} \times \frac{\partial y_{h1}}{\partial w_{ij}} \quad \begin{matrix} \text{here we need to} \\ \text{consider OLP} \\ \text{layer neurons} \\ \text{too} \end{matrix}$$

$\hookrightarrow (4)$

$$E_{01} = \frac{1}{2} \sum_{i=1}^n e_{01}^i \rightarrow \boxed{\frac{\partial E_{01}}{\partial e_{01}} = e_{01}}$$

$$\frac{\partial E_{01}}{\partial e_{01}} = e_{01}, \quad \frac{\partial E_{02}}{\partial e_{02}} = e_{02} \quad \boxed{y_{01} = \psi(v_1)}$$

$$e_{01} = d_{01} - y_{01}$$

$$\boxed{\frac{\partial e_{01}}{\partial y_{01}} = -1}$$

$$\boxed{\frac{\partial y_{01}}{\partial v_{01}} = y_{01}(1-y_{01})}$$

(for Sigmoid activation fun.)

$$\Rightarrow \boxed{\frac{\partial v_{01}}{\partial y_i} = w_{ji}}$$

where

$$v_{01} = \sum_{j=1}^m w_{ji} y_j + b_0$$

From ①

$$\begin{aligned} \frac{\partial E_{01}}{\partial y_i} &= \frac{\partial E_{01}}{\partial e_{01}} \times \frac{\partial e_{01}}{\partial y_{01}} \times \frac{\partial y_{01}}{\partial v_{01}} \times \frac{\partial v_{01}}{\partial x_i} \\ &= e_{01} \times -1 \times y_{01} (1-y_{01}) \times w_{ji} \\ &= -y_{01}(1-y_{01}) e_{01} w_{ji} \end{aligned}$$

$$\boxed{\frac{\partial y_i}{\partial x_i} = y_i(1-y_i)}$$

$$\boxed{\frac{\partial v_j}{\partial w_{ji}} = x_i}$$

From ②

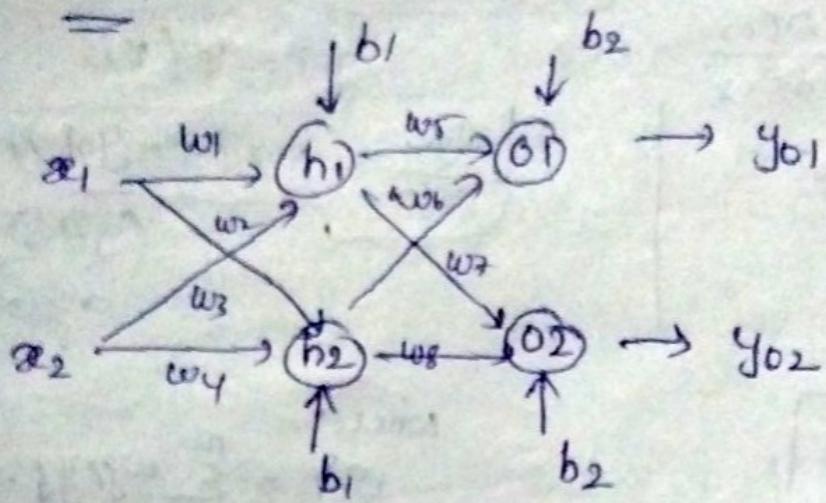
$$\begin{aligned} \frac{\partial E_{\text{total}}}{\partial w_{ji}} &= -y_{01}(1-y_{01}) e_{01} w_{ji} y_j (1-y_j) x_i \\ &= -y_{01}(1-y_{01}) y_j (1-y_j) e_0 w_{ji} x_i \end{aligned}$$

$$g(m) = -y_{01}(1-y_{01}) y_j (1-y_j) e_0 w_{ji} x_i$$

$$w_{\text{new}} = w_{\text{old}} + \eta y_{01}(1-y_{01}) y_j (1-y_j) e_0 w_{ji} x_i$$

↳ at hidden layers.

Example



Q1

$$v_{o1} = y_{h1} \times w_5 + y_{h2} \times w_6 + b_2$$

$$y_{o1} = \text{Sigmoid}(v_{o1}) = \frac{1}{1+e^{-v_{o1}}}$$

$$e_{o1} = d_{o1} - y_{o1}$$

$$E_{o1} = \gamma_2 \sum_{P=1}^n e_{o1}(P)$$

$$= \gamma_2 \sum_{P=1}^n (d_{o1}(P) - y_{o1}(P))^2$$

$$w_5(\text{new}) = w_5(\text{old}) - \eta g(n)$$

$$g(n) = -e_{o1} y_{o1} (1-y_{o1}) y_{h1}$$

$$w_5(\text{new}) = w_5(\text{old}) + \eta e_{o1} y_{o1} (1-y_{o1}) y_{h1}$$

$$= \eta w_5(1^{P-1}) \eta^P (1^{P-1}) \dots$$

$$w_6(\text{new}) = w_6(\text{old}) - \eta g(n)$$

$$g(n) = -e_{o1} y_{o1} (1-y_{o1}) y_{h2}$$

$$\underline{\underline{v}}_{02} = y_{h1} \times w_7 + y_{h2} \times w_8 + b_2$$

$$y_{02} = \text{Sigmoid}(v_{02}) = \frac{1}{1+e^{-v_{02}}}$$

$$e_{02} = d_{02} - y_{02}$$

$$E_{02} = \eta \sum_{p=1}^n e_{02}^{(p)} = \eta \sum_{p=1}^n (d_{02}^{(p)} - y_{02}^{(p)})^2$$

$$w_7(\text{new}) = w_7(\text{old}) - \eta g(n)$$

$$g(n) = -e_{02} y_{02} (1-y_{02}) y_{h1}$$

$$w_7(\text{new}) = w_7(\text{old}) + \eta e_{02} y_{02} (1-y_{02}) y_{h1}$$

$$w_8(\text{new}) = w_8(\text{old}) - \eta g(n)$$

$$g(n) = -e_{02} y_{02} (1-y_{02}) y_{h2}$$

23/2/22

Example

Iteration 1:

Forward Pass

at h_1

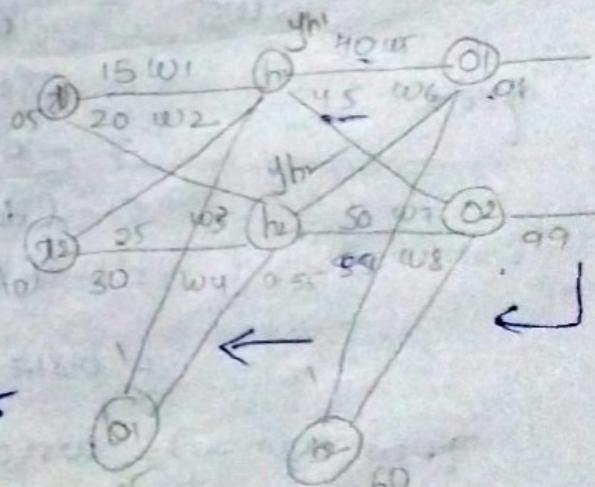
$$\begin{aligned} v_{h1} &= w_1 x_1 + w_2 x_2 + b_1 \\ &= 0.15 \times 0.05 + 0.20 \times 0.10 + 0.35 \\ &= 0.3825 \end{aligned}$$

$$y_{h1} = \psi(v_{h1}) = \frac{1}{1+e^{-v_{h1}}} = \frac{1}{1+e^{-0.3825}} = 0.5935$$

at h_2

$$\begin{aligned} v_{h2} &= w_3 x_1 + w_4 x_2 + b = 0.05 \times 0.05 + 0.30 \times 0.10 + 0.35 \\ &= 0.39 \end{aligned}$$

$$y_{h2} = \psi(v_{h2}) = \frac{1}{1+e^{-v_{h2}}} = \frac{1}{1+e^{-0.39}} = 0.5962$$



at 01

$$v_{01} = y_{h1}w_5 + y_{h2}w_7 + b_2 \\ = 0.594 \times 0.40 + 0.5962 \times 0.50 + 0.60 \\ = 1.1859$$

$$y_{01} = \psi(v_{01}) = \frac{1}{1+e^{-v_{01}}} = \frac{1}{1+e^{-1.1859}} = 0.7569$$

$$v_{02} = y_{h1}w_8 + y_{h2}w_9 + b_2 = 0.594 \times 0.45 + 0.5962 \\ = 1.1954 \quad \times 0.55 + 0.60$$

$$y_{02} = \psi(v_{02}) = \frac{1}{1+e^{-v_{02}}} = 0.767$$

	d	y
01	0.01	0.7569
02	0.99	0.767

Backward pass

$$w_5(\text{new}) = w_5(01d) - \eta g(m)$$

$$g(m) = \frac{\partial E_{\text{total}}}{\partial w_5} = \frac{\partial E_{\text{total}}}{\partial v_{01}} \times \frac{\partial v_{01}}{\partial y_{01}} \times \frac{\partial y_{01}}{\partial w_5} \times \frac{\partial w_5}{\partial w_3} \\ = -e_0 \times -1 \times y_{01}(1-y_{01}) \times y_{h1} \\ = -e_{01} \times y_{01} \times (1-y_{01}) \times y_{h1} \\ = -(d_{01} - y_{01})(y_{01}(1-y_{01})) \times y_{h1} \\ = -0.8(0.01 - 0.7569)(1 - 0.7569) \\ = 0.817 \times 0.5945$$

$$w_5(\text{new}) = w_5(01d) - \eta g(m)$$

$$= 0.40 - 0.5 \times 0.817 \\ = 0.359$$

$$w_7(\text{new}) = w_7(01d) - \eta g(m)$$

$$g(m) = \frac{\partial E_{\text{total}}}{\partial w_7} = \frac{\partial E_{\text{total}}}{\partial v_{01}} \times \frac{\partial v_{01}}{\partial y_{01}} \times \frac{\partial y_{01}}{\partial w_1} \times \frac{\partial w_1}{\partial w_7} \\ = e_{01} \times -1 \times y_{01} \times (1-y_{01}) \times y_{h2}$$

$$= -(d_{01} - y_{01}) y_{01} (1 - y_{01}) y_{h2}$$

$$= -(0.01 - 0.7569) 0.7569 (1 - 0.7569) 0.5962$$

$$= 0.0819$$

$$w_7(\text{new}) = w_7(\text{old}) - \eta g(n)$$

$$= 0.50 - (0.5)(0.0819) = 0.459$$

$$v_{02} = y_{h1} w_6 x y_{h2}^{w_8 + b_2}$$

$$\overline{w_6(\text{new})} = w_6(\text{old}) - \eta g(n)$$

$$g(n) = \frac{\partial E_{\text{total}}}{\partial w_6} = \frac{\partial E_{\text{total}}}{\partial e_{02}} \times \frac{\partial e_{02}}{\partial y_{02}} \times \frac{\partial y_{02}}{\partial w_6} \times \frac{\partial v_{02}}{\partial w_6}$$

$$= e_{02} \times -1 \times y_{02} (1 - y_{02}) \times y_{h1}$$

$$= -e_{02} (y_{02}) (1 - y_{02}) y_{h1}$$

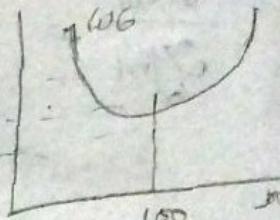
$$= -(d_{02} - y_{02}) y_{02} (1 - y_{02}) y_{h1}$$

$$= -(0.99 - 0.7569) (0.7569) (1 - 0.7569) 0.5962$$

$$= -0.023$$

$$w_6(\text{new}) = w_6(\text{old}) - \eta g(n)$$

$$= 0.45 - (0.5)(-0.023) = 0.46175$$



$$\overline{w_8(\text{new})} = w_8(\text{old}) - \eta g(n)$$

~~$$g(n) = \frac{\partial E_{\text{total}}}{\partial w_8} \times \frac{\partial v_{02}}{\partial e_{02}} \times \frac{\partial e_{02}}{\partial y_{02}} \times \frac{\partial y_{02}}{\partial w_8} \times \frac{\partial v_{02}}{\partial w_8}$$~~

$$= e_{02} \times -1 \times y_{02} (1 - y_{02}) \times y_{h1} = -(d_{02} - y_{02}) y_{02} (1 - y_{02}) y_{h1}$$

$$= -(0.99 - 0.7569) (0.7569) (1 - 0.7569) 0.5962$$

$$= -0.0234$$

$$w_8(\text{new}) = w_8(\text{old}) - \eta g(n) = 0.55 - (0.5)(-0.0234)$$

$$= 0.5617$$

at hidden layers

at h_1

update of w_1

$$w_{1,(\text{new})} = w_{1,(\text{old})} - \eta g(n)$$

$$g(n) = \frac{\partial E_{\text{total}}}{\partial w_1}$$

$$\frac{\partial E_{\text{total}}}{\partial w_1} = \frac{\partial E_{\text{total}}}{\partial h_1} \times \frac{\partial y_{h1}}{\partial v_{h1}} \times \frac{\partial v_{h1}}{\partial w_1}$$

$$\frac{\partial E_{\text{total}}}{\partial h_1} = \frac{\partial E_{01}}{\partial y_{h1}} + \frac{\partial E_{02}}{\partial y_{h1}} \quad \text{--- (2)}$$

$$\frac{\partial E_{01}}{\partial h_1} = \frac{\partial E_{01}}{\partial e_{01}} \times \frac{\partial e_{01}}{\partial y_{01}} \times \frac{\partial y_{01}}{\partial v_{01}} \times \frac{\partial v_{01}}{\partial y_{h1}}$$

$$\frac{\partial E_{02}}{\partial h_1} = \frac{\partial E_{02}}{\partial e_{02}} \times \frac{\partial e_{02}}{\partial y_{02}} \times \frac{\partial y_{02}}{\partial v_{02}} \times \frac{\partial v_{02}}{\partial y_{h1}}$$

$$\Rightarrow \frac{\partial E_{01}}{\partial y_{h1}} = \frac{\partial E_{01}}{\partial e_{01}} \times \frac{\partial e_{01}}{\partial y_{01}} \times \frac{\partial y_{01}}{\partial v_{01}} \times \frac{\partial v_{01}}{\partial y_{h1}}$$

$$= e_{01} \times -1 \times y_{01}(1-y_{01}) \times w_5$$

$$y_{01} = \frac{1}{1+e^{-v_{01}}}$$

$$= -e_{01} y_{01}(1-y_{01}) w_5$$

$$\begin{aligned} E_{01} &= \frac{1}{2} \sum_{p=1}^2 e_{01}^p \\ \frac{\partial E_{01}}{\partial e_{01}} &= e_{01} \end{aligned}$$

$$\frac{\partial y_{01}}{\partial v_{01}} = y_{01}(1-y_{01})$$

$$v_{01} = y_{h1} \times w_5 + y_{h2} \times w_7 + b_2$$

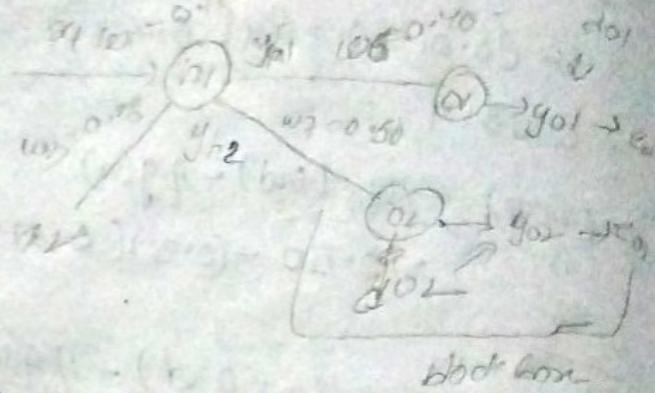
$$e_{01} = d_{01} - y_{01}$$

$$\frac{\partial v_{01}}{\partial y_{01}} = -1$$

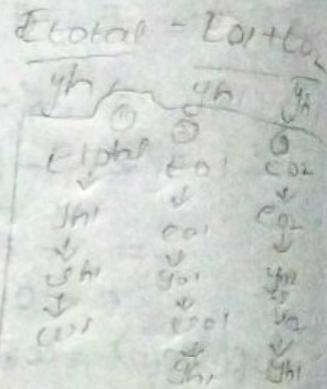
$$\frac{\partial v_{01}}{\partial y_{h1}} = w_5$$

$$\begin{aligned} \text{Similarly, } \frac{\partial E_{02}}{\partial y_{h1}} &= \frac{\partial E_{02}}{\partial e_{02}} \times \frac{\partial e_{02}}{\partial y_{02}} \times \frac{\partial y_{02}}{\partial v_{02}} \times \frac{\partial v_{02}}{\partial y_{h1}} \\ &= e_{02} \times -1 \times y_{02}(1-y_{02}) \times w_6 \end{aligned}$$

$$\begin{aligned} \frac{\partial v_{02}}{\partial y_{h1}} &= -e_{02} y_{02}(1-y_{02}) w_6 \end{aligned}$$



$$v_{01} = \sum_{j=1}^2 v_{hj} w_{0j} + b_0$$



$$\begin{aligned}\frac{\partial E_02}{\partial y_{h1}} &= -e_{01} \times y_{01} \times (1-y_{01}) \times w_5 \\ &= -(d_{01}-y_{01}) \times y_{01} \times (1-y_{01}) \times w_5 \\ &= -(0.01-0.7569) \times 0.4569 (1-0.4569) \times 0.40 \\ &= 0.05497\end{aligned}$$

$$\begin{aligned}\frac{\partial E_02}{\partial y_{h1}} &= -e_{02} \times y_{02} \times (1-y_{02}) \times w_6 \\ &= -(d_{02}-y_{02}) \times y_{02} \times (1-y_{02}) \times w_6 \\ &= -(0.99-0.467) \times 0.467 (1-0.467) \times 0.45 \\ &= -0.0878 - 0.01493\end{aligned}$$

$$\frac{\partial E_{\text{total}}}{\partial h_1} = \frac{\partial E_{01}}{\partial y_1} + \frac{\partial E_{02}}{\partial y_{h1}} = 0.05497 - 0.01493 = 0.034$$

$$w_1(\text{new}) = w_1(\text{old}) - \eta g(n)$$

$$D_h = w_1 z_1 + w_3 z_2 + b$$

$$\begin{aligned}\frac{\partial E_{\text{total}}}{\partial w_1} &= \frac{\partial E_{\text{total}}}{\partial y_{h1}} \times \frac{\partial y_{h1}}{\partial w_1} \times \frac{\partial w_1}{\partial w_1} \\ &= 0.034 \times y_{h1} (1-y_{h1}) \times 2, \\ &= 0.034 \times (0.594) (1-0.594) (0.05) \\ &= 4.465 \times 10^{-4} = 0.0004465\end{aligned}$$

$$\begin{aligned}w_1(\text{new}) &= w_1(\text{old}) - \eta g(n) \quad \eta = 0.5 \\ &= 0.15 - (0.5)(0.0004465) \\ &= 0.1497\end{aligned}$$

• Updation of w_3

$$w_3(\text{new}) = w_3(\text{old}) - \eta g(n)$$

$$g(n) = \frac{\partial E_{\text{total}}}{\partial w_3} = \frac{\partial E_{\text{total}}}{\partial h_1} \times \frac{\partial y_{h1}}{\partial w_3} \times \frac{\partial w_3}{\partial w_3}$$

$$\frac{\partial E_{\text{total}}}{\partial h_1} = \frac{\partial E_{01}}{\partial y_{h1}} + \frac{\partial E_{02}}{\partial h_2} \quad \text{--- (2)}$$

$$\begin{aligned}\frac{\partial E_{01}}{\partial h_1} &= \frac{\partial E_{01}}{\partial e_{01}} \times \frac{\partial e_{01}}{\partial y_{01}} \times \frac{\partial y_{01}}{\partial w_0} \times \frac{\partial w_0}{\partial y_{h1}} = e_{01} \times -1 \times y_{01} (1-y_{01}) \times w_5 \\ \frac{\partial E}{\partial h_1} &= -(0.01-0.7569) \times 0.4569 (1-0.4569) \times 0.40 \\ &= 0.05497\end{aligned}$$

$$\begin{aligned}\frac{\partial E_{O_2}}{\partial y_{h1}} &= \frac{\partial E_{O_2}}{\partial e_{O_2}} \times \frac{\partial e_{O_2}}{\partial y_{O_2}} \times \frac{\partial y_{O_2}}{\partial w_{O_1}} \times \frac{\partial w_{O_1}}{\partial y_{h1}} \\ &= e_{O_2} \times -1 \times y_{O_2} (1-y_{O_2}) \times w_7 \\ &= -0.01793\end{aligned}$$

$$\frac{\partial E_{\text{total}}}{\partial y_{h1}} = \frac{\partial E_{O_1}}{\partial y_{h1}} + \frac{\partial E_{O_2}}{\partial y_{h1}} = 0.05497 - 0.01793 = 0.034$$

$$\frac{\partial E_{\text{total}}}{\partial w_3} = \frac{\partial E_{\text{total}}}{\partial y_{h1}} \times y_{h1} \times (1-y_{h1}) \times y_{h2} = 0.034 \times 0.594 \\ (1-0.594) \times 0.10 = 0.00089$$

$$\begin{aligned}\rightarrow w_3(\text{new}) &= w_3(\text{old}) - 2g(n) \\ &= 0.25 - 0.00089 = 0.249\end{aligned}$$

At h₂

updation of w₂

$$e_{O_2}(\text{new}) = w_2(\text{old}) - 2g(n)$$

$$g(n) = \frac{\partial E_{\text{total}}}{\partial w_2}$$

$$= \frac{\partial E_{\text{total}}}{\partial y_{h2}} \times \frac{\partial y_{h2}}{\partial w_2} \times \frac{\partial w_2}{\partial e_{O_2}} \quad \text{--- (1)}$$

$$\frac{\partial E_{\text{total}}}{\partial y_{h2}} = \frac{\partial E_{O_1}}{\partial y_{h2}} * \frac{\partial E_{O_2}}{\partial y_{h2}} \quad \text{--- (2)}$$

$$\frac{\partial E_{O_1}}{\partial y_{h2}} = \frac{\partial E_{O_1}}{\partial e_{O_1}} \times \frac{\partial e_{O_1}}{\partial y_{O_1}} \times \frac{\partial y_{O_1}}{\partial w_{O_1}} \times \frac{\partial w_{O_1}}{\partial y_{h2}}$$

$$= e_{O_1} \times -1 \times y_{O_1} (1-y_{O_1}) \times w_8$$

$$= -(e_{O_1} - y_{O_1}) \times y_{O_1} (1-y_{O_1}) \times w_8$$

$$= -(0.01 - 0.7569) \times 0.7569 (1-0.7569) \times 0.50$$

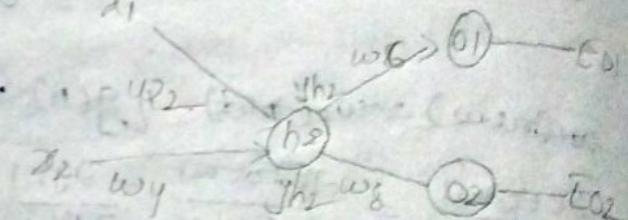
$$= 0.068$$

$$\frac{\partial E_{O_2}}{\partial y_{h2}} = e_{O_2} \times -1 \times y_{O_2} (1-y_{O_2}) \times w_8 \quad \text{--- (3)}$$

$$= -(0.99 - 0.467) (1-0.467) (0.467) \times 0.50 \quad \text{--- (4)}$$

$$= -0.02147$$

$$\begin{aligned}\frac{\partial E_{\text{total}}}{\partial y_{h2}} &= 0.068 - 0.02147 \\ &= 0.047\end{aligned}$$



$$\begin{aligned}w_{O_1} &= y_{h1} \times w_5 \\ &\quad + y_{h1} \times w_8 \\ &\quad + b_2 \\ &\quad + b_3\end{aligned}$$

$$\begin{aligned}e_{O_2} &= y_{h2} \times w_6 \\ &\quad + y_{h2} \times w_8 \\ &\quad + b_2 \\ &\quad + b_3\end{aligned}$$

$$\frac{\partial E_{\text{total}}}{\partial w_2} = 0.047 \times y_{h_2} \times (1 - y_{h_2}) \times x_2$$

$$= 0.047 \times 0.596 \times (1 - 0.596) \times 0.05 = 0.005$$

$$w_2(\text{new}) = w_2(\text{old}) - \eta g(n)$$

$$= 0.20 - 0.5 \times 0.005 = 0.195$$

$$= 0.199$$

$$w_4(\text{new}) = w_4(\text{old}) - \eta g(n)$$

$$= 0.30 - 0.5 \times 0.00011$$

$$= 0.299$$

~~$\frac{\partial E_{\text{total}}}{\partial w_2} / \frac{\partial E_{\text{total}}}{\partial w_4}$~~

$$g(n) = \frac{\partial E_{\text{total}}}{\partial w_4}$$

$$= \frac{\partial E_{\text{total}}}{\partial y_{h_2}} \times \frac{\partial y_{h_2}}{\partial w_4}$$

$$= 0.047 \times y_{h_2} \times (1 - y_{h_2}) \times x_L$$

$$= 0.047 \times 0.596 \times (1 - 0.596) \times 0.10$$

$$= 0.00011$$

After Iteration 1

$$\begin{aligned} w_1 &= 0.149 & w_5 &= 0.359 \\ w_2 &= 0.199 & w_6 &= 0.5614 \\ w_3 &= 0.249 & w_7 &= 0.459 \\ w_4 &= 0.299 & w_8 &= 0.5614 \end{aligned}$$

$$\begin{cases} x_1 = 0.05 \\ x_2 = 0.10 \end{cases}$$

Iteration 2
Forward pass is done

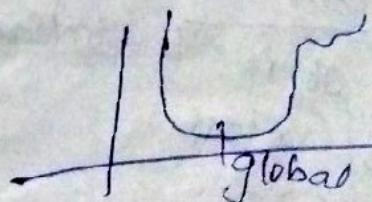
Back propagation algorithm.

① Rate of Learning
algorithm specifies a step size taken by learning
takes a value in between 0-1 $0 \leq \eta \leq 1$

2. If η is very small:

$$w(\text{new}) = w(\text{old}) - \eta g(n)$$

\downarrow
Learning rate



→ slowly converges to the global minimum
 - if increase no. of iterations required for convergence
 - algorithm can trapped into the local minimum point.

η very large

→ algorithm can overshoot global minimum point.

→ can oscillate

→ avoids getting trap into local minima point.

So,

Momentum :-

$$-\Delta w_{ij}(n) = \alpha \Delta w_{ij}(n-1) + \Delta w_{ij}(n)$$

- provides solⁿ to problems associated with learning rate.

- adds small amount of update from previous iteration controlled by momentum term (α)

→ Sequential or Batch mode :-

epoch

↓
presentation of complete set of examples or

examples of samples

- also called online mode

- update for weight vector computed sequentially

by considering only one sample at a time.

It requires less memory space as it slowly

converges towards global minima point.

drawback

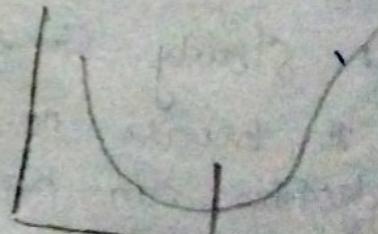
it is theoretically difficult to determine stopping

criteria.

* Design Issues in Backpropagation Algorithm

$$w(\text{new}) = w(\text{old})$$

$$- \eta g(n)$$



- ① momentum
 ② stopping criteria

→ Backpropagation Alg
 ↳ Learning Alg

- (1) when model produces error value that is acceptable. (e.g. of accuracy)
 (2) when norm of the difference b/w y & d reaches Euclidean to minimum value

- ③ New sample

- not there in the training set.

↳ Generalization

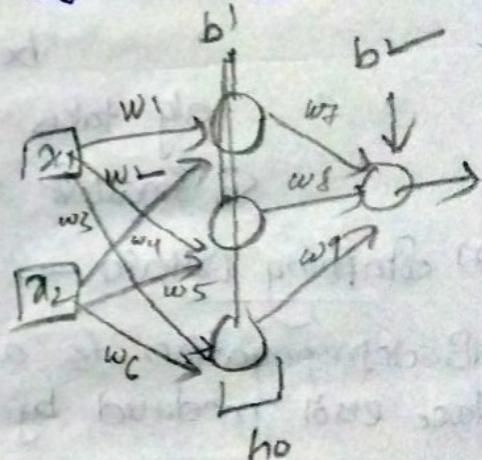
Creates
Power

* Implementation of Backpropagation Algorithm

- 1) Initialize free parameters

$$hw = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \\ w_6 \end{bmatrix} \quad C \times 1$$

$$h_0 = \begin{bmatrix} w_7 \\ w_8 \\ w_9 \end{bmatrix} \quad 3 \times 1$$



$$v_{h_1} = x_1 w_1 + x_2 w_4$$

$$v_{h_2} = x_1 w_2 + x_2 w_5$$

$$v_{h_3} = x_1 w_3 + x_2 w_6$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad 2 \times 1$$

$$hw = \begin{bmatrix} w_1 & w_4 \\ w_2 & w_5 \\ w_3 & w_6 \end{bmatrix}$$

random packag

↳ random

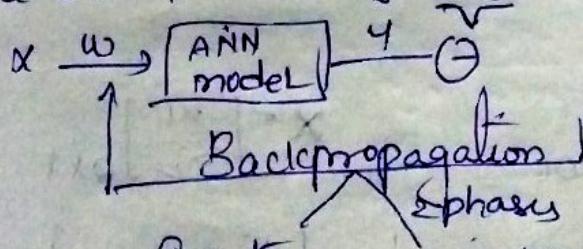
↳ to fixate values
for w, b

Design Issues in Back Propagation Algorithm

- 1) Learning rate (η) - due to this algorithm converges slowly or at fast rate. $w(\text{new}) = w(\text{old}) - \eta g(n)$ (i.e., update value of previous iteration is added to update formula)
- alg stuck in local min point
Momentum is applied (it accelerates alg and alg move forward from local min point and to reach global min point) then update formula changes
- $$\Delta w = \Delta w(n) + \alpha \Delta w(n-1) - \eta g(n)$$
- α -momentum term
 - lies b/w 0 to 1 mostly
 - $\alpha = 0.001$ or 0.01
 - α is multiplied with $\Delta w(n-1)$ because if $\Delta w(n-1)$ is too high then alg takes few & faster steps to reach GMP.
 - So, smallest value α is multiplied with $\Delta w(n-1)$
-

(2) Stopping Criteria —

Backpropagation is a learning / optimization alg used to reduce error produced by AN



$d-y = e \rightarrow$ Simple Gradient descent Alg

Goes to goes to feed forward phase
iteration process (read s/p produces O/P & identifies error ~~rate~~)

This process occurs iteratively.
If no stopping criteria exists, then Back propagation iterates infinitely.

when $y \neq d$ - but leads to infinite loop
 General stopping criteria used is to be BP accepted
 AIG is stopped when error values reaches a part threshold
 value (98%).
 correct values should predicted
 Stopping criteria = 5%
 95% accurate

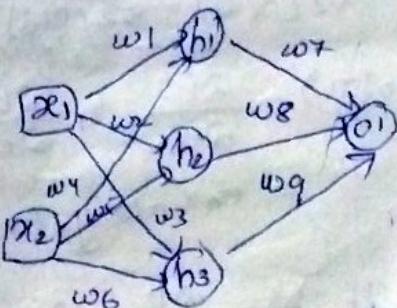
- when Euclidean norm of d & y reaches to a min. value
- new sample not there in training set
 - Backpropagation goes through generalization.
- process (it tries to compare data with already trained! (given sample) previous samples)
 data & make conclusions.

* poweret is created for the training data using which model is trained. poweret contains all features of dataset so that based on it model can make predictions for new data to predict its next for generalization.

12/10/22

Implementing Backpropagation

- Initialize parameters
 - weights
 - bias
 - Learning parameter
- Construct neural network
 - no. of IIPs
 - no. of hidden layers
 - no. of neurons
 - no. of neurons in OLP layer



$$\begin{aligned}
 v_{h1} &= w_1 x_1 + w_4 x_2 + b \\
 v_{h2} &= w_2 x_1 + w_5 x_2 + b \\
 v_{h3} &= w_3 x_1 + w_6 x_2 + b
 \end{aligned}$$

$$\begin{aligned}
 u_q &= \begin{bmatrix} w_1 & w_4 \\ w_2 & w_5 \\ w_3 & w_6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} w_1 x_1 + w_4 x_2 \\ w_2 x_1 + w_5 x_2 \\ w_3 x_1 + w_6 x_2 \end{bmatrix} \\
 &\quad \downarrow \qquad \qquad \qquad n \times 1 \\
 &\text{Size of IIP} \times \text{no. of neurons in h1}
 \end{aligned}$$

$$w_2 = \begin{bmatrix} w_7 \\ w_8 \\ w_9 \end{bmatrix}_{3 \times 1}^T \Rightarrow \begin{bmatrix} w_7 & w_8 & w_9 \end{bmatrix}_{3 \times 1} \Rightarrow \text{for OLP layer} \quad \begin{bmatrix} y_{h1} \\ y_{h2} \\ y_{h3} \end{bmatrix}$$

Size of hidden layer \times Size of OLP layer

$$\Rightarrow w_7 y_{h1} + w_8 y_{h2} + w_9 y_{h3}$$

We use random module for creating w, b values.

If we want the more bias values

$$\text{bias} \cdot b_n = \begin{bmatrix}] \end{bmatrix}_{1 \times 1}, \quad (n = \text{no of hidden layer neurons}),$$

$$b_0 = \begin{bmatrix}] \end{bmatrix}_{1 \times 1}$$

Step 3: Compute the error value

$$e = d - y$$

Since error occurs back propagation updates the weights

$$w_1(\text{new}) = w_1(\text{old}) - \eta g_h(n)$$

$$w_2(\text{new}) = w_2(\text{old}) - \eta g(n)$$

$$g_0(n) = \frac{\partial E}{\partial w_2} = \frac{\partial E}{\partial e_0} \frac{\partial e_0}{\partial y_0} \frac{\partial y_0}{\partial v_0} \frac{\partial v_0}{\partial w_2}$$

$$= e_0 x - 1 \times y_0(1-y_0) y_h$$

First Compute δ value

$$g_h(n) = \frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial y_h} \times \frac{\partial y_h}{\partial v_h} \times \frac{\partial v_h}{\partial w_1}, \quad \text{propagation of error by hidden layer}$$

$$= \frac{\partial E}{\partial e_0} \frac{\partial e_0}{\partial y_0} \frac{\partial y_0}{\partial v_0} \frac{\partial v_0}{\partial y_h} \times \frac{\partial y_h}{\partial v_h} \times \frac{\partial v_h}{\partial w_1}$$

$$= \frac{\delta \times w_2 \times y_h(1-y_h) \times x}{w_2}$$

$$D_{\text{dataset}} = \delta \times w_2 \times y_h(1-y_h) \times$$

	hours studied	have enjoyed	result
1	1	2	80
2	2	6	50
3	9	1	95

Initializing Parameters. import numpy as np

inputSize = 2

hiddenSize = 3

outputSize = 1

lr = 0.1
w1 = np.random.randn(inputSize, hiddenSize) * 0.01

w2 = np.random.randn(hiddenSize, outputSize) * 0.01

b1 = np.random.randn(hiddenSize, 1) * 0.01

b2 = np.random.randn(outputSize, 1) * 0.01

define activation fun

def Sigmoid(v):

$$y = 1 / (1 + np.exp(-v))$$

return y

def derivative(y):

$$\text{return } y * \text{exp}(1-y)$$

y_p = predicted dp

def mseError(yP, y)

$$E = ((yP - y) ** 2).sum() / 2$$

return E

define a dataset

x = np.array([[7, 2], [2, 6], [9, 17]])

y = np.array([80, 50, 95])

Normalize

x = x / np.amax(x, axis=0)

y = y / 100

Create a neural network and train

def train(x, y):

global w1, w2, b1, b2, ls

forward phase.

$$v_h = \text{np.dot}(x, w_1) + b_1 \quad \text{(forward pass)} \quad v_h(w)$$

$$y_h = \text{Sigmoid}(v_h)$$

$$v_o = \text{np.dot}(y_h, w_2) + b_2 \quad (2) \quad v_o @ w_2$$

$$y_o = \text{Sigmoid}(v_o)$$

return y_o

train(x, y)

Back propagation phase

$$e_1 = y_o - y$$

$$\delta_{\text{out}} = e_1 * y_o * (1 - y_o)$$

$$e_2 = \delta_{\text{out}} @ w_2 \text{ T}$$

$$\delta_{\text{hidden}} = e_2 * \text{Sigmoid}(y_h(1 - y_h))$$

$$\delta_o = y_h @ \delta_{\text{hidden}}$$

$$\delta_h = x.T @ \delta_{\text{hidden}}$$

$$w_1 = w_1 - \eta_r * \delta_h$$

$$w_2 = w_2 - \eta_r * \delta_o$$

$$b_1 = b_1 - \eta_r * \delta_{\text{hidden}}$$

$$b_2 = b_2 - \eta_r * \delta_o$$

, # Back propagation is iterative algorithm

-for i in range(2000):

train(x, y)

def forwardtest(x, y):

global w1, w2, b1, b2

forward phase

$$v_h = \text{np.dot}(x, w_1) + b_1 \quad (3) \quad (x @ w_1) + b_1$$

$$y_h = \text{Sigmoid}(v_h)$$

$$v_o = \text{np.dot}(y_h, w_2) + b_2 \quad (4) \quad (y_h @ w_2) + b_2$$

$$y_o = \text{Sigmoid}(v_o)$$

return y_o

$$\begin{aligned} \text{forward} &= \text{np.dot}(x, w_1) + b_1 \\ \text{backward} &= \text{np.dot}(y_h, w_2) + b_2 \end{aligned}$$

forward test (x, y)

$L_{dim} = [w_{1, dim}]$
for $i \in \text{range}(dim)$
 $w[i] \leftarrow \text{rand}(1)$ = np.
random randn(dim[1])