

26/8/20

Stochastic gradient descent:-

-drawbacks of gradient descent:

*trap into the local minima point.

*convergence takes a long time when data is large.

→ uses only one data sample to compute update.

update rule

$$w(n+1) = w(n) - \eta e(n) X$$

→ slowly converges to the global minimum point.

Newton's Method

- Approximation of quadratic cost function, then computes gradient of the approximated cost function.

- uses second order Taylor series expansion to approximate the cost function.

$$\Delta w(n) = w(n+1) - w(n)$$

$$\therefore f(x) = f(a) + f'(a)x + \frac{1}{2} f''(x) H(x) f(x) + \dots$$

$$\Delta w(n) = w(n+1) - w(n)$$

$$\approx g(n) \Delta w(n) + \frac{1}{2} \Delta w(n) H(n) \Delta w(n)$$

$H(n) \rightarrow$ Hessian matrix

$$H = \begin{bmatrix} \frac{\partial^2 C}{\partial w_1^2} & \frac{\partial^2 C}{\partial w_1 \partial w_2} & \frac{\partial^2 C}{\partial w_1 \partial w_3} & \cdots & \frac{\partial^2 C}{\partial w_1 \partial w_n} \\ \frac{\partial^2 C}{\partial w_2 \partial w_1} & \frac{\partial^2 C}{\partial w_2^2} & \frac{\partial^2 C}{\partial w_2 \partial w_3} & \cdots & \frac{\partial^2 C}{\partial w_2 \partial w_n} \\ \vdots & & & & \\ \frac{\partial^2 C}{\partial w_n \partial w_1} & \frac{\partial^2 C}{\partial w_n \partial w_2} & \cdots & \cdots & \frac{\partial^2 C}{\partial w_n^2} \end{bmatrix}$$

$$\frac{\partial f}{\partial \Delta w(n)} g(n) + H(n) \Delta w(n) = 0$$

$$H(n) = -g(n) \Rightarrow \Delta w(n) = \frac{-g(n)}{H(n)}$$

$$\Delta w(n) = -g(n) H^{-1}(n)$$

$$w(n+1) = w(n) + \Delta w(n) \quad \text{update amount}$$

$$w(n+1) = w(n) - \eta g(n) H^{-1}(n)$$

6/9/22

* Gauss-Newton Method:

$$- c(n) = \frac{1}{2} \sum_{i=1}^n e^2(i)$$

$$- \text{Linearize cost fn } (e) \rightarrow$$

$$- e(i; w) = e(i) + \left(\frac{\partial e(i)}{\partial w} \right)_{w=w(n)} (w - w(n))$$

Linearization

$$L(x) = f(a) + f'(a)(x-a)$$

$\downarrow w(n)$

Matrix notation \rightarrow ① where $i = 1, 2, 3, \dots, n$

$$- e(n; w) = e(n) + J(n) (w - w(n)) \rightarrow ②$$

$e(n; w) = e(n) + J(n) \text{ vector } x$

$$- e(n) \rightarrow \begin{bmatrix} e(1), e(2), \dots, e(n) \end{bmatrix}^T$$

$$- e(n) = [e(1), e(2), \dots, e(n)]^T$$

$J(n)$ is a jacobian matrix

$$\frac{\partial e(i)}{\partial w} = \left[\frac{\partial e(i)}{\partial w_1}, \frac{\partial e(i)}{\partial w_2}, \dots, \frac{\partial e(i)}{\partial w_m} \right]^T$$

$$J(n) = \begin{bmatrix} \frac{\partial e(1)}{\partial w_1} & \frac{\partial e(1)}{\partial w_2} & \dots & \frac{\partial e(1)}{\partial w_m} \\ \frac{\partial e(2)}{\partial w_1} & \frac{\partial e(2)}{\partial w_2} & \dots & \frac{\partial e(2)}{\partial w_m} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial e(n)}{\partial w_1} & \frac{\partial e(n)}{\partial w_2} & \dots & \frac{\partial e(n)}{\partial w_m} \end{bmatrix}$$

$$w(n+1) = \arg \min_w \frac{1}{2} \|e(n; w)\|^2$$

Substitute in ②

$$\frac{1}{2} \|e(n; w)\|^2 = \frac{1}{2} (e^2(n) + e(n) J(n) (w - w(n)))^T J^T(n) J(n) + \frac{1}{2} (w - w(n))^T (w - w(n)) \rightarrow ③$$

$$\begin{aligned}
 & \text{Take derivative of (1) w.r.t } w \\
 & e(n) J(n) + J^T(n) J(n) (w - w(n)) = 0 \\
 & e(n) J(n) + J^T(n) J(n) (w - w(n)) = -e(n) J(n) \\
 & w - w(n) = -\frac{e(n) J(n)}{J^T(n) J(n)} \\
 & (2) \\
 & w - w(n) = -e(n) J(n) (J^T(n) J(n))^{-1} \\
 & \boxed{w = w(n) - e(n) J(n) (J^T(n) J(n))^{-1}} \rightarrow (3)
 \end{aligned}$$

$$\underline{\underline{w(n+1)}}$$

* Perception convergence:-

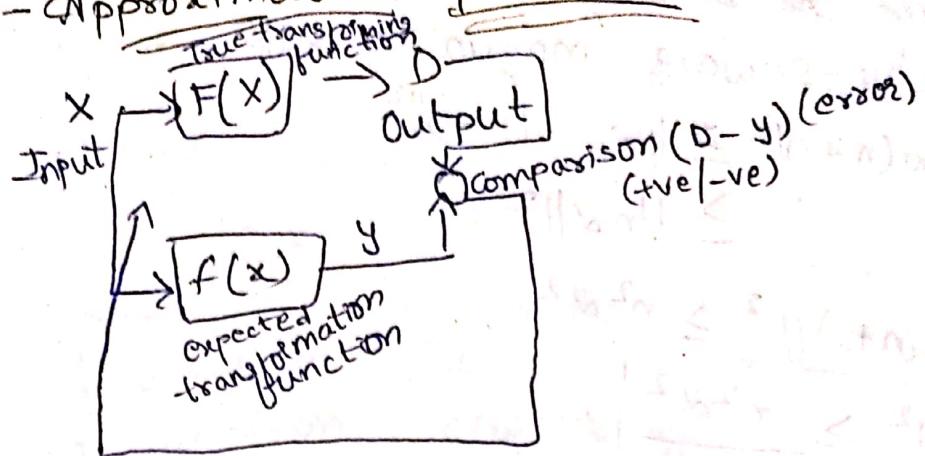
- perception is used to perform classification task.
- simple perception can perform binary classification only
- It can classify data if there is linearly separable boundary.
- Error $e_i = \hat{y}_i - y_i$
target output
actual output
- $y = w^T x + b \Rightarrow y > 0 \quad C_1 \text{ (Assumption)}$
 $y \leq 0 \quad C_2$
- when $\langle x, c_1 \rangle$ is data we have, if perception output for input x as C_1 then we don't change weight.
 $w(n+1) = w(n)$
- If there is a misclassification, then we need to update free parameters
 $w(n+1) = w(n) - \eta(n) \cdot x(n)$

sample $\langle x, c_1 \rangle$
 Perception $\langle x, c_2 \rangle$

- Sample $\langle x, c_2 \rangle$ perception op $\langle x, c_1 \rangle$

$$w(n+1) = w(n) + \eta(n) \cdot x(n)$$

- Approximation procedure:-



q9/22

$$\text{Let } \eta = 1$$

$$w(n+1) = w(n) + x(n)$$

$$w(n+1) = w(n) + x(n)$$

Initially we assume that

$n=1$, $w(0) = 0$ where perception gives correct output.

$$w(n+1) = w(n) + x(n) \quad \text{for } n=1, 2, \dots, m$$

$$w(1) = w(0) + x(0)$$

$$w(2) = x(0) + x(1)$$

$$w(n+1) = x(0) + x(1) + x(2) + \dots + x(n) \rightarrow ③$$

$$\text{Let } x(0) = 0$$

$$w(n+1) = x(1) + x(2) + \dots + x(n) \rightarrow ④$$

Multiply ④ with w_0

$$w_0 w(n+1) = w_0 x(1) + w_0 x(2) + w_0 x(3) + \dots + w_0 x(n)$$

Let α be positive constant.

$$\alpha = \min w_0 x(n) \text{ for } n \in \{1, 2, \dots, m\}$$

$$w_0 w(n+1) \geq n\alpha$$

Now Compute the boundary for data items.

Apply cauchy-schwarz inequality rule

$$\begin{aligned} \|w_0\|^2 \|w(n+1)\|^2 &\geq \|w_0 w(n+1)\|^2 \\ &\geq \|n\alpha\|^2 \end{aligned}$$

$$\|w_0\|^2 \|w(n+1)\|^2 \geq n^2 \alpha^2$$

$$\boxed{\|w(n+1)\|^2 \geq \frac{n^2 \alpha^2}{\|w_0\|^2}}$$

- Alternative method:-

$$w(n+1) = w(n) + x(n) \rightarrow \textcircled{5}$$

\Rightarrow Apply euclidean norm to $\textcircled{5}$

$$\begin{aligned} \|w(n+1)\|^2 &= \|w(n) + x(n)\|^2 \\ &= \|w(n)\|^2 + \|x(n)\|^2 + 2(w(n) \cdot x(n)) \\ &= \|w(n)\|^2 + \|x(n)\|^2 \end{aligned}$$

$$\|w(n+1)\|^2 = \|x(n)\|^2 \text{ for } n = 1, 2, \dots, m$$

$\|w(n+1)\|^2 \leq \|x(n)\|^2$ which is $\max x(n)$

Let β be constant

$$\boxed{\|w(n+1)\|^2 \leq n\beta}$$

- perceptron produces correct output after n_{\max} iterations.

At $n_{\max} + 1$ iteration

$$\boxed{\frac{n_{\max} \alpha^2}{\|w_0\|^2} = n_{\max} \beta}$$

13/9/22

*Implementing gradient descent algorithm.

Residual error

$$e = d - y$$

$$\Rightarrow y = \varphi \left(\sum_{i=1}^n x_i w_i + b \right)$$

$$y = \sum_{i=1}^n x_i w_i + b$$

$$\boxed{y = w \cdot x + b} \rightarrow 0$$

$$\text{cost function } C(w) = \frac{1}{2} \sum_{i=1}^n e_i^2$$

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

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

$$\rightarrow g(n) = - \sum_{i=1}^n e_i x_i$$

$$= - e_i x_i$$

$$w_{\text{new}} = w_{\text{old}} + \eta e_i x_i$$

$$\boxed{w(n+1) = w(n) + \eta e_i x_i}$$

$$x = df[[c_1, c_2, \dots, c_n]]$$

$$y = df[c_j]$$

code def predict(x, w, b):

$$y = np.dot(x, w) + b$$

return y

def update(x, w, b, y_pred, y, lr_rate):

$$gw = np.dot((y - y_pred))$$

$$w_{\text{new}} = w + lr_rate * gw$$

return wnew

def gradient_descent(x, y, lr_rate, niter):

initialize b, w

100 rows
5 columns
(5w)

100x5
5x1

$\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_5 \end{bmatrix}$

```

b = random.random()
w = np.random.rand(x.shape[1])
w = w.reshape((x.shape, 1))
for i in range(niter):
    ypre = predict(x, w, b)
    w = update(x, w, ypre, y, lrate)
    print(w)
    print(y - ypre)

```

- Implement gradient descent algorithm
for given activation function

$$\begin{aligned}
 & \text{def gradient_descent}(x, y, w, b, lrate, niter) \\
 & \quad \text{for } i \in \text{range}(niter): \\
 & \quad \quad \text{ypre} = \text{forward}(x, w, b) \\
 & \quad \quad \text{err} = y - \text{ypre} \\
 & \quad \quad \text{grad_w} = \frac{1}{m} \sum_{j=1}^m \text{err}_j \cdot x_j \\
 & \quad \quad \text{grad_b} = \frac{1}{m} \sum_{j=1}^m \text{err}_j \\
 & \quad \quad w = w - lrate * grad_w \\
 & \quad \quad b = b - lrate * grad_b
 \end{aligned}$$

14/9/22

UNIT-IV

Multi Layer Perceptron

- Multi layer perception is also known as feed forward neural network.
- It contains one or more hidden layers between inputs and output layer.
- It can transform non-linear separable to linear separable, with the help of hidden layers.
- It can perform low-level feature extraction using hidden layers.
- Some of the tasks that can be performed using multilayer perception are object detection, handwritten digits recognition etc..
- Each hidden layer extracts some features such as in object detection 1st hidden layer detects horizontal line & 2nd hidden layer detects vertical lines and so on.
- The output of each hidden layer is passed to next hidden layer. The output layer gets aggregation of all these and it recognises the object.
- Advantages:-
 - MLP can perform complex tasks such as multiclass classification.
 - It can perform transformations.
 - Importance of hidden layers:-
 - Hidden layers in MLP let it to apply non-linear transformations so that the data is linear separable.
 - can change dimensionality of data.
 - extract low level features.

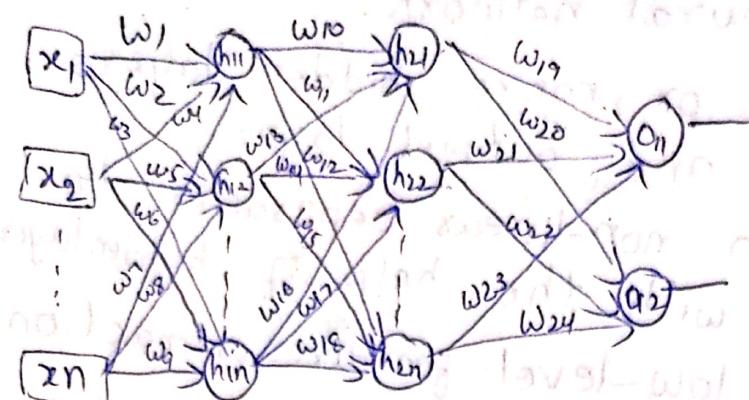
(Θ_1, Θ_2)
Single layer X
MLP

compute gradient

VI-TM0

-structure of MLP:-

Input layer hidden layer1 hidden layer2 Output layer



-conventions:-

$$i \rightarrow j$$

$\rightarrow w_{ji}$ → weight of the input from neuron i to j

$\rightarrow y_j$ → actual output of neuron j

$\rightarrow d_j$ → target output associated with neuron j

$\rightarrow e_j$ → error value for neuron j

$\rightarrow \psi_j(\cdot)$ → activation function for neuron j

$\rightarrow v_j$ → induced local field value for neuron j

$$\Rightarrow y_j = \psi_j(v_j)$$

$$v_j = \sum_{i=1}^n (w_{ji} \bullet x_i) + b$$

$$e_j = d_j - y_j$$

$$c(w_{ji}) = \frac{1}{2} \sum_{i=1}^n e_j^2$$

$$C_{avg}(w_{ji}) = \frac{1}{N} \sum_{i=1}^n e_j^2$$

* Implementation of MLP

sklearn

 | neural-network

 | | MLPClassifier

 | | | MLPClassifier

constructors of MLPClassifier

parameters

hidden_layer_sizes = tuple
① hidden_layer_sizes is $(100,)$ \Rightarrow 100 neurons & only 1 hidden layer by default.
default tuple is $(10, 5, 3)$ \Rightarrow then it creates 3 hidden layers in which 1st hidden layer has 10 neurons, 2nd has 5 neurons & 3rd has 3 neurons.
No. of neurons in output layer = No. of classes
No. of neurons in each hidden layer = No. of neurons in target variable.

No. of neurons in each hidden layer = No. of neurons in target variable.

To decide no. of neurons in each hidden layer - rule of thumb

hidden layers can have 70% (0.7) neurons

1st hidden layer = no. of input features

1st hidden layer = 70% of input features

2nd hidden layer = 70% of 1st hidden layer

2nd hidden layer = 70% of 1st hidden layer neurons (in 1st hidden layer)

activation

 | identity

 | logistic

 | tanh

 | relu (default)

③ solves

 | optimization algorithm.

 | To specify optimization algorithm.

 | lbgfs - Newton's

 | sgd - stochastic gradient descent

 | adam - variant of sgd (default)

④ max_iter

Specifies maximum no. of iterations
default it is 200

⑤ tol (tolerance value)

(Specifies stopping criteria)

code

```
from sklearn import datasets
```

```
d=datasets.load_iris()
```

```
d.data #returns array present in d  
(input features)
```

```
d.target #returns output feature array
```

```
d.feature_names #returns names of input feature
```

```
import pandas as pd
```

```
from sklearn.model_selection import train_test_split
```

```
from sklearn.neural_network import MLPClassifier
```

```
df=pd.DataFrame(d.data,columns=
```

```
["sepal length", "sepal width", "petal length",  
 "petal width"])
```

```
df.head()
```

```
df.info()
```

```
df.columns
```

```
x=df[["sepal length", "sepal width", "petal length",  
 "petal width"]]
```

```
y=d.target
```

```
xtr,xte,ytr,yte=train_test_split(x,y,
```

```
test_size=0.2,random_state
```

$mc = \text{MLPClassifier}(\text{hidden_layer_sizes}=(4, 3), \text{max_iters}=500,$
 $\text{activation}=\text{"identity"}, \text{solver}=\text{"sgd"})$

$mc.\text{fit}(x_{tr}, y_{tr})$

$mc.\text{predict}(x_{te})$

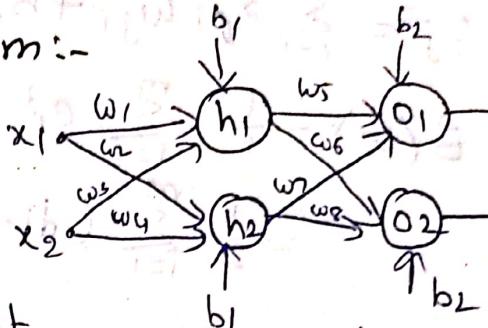
$mc.\text{score}(x_{te}, y_{te})$

~~17/9/22~~
* Back propagation algorithm:-

- 2 phases

① Forward pass

② Backward pass



① \Rightarrow network fed with input

- Neurons in hidden layers process the input and pass to next layers.

② \Rightarrow It is the error correction pass

- error value is computed.

- when there is an error, then back propagation algorithm uses gradient descent to minimize error value.

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

Computing gradient at o/p layer.

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

$$E_{O1} = \frac{1}{2} \sum_{i=1}^m e_i^2$$

$m = \text{Size of training set.}$

$$= \frac{1}{2} \sum_{i=1}^m (d_{O1} - y_{O1})^2$$

$$E_{O2} = \frac{1}{2} \sum_{i=1}^m e_i^2 = \frac{1}{2} \sum_{i=1}^m (d_{O2} - y_{O2})^2$$

20/9/22

At output layeroutput layer neuron j

$$e_j = d_j - y_j \rightarrow ①$$

$$y_j = \psi(v_j)$$

$$y_j = \text{sigmoid}(v_j) = \frac{1}{1+e^{-v_j}} \quad \begin{array}{l} \text{(If sigmoid function is used)} \\ \text{then, } v_j = \dots \end{array} \rightarrow ②$$

$$v_j = \sum_{i=1}^m x_i w_{ji} + b_0 \rightarrow ③$$

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

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

update rule

- as per gradient descent

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

(or)

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

g(n) is gradient of cost function

- a neural network with k neurons in output

$$\mathcal{E}_{\text{total}}(w) = \sum_{i=1}^k \mathcal{E}_i(w)$$

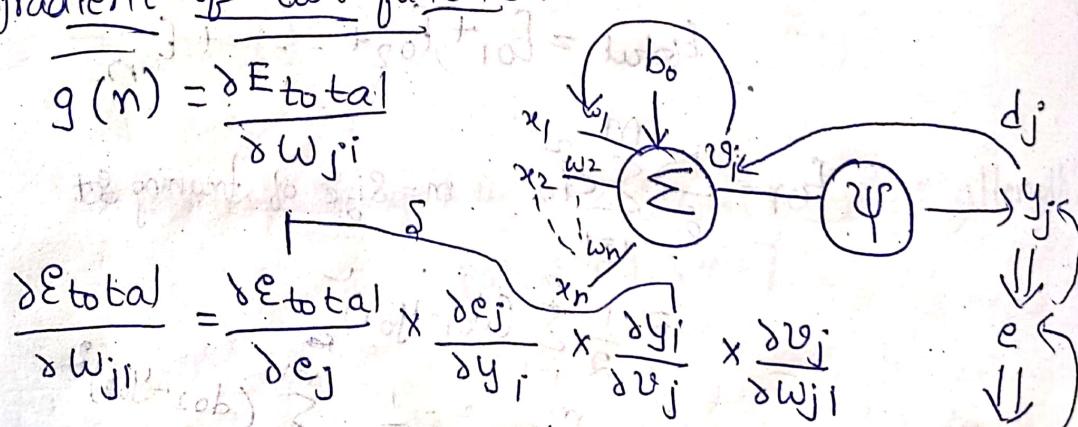
$$= \mathcal{E}_1(w) + \mathcal{E}_2(w) + \dots + \mathcal{E}_k(w) \rightarrow ⑤$$

gradient of cost function

$$g(n) = \frac{\partial \mathcal{E}_{\text{total}}}{\partial w_{ji}}$$

$$\frac{\partial \mathcal{E}_{\text{total}}}{\partial w_{ji}} = \frac{\partial \mathcal{E}_{\text{total}}}{\partial e_j} \times \frac{\partial e_j}{\partial y_i} \times \frac{\partial y_i}{\partial v_j} \times \frac{\partial v_j}{\partial w_{ji}}$$

(by applying the chain rule)

 E_{total} 

$$\frac{\partial \varepsilon_{\text{total}}}{\partial e_j} = \frac{\partial}{\partial e_j} (\varepsilon_1(w) + \varepsilon_2(w) + \dots + \varepsilon_k(w))$$

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

$$= e_j$$

$$\frac{\partial e_j}{\partial y_j} = \frac{\partial}{\partial y_j} (d_j - y_j) = -1$$

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

$$\frac{\partial y_i}{\partial v_j} = \frac{\partial}{\partial v_j} \left(\frac{1}{1+e^{-v_j}} \right) = \frac{(1+e^{-v_j})(0) - 1(e^{-v_j})}{(1+e^{-v_j})^2}$$

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

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

$$g(n) = \frac{\partial \varepsilon_{\text{total}}}{\partial w_{ji}} = e_j \times -1 \times y_j(1-y_j) \times x_i$$

$$= -e_j y_j(1-y_j) x_i$$

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

$$w_{\text{new}} = w_{\text{old}} + \eta y(1-y) e_j x_i$$

⑥

21/9/22

At hidden layer

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

$\Rightarrow g(n) \rightarrow$ gradient of cost function w.r.t weight t

$\epsilon_{\text{total}} =$ error caused by hidden layer neuron i at output layer neuron j

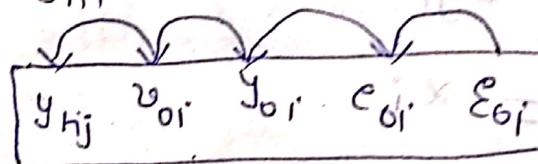
$$g(n) = \frac{\partial \epsilon_{\text{total}}}{\partial w_{ji}} \rightarrow ①$$

$$\frac{\partial \epsilon_{\text{total}}}{\partial w_{ji}} = \frac{\partial \epsilon_{\text{total}}}{\partial y_j} * \frac{\partial y_j}{\partial v_j} * \frac{\partial v_j}{\partial w_{ji}} \rightarrow ②$$

$$\epsilon_{\text{total}} \rightarrow y_j \rightarrow v_j \rightarrow w_{ji}$$

$$\frac{\partial \epsilon_{\text{total}}}{\partial y_j} = \frac{\partial \epsilon_{01}}{\partial y_j} + \frac{\partial \epsilon_{02}}{\partial y_j} + \dots + \frac{\partial \epsilon_{0k}}{\partial y_j} \rightarrow ③$$

$$\begin{matrix} \downarrow \\ \text{hidden layer} \\ \text{output} \end{matrix} = \frac{\partial \epsilon_{01}}{\partial y_{h1}} + \frac{\partial \epsilon_{02}}{\partial y_{h1}}$$



$$\frac{\partial \epsilon_{0i}}{\partial y_j} = \frac{\partial \epsilon_{0i}}{\partial v_{0i}} * \frac{\partial v_{0i}}{\partial y_{0i}} * \frac{\partial y_{0i}}{\partial w_{0i}} * \frac{\partial w_{0i}}{\partial y_j} \rightarrow ④$$

$$\epsilon_{0i} = \frac{1}{2} \sum_{i=1}^n c_{0i}^2 \quad \left| \begin{array}{l} \frac{\partial \epsilon_{0i}}{\partial y_{0i}} = \frac{\partial}{\partial y_{0i}} (\frac{1}{2} (y_{0i} - y_{0i})^2) = -1 \\ \frac{\partial y_{0i}}{\partial v_{0i}} = \frac{\partial}{\partial v_{0i}} (\varphi(v_{0i})) = y_{0i}(1 - y_{0i}) \end{array} \right.$$

$$\frac{\partial v_{0i}}{\partial w_{0i}} = c_{0i} \quad \left| \begin{array}{l} \frac{\partial v_{0i}}{\partial w_{0i}} = \frac{\partial}{\partial w_{0i}} (\varphi(w_{0i} y_{0i} + b_0)) = y_{0i} \end{array} \right.$$

$$\frac{\partial w_{0i}}{\partial y_j} = \frac{\partial}{\partial y_j} \left(\sum_{j=1}^m w_{ji} y_j + b_0 \right) = w_{ji}$$

$$\frac{\partial \varepsilon_{0i}}{\partial y_j} = \varepsilon_{0i} (-1) y_{0i} (y - y_{0i}) w_{ji}$$

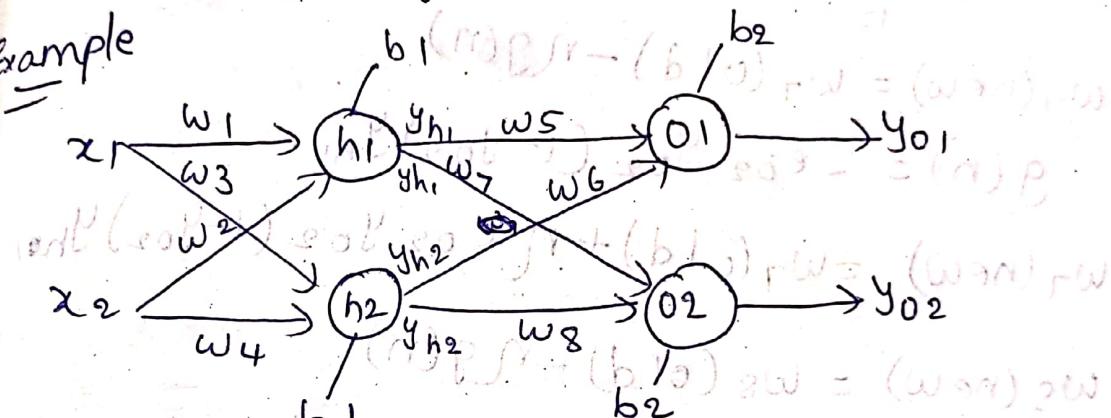
$$\cancel{\frac{\partial \varepsilon_{\text{total}}}{\partial y_j}} \quad \frac{\partial \varepsilon_{\text{total}}}{\partial y_j} = -\varepsilon_{01} y_{01} (y - y_{01}) w_{h11} \\ - \varepsilon_{02} y_{02} (y - y_{02}) w_{h12}$$

$$\frac{\partial y_j}{\partial \vartheta_j} = y_j (1 - y_j)$$

$$\frac{\partial \vartheta_j}{\partial w_{ji}} = x_i$$

$$\frac{\partial \varepsilon_{\text{total}}}{\partial y_j} = (-\varepsilon_{01} y_{01} (y - y_{01}) w_{h11} - \varepsilon_{02} y_{02} (y - y_{02}) w_{h12}) y_j (1 - y_j) x_i$$

Example



$$y_{01} = y_{h1} \times w_5 + y_{h2} \times w_6 + b_2$$

$$y_{01} = \text{sigmoid}(y_{01}) = \frac{1}{1 + e^{-y_{01}}}$$

$$\varepsilon_{01} = d_{01} - y_{01}$$

$$E_{01} = \frac{1}{2} \sum_{i=1}^n \varepsilon_{01}(i)$$

$$= \frac{1}{2} \sum_{i=1}^n (d_{01}(i) - y_{01}(i))^2$$

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

$$g(n) = -e_{01} y_{01} (1-y_{01}) y_{h1}$$

$$w_5(\text{new}) = w_5(\text{old}) + \eta e_{01} y_{01} (1-y_{01}) y_{h1}$$

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

$$g(n) = -e_{01} y_{01} (1-y_{01}) y_{h2}$$

$$v_{02} = y_{h1} \times w_7 + y_{h2} \times w_8 + b_2$$

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

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

$$\sum_{i=1}^n e_{02}^2(i) = \frac{1}{n} \sum_{i=1}^n (d_{02}^{(i)} - y_{02}^{(i)})^2$$

$$= \frac{1}{2} \sum_{i=1}^n (d_{02}^{(i)} - y_{02}^{(i)})^2$$

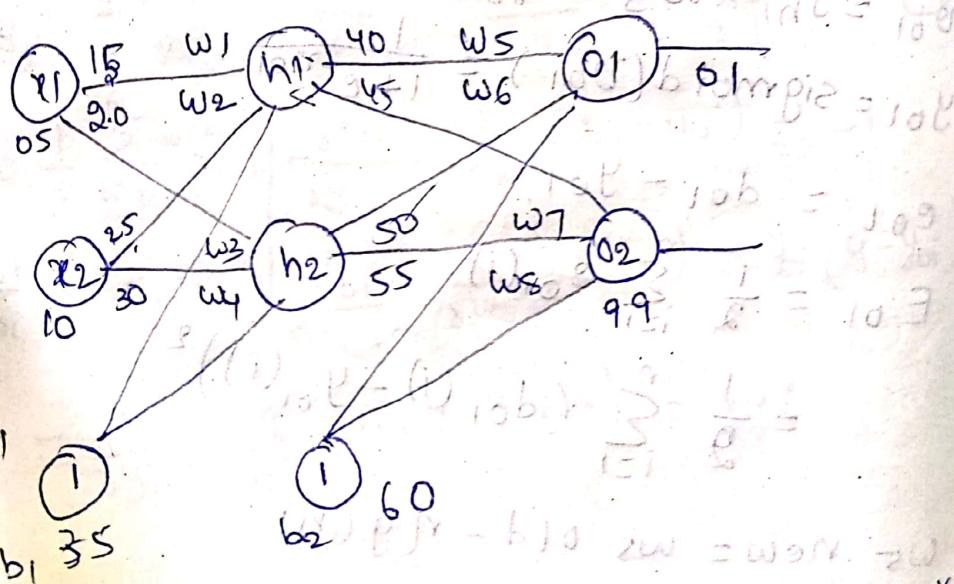
$$w_7(\text{new}) = w_7(\text{old}) - \eta g(n)$$

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

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

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

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



Iteration 1

Forward pass

at h_1

$$v_{h_1} = w_1 x_1 + w_2 x_2 + b_1$$

$$= 0.15 \times 0.05 + 0.25 \times 0.10 + 0.35$$

$$= 0.3825$$

$$y_{h_1} = \varphi(v_{h_1}) = \frac{1}{1 + e^{-v_{h_1}}} = 0.5965$$

at h_2

$$v_{h_2} = w_3 x_1 + w_4 x_2 + b_1$$

$$= 0.20 \times 0.05 + 0.30 \times 0.10 + 0.35$$

$$= 0.5980.39$$

$$y_{h_2} = \varphi(v_{h_2}) = 0.5912$$

	d	y
01	0.01	0.7569
02	0.99	0.767

at o1

$$v_{o1} = y_{h_1} w_5 + y_{h_2} w_7 + b_2$$

$$= 0.5962 \times 0.40 + 0.5962 \times 0.5 + 0.60$$

$$= 1.1359$$

$$y_{o1} = \varphi(v_{o1}) = 0.7569$$

at o2

$$v_{o2} = y_{h_1} w_6 + y_{h_2} w_8 + b_2$$

$$= 0.5962 \times 0.45 + 0.5962 \times 0.55 + 0.60$$

$$= 1.1954$$

$$y_{o2} = \varphi(v_{o2}) = 0.767$$

Backward pass

Iteration 1 Corp (Layer 1)

VSK

e_f

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

$$g(n) = \frac{\partial E_{\text{total}}}{\partial w_5}$$

$$\frac{\partial E_{\text{total}}}{\partial w_5} = \frac{\partial E_{\text{total}}}{\partial e_{01}} \times \frac{\partial e_{01}}{\partial y_{01}} \times \frac{\partial y_{01}}{\partial v_{01}} \times \frac{\partial v_{01}}{\partial w_5}$$

$$= -e_{01} \times -1 \times y_{01} \times (1 - y_{01}) \times y_{h1}$$

$$= -e_{01} \times y_{01} \times (1 - y_{01}) \times y_{h1}$$

$$= -(d - y) (y_{01} \times (1 - y_{01})) y_{h1}$$

$$= -(0.01 - 0.7569) (0.7569 \times (1 - 0.7569)) (0.5945)$$

$$\approx 0.0817$$

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

$$= 0.40 - 0.05 \times 0.0817 = 0.3591$$

Updating w_7

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

$$g(n) = \frac{\partial E_{\text{total}}}{\partial w_7}$$

$$\frac{\partial E_{\text{total}}}{\partial w_7} = \frac{\partial E_{\text{total}}}{\partial e_{01}} \times \frac{\partial e_{01}}{\partial y_{01}} \times \frac{\partial y_{01}}{\partial v_{01}} \times \frac{\partial v_{01}}{\partial w_7}$$

$$= e_{01} \times (-1) \times (y_{01} \times (1 - y_{01})) \times y_{h2}$$

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

$$= -(0.01 - 0.7569) \times 0.7569 \times (1 - 0.7569) \times 0.59612$$

$$= 0.0819$$

$$w_6(\text{new}) = 0.50 - 0.5 \times 0.0819$$

$$= 0.4590$$

At O2

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

$$g(n) = \frac{\partial \varepsilon_{\text{total}}}{\partial w_6}$$

$$\frac{\partial \varepsilon_{\text{total}}}{\partial w_6} = \frac{\partial \varepsilon_{\text{total}}}{\partial e_{02}} \times \frac{\partial e_{02}}{\partial y_{02}} \times \frac{\partial y_{02}}{\partial v_{02}} \times \frac{\partial v_{02}}{\partial w_6}$$

$$= e_{02} \times (-1) \times y_{02} (1-y_{02}) \times y_{h2}$$

$$= (d_{02} - y_{02}) (-1) y_{02} (1-y_{02}) \times y_{h2}$$

$$= -(0.99 - 0.767) (0.767) (1-0.767) (0.5945)$$

$$= -0.0235$$

$$w_6(\text{new}) = 0.45 - (0.5) (-0.0235)$$

$$= 0.45 + 0.01175$$

$$= 0.46175$$

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

$$g(n) = \frac{\partial \varepsilon_{\text{total}}}{\partial w_8}$$

$$\frac{\partial \varepsilon_{\text{total}}}{\partial w_8} = \frac{\partial \varepsilon_{\text{total}}}{\partial e_{02}} \times \frac{\partial e_{02}}{\partial y_{02}} \times \frac{\partial y_{02}}{\partial v_{02}} \times \frac{\partial v_{02}}{\partial w_8}$$

$$= -(d_{02} - y_{02}) y_{02} (1-y_{02}) y_{h2}$$

$$= -0.0234$$

$$w_8(\text{new}) = 0.55 - (0.5) (-0.0234)$$

$$= 0.5617$$

At hidden layer

at h)

=
update of w)

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

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

$$\frac{\partial E_{\text{total}}}{\partial w_i} = \frac{\partial E_{\text{total}}}{\partial y_{h_i}} \times \frac{\partial y_{h_i}}{\partial \theta_{h_i}} \times \frac{\partial \theta_{h_i}}{\partial w_i}$$

$$\frac{\partial E_{\text{total}}}{\partial y_{h_i}} = \frac{\partial E_01}{\partial y_{h_i}} + \frac{\partial E_02}{\partial y_{h_i}}$$

$$\frac{\partial E_01}{\partial y_{h_i}} = \frac{\partial E_01}{\partial e_{01}} \times \frac{\partial e_{01}}{\partial y_{01}} \times \frac{\partial y_{01}}{\partial v_{01}} \times \frac{\partial v_{01}}{\partial y_{h_i}}$$

$$\frac{\partial E_02}{\partial y_{h_i}} = \frac{\partial E_02}{\partial e_{02}} \times \frac{\partial e_{02}}{\partial y_{02}} \times \frac{\partial y_{02}}{\partial v_{02}} \times \frac{\partial v_{02}}{\partial y_{h_i}}$$

Continuation in papers

* Practical & design issues of Back propagation

algorithm:-

① Rate of learning

- specifies a step size taken by learning

- algorithm

- takes a value in between 0-1

$0 \leq \eta \leq 1$

η is very small

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

learning rate

slowly converges to global minimum

It increase no. of iterations required for converging.

- algorithm can trapped into local minimum point.

n very large
algorithm can overshoot global minimum point

- can oscillate

- avoids getting trap into local minima point.

solution

Momentum:

$$\Delta w_{ji}(n) = \alpha \Delta w_{ji}(n-1) + \Delta w_{ji}^*(n)$$

- provides solution 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 samples

examples or samples

also called online mode.

- update for weight vector computed

sequentially by considering only

one sample at a time.

less memory space as it

- It requires slow convergence towards global

minima point.

drawback

- It is theoretically difficult to determine stopping criteria.

1/10/22

② Stopping Criteria

- When model produces errors that is accepted.

- Euclidean norm of $y-d$ reaches to minimum value.

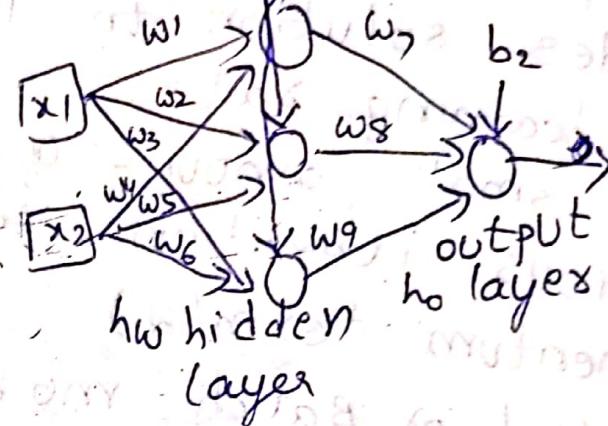
- Training set.

③ New sample not there in training set.

- Generalization. (By creating power set)

Implementation of Backpropagation algorithm.

$$h_w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \\ w_6 \end{bmatrix} \quad h_o = \begin{bmatrix} w_7 \\ w_8 \\ w_9 \end{bmatrix} \quad 3 \times 1$$



$$v_{h1} = x_1 w_1 + x_2 w_4$$

$$v_{h2} = x_1 w_2 + x_2 w_5$$

$$v_{h3} = x_1 w_3 + x_2 w_6$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad h_w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \\ w_6 \end{bmatrix}$$

11/10/22 Design Issues in Backpropagation Algorithm

1) Learning rate - due to this algorithm converges slowly or at fast rate

low

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

* alg stuck in local min point

* momentum's } applied to alg and alg move forward

from local min point and to reach global min point

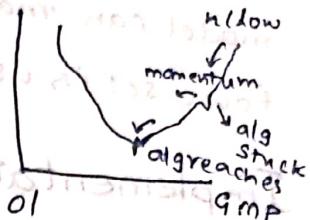
(i.e., update value of previous iteration is added to update formula changes)

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

α -momentum term

lies b/w 0 to 1

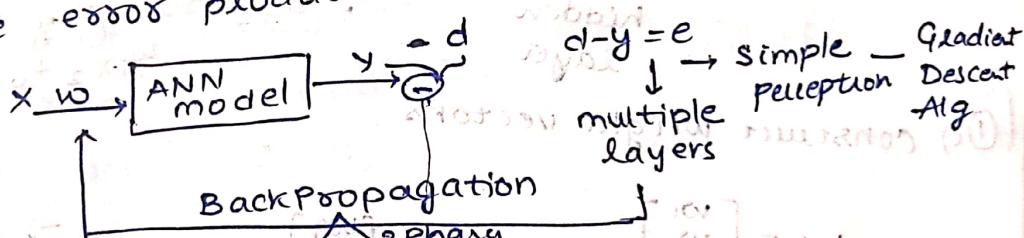
mostly $\alpha = 0.001$ or 0.01



α is multiplied with $\Delta w(n-1)$ because if $\Delta w(n-1)$ is too high then alg takes fewer & 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.



goes to goes to
iteration feed forward Phase
Process. (read IIP, produces OLP &
{identifies error})

This process occurs iteratively.

If no stopping criteria exists, then Back propagation occurs infinitely.

* when $y=d$ - but leads to infinite loop i.e., BP General Stopping Criteria used is to stop

Alg is stopped when error value reaches a particular threshold value. (98% accepted, 95% accurate.)

When euclidean norm

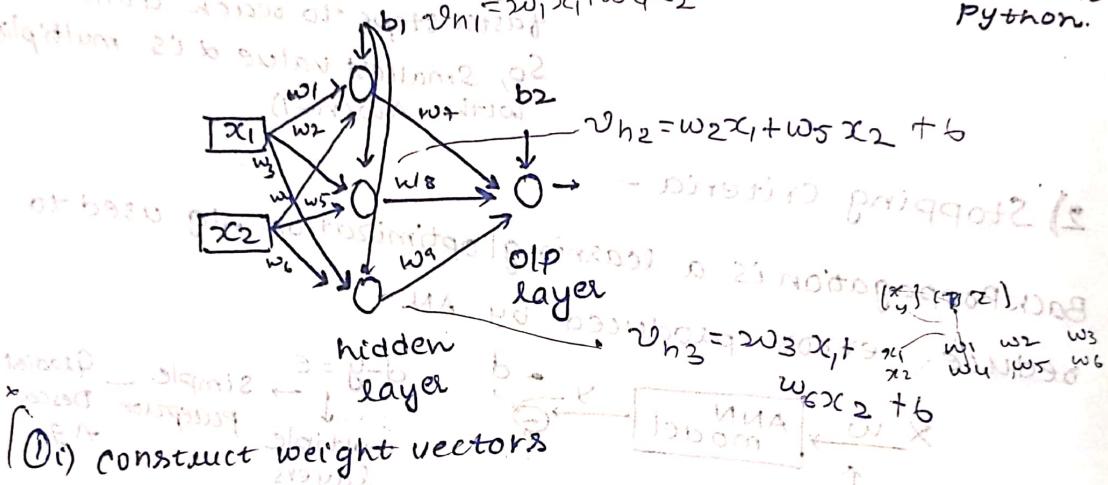
of d & y reaches to a min. value.

- 3) new samples not there in training set
- back propagation goes through Generalization
 - Process (It tries to match data with already trained data and conclusions)
 - (compares) given sample (with previous sample)
 - make conclusions based on it.
 - * Powerset is created for the training data using which model is trained. Powerset contains all features of dataset so that based on it, model can make predictions for new data to.
 - Powerset is used in Generalisations.

Implementation of Backpropagation Algorithm

Before implementation following tasks must be done-

- ① initialize free parameters (w, b, η) # random package randint method in Python.



- ② construct weight vectors

$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \\ w_6 \\ w_7 \\ w_8 \\ w_9 \end{bmatrix} \quad h_0 = \begin{bmatrix} u_{h1} \\ u_{h2} \\ u_{h3} \end{bmatrix} \quad h_1 = \begin{bmatrix} u_{h4} \\ u_{h5} \end{bmatrix}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

- ③ Construct neural network models.
- no. of inputs
 - no. of hidden layers
 - no. of neurons in each h.l.
 - no. of neurons in o.p. layer.

- ④ Compute error value

$$e = d - y$$

Update the weights.

$$\begin{bmatrix} w_1 & w_4 \\ w_2 & w_5 \\ w_3 & w_6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

ip vector/matrix
size: 3x2
no. of ip's
no. of neurons
in h.l x no. of ip's

$$w_l = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}_{3x1}$$

$$\begin{bmatrix} y_{h1} \\ y_{h2} \\ y_{h3} \end{bmatrix}_{3x1}$$

$$(w_7 w_8 w_9)_{3x1}$$

Program:

Initializing parameters:

import numpy as np

inputsize=2

hiddensize=3

outputszie=1

lr=0.1

w1=np.random.randn(inputsize, hiddensize)

w2=np.random.randn(hiddensize, outputszie)

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

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

define activation fn:

```
def sigma(z):
    y = 1/(1+np.exp(-z))
    return y.
```

def derivate(y):

return y*(1-y)

def mapear(yP, y):

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

return E

for loop

define a dataset

```
x = np.array([[7, 2], [2, 6], [9, 1]])  
y = np.array([80, 50, 95])
```

hrs studied	hrs enjoyed	result
7	2	80
2	6	50
9	1	95

normalise

$$x = x / \text{np. amax}(x, \text{axis}=0)$$

$$y = y / 100$$

create a neural net and train

```
def train(x, y):
```

global w1, w2, b1, b2, l8

forward phase

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

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

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

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

return y_o

train(x, y)

Back propagation phase:

find error value $(d - y)$:

$$e_1 = y_o - y$$

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

$$e_2 = \delta_{-o} @ w_2.T$$

$$\delta_{-h} = e_2 * y_h * (1 - y_h)$$

$$g_o = y_h @ \delta_{-o}$$

final gradient

```

 $g_n = x \cdot T @ \delta_{\text{tanh}}$ 
 $w_1 = w_1 - \eta \cdot g_n$ 
 $w_2 = w_2 - \eta \cdot g_o$ 
 $b_1 = b_1 - \eta \cdot \delta_{\text{tanh}}$ 
 $b_2 = b_2 - \eta \cdot \delta_{\text{o}}$ 

# updating weights
& bias

for i in range(2000):
    train(x,y)

def forwardtest(x,y):
    global w1, w2, b1, b2
    # forward phase
     $v_n = \text{np.dot}(x, w_1) + b_1$  ( $x @ w_1$ ) +  $b_1$ 
     $y_n = \text{sigm}(v_n)$ 
     $v_o = \text{np.dot}(y_n, w_2) + b_2$  ( $y_n @ w_2$ ) +  $b_2$ 
     $y_o = \text{sigm}(v_o)$ 
    return y_o

forwardtest(x,y)

```

UNIT-IV Regression using MLP

- independent variables - o/p variables

- dependent variables - i/p variables.

Regression-

• It is a technique used to determine relationship b/w independent and dependent variables.

• 3 techniques -

① simple linear regression:

• determines relationship between independent & dependent variables.

$$y = \beta_0 + \beta_1 x$$

x - independent variable

y - dependent variable

β_0 - y intercept

β_1 - slope

Eg:	\underline{x}	\underline{y}	$\underline{x-\bar{x}}$	$\underline{y-\bar{y}}$	$\frac{(x-\bar{x})(y-\bar{y})}{\sum(x-\bar{x})(y-\bar{y})}$	$\frac{(x-\bar{x})^2}{\sum(x-\bar{x})^2} (y-\bar{y})^2$
	1	2	-2	-2	4	4
	2	4	-1	0	0	0
	3	5	0	1	0	1
	4	4	1	0	0	0
	5	5	2	1	2	1
	$\bar{x} = 3$		$\bar{y} = 4$			

$$\beta_1 = \frac{\sum(x-\bar{x})(y-\bar{y})}{\sum(x-\bar{x})^2} = \frac{6}{10} = 0.6$$

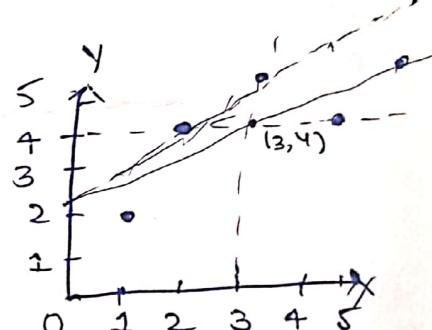
$$\bar{y} = \beta_0 + \beta_1 \bar{x}$$

$$4 = \beta_0 + 0.6 \times 3$$

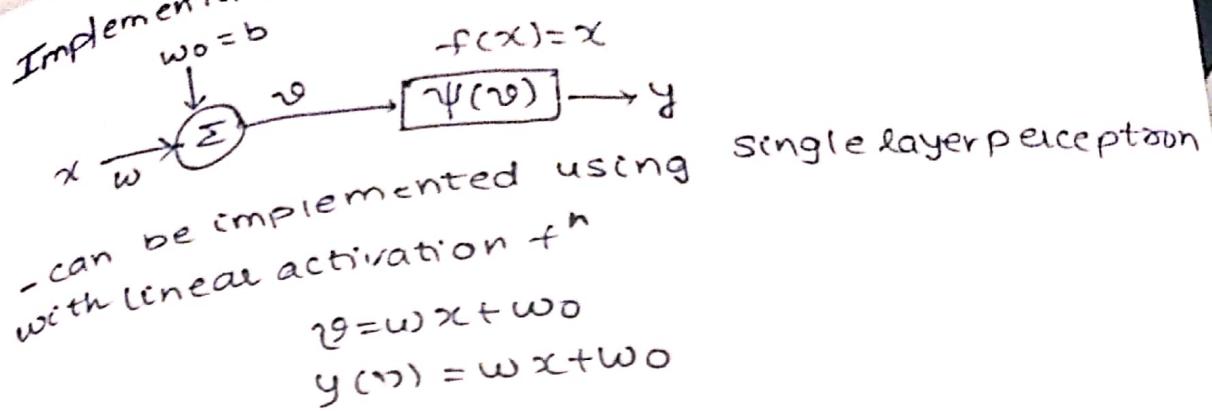
$$\beta_0 = 4 - 0.6 \times 3$$

$$= 4 - 1.8$$

$$\beta_0 = 2.2$$



Implementation using perceptron -



- can be implemented using
with linear activation +
 $\varphi = w^T x + w_0$
 $y(\varphi) = w^T x + w_0$