

Decision trees

①

Classification is the process of dividing the datasets into different categories or groups by adding label.

Eg:- classifying spam or not spam mail.

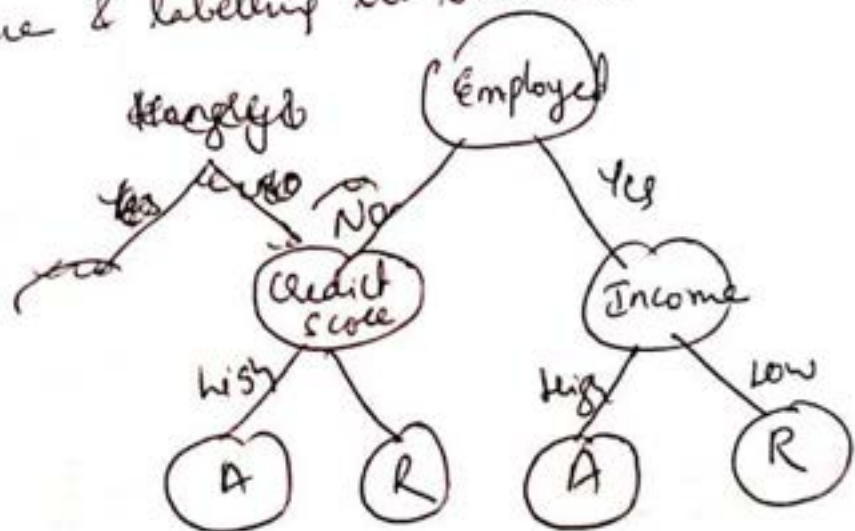
There are different types of classification Algos

- 1) Decision tree
- 2) Naive Bayes classifier
- 3) KNN Classifier
- 4) Random forest

A decision tree is a graphical representation of all possible solutions to a decision based on certain conditions.

A decision tree is a hierarchical model for Supervised Learning. A decision tree is composed of internal decision nodes & terminal leaves. Each decision node implements a test function with discrete outcome & labelling the branches.

Eg:-
To sanction
loan or not



Decision tree is a tree where each node represents ⁽²⁾ a feature (attribute), each link (branch) represents a decision rule and each leaf represent an outcome.

Given a dataset and we have to make a decision tree that predicts whether Eg:- tennis will be played on the day?

Step 1:- create a root node

How to choose the root node?

The attribute that best classifies the training data, use this attribute at the root of the tree.

Step 2 How to choose the best attribute

~~then~~ - Use the ID3 Algm

It has two Important concepts.

(a) Entropy:- which calculates amount of uncertainty in dataset [It means number of positive and number of negative examples]

Note:- If we have equal number of positive and negative examples then Entropy = 1

If we have only positive example or only negative example then entropy = 0

$$\text{Entropy} = -\frac{P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

Weighted Information Entropy

(3)

$$I(\text{Avg Information Entropy})(\text{attribute}) = \sum \frac{P_i + n_i}{P + N} \text{Entropy}(A)$$

Information gain :- Difference in Entropy before & after splitting dataset on attribute A.

$$\text{Gain} = \text{Entropy}(S) - I(\text{attribute}).$$

good quantitative measure that measures how well given attribute separates the training examples according to their target classification.

Steps:

- 1) Compute the entropy for dataset $\text{Entropy}(S)$
- 2) for every attribute/feature
 - (a) Calculate entropy for all other values
 - (b) Take average Information Entropy for current attribute
 - (c) Calculate gain for current attribute.
 - (d) Pick the highest gain attribute.
 - (e) Repeat steps (b)(c)(d) until we get the tree we desired.

ed on the attribute (features) given whether we get profit in the company or not (up or down) (4)

	Age	Competation	Type	profit
1)	old	Yes	Software	down
2)	old	NO	Software	down
3)	old	No	Hardware	down
4)	mid	Yes	Software	down
5)	mid	Yes	Hardware	down
6)	mid	NO	Hardware	up
7)	mid	NO	Software	up
8)	New	Yes	Software	up
9)	New	NO	Hardware	up
10)	New	NO	Software	up

No of Positive examples \rightarrow profit goes up $P=5$
 No of Negative examples \rightarrow profit goes down $N=5$

Entropy of class

$$\begin{aligned} \text{Entropy (S)} &= -\frac{P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{N+N} \log_2 \left(\frac{N}{P+N} \right) \\ &= -\frac{5}{10} \log_2 \left(\frac{5}{10} \right) - \frac{5}{10} \log_2 \left(\frac{5}{10} \right) = 1 \end{aligned}$$

To find which attribute from the dataset becomes the root node find.

for each attribute find

\rightarrow Entropy

\rightarrow Average Information Entropy

\rightarrow Gain.

(4)

	Pi	Ni	Entropy
old	0	3	0
New	3	0	0
Mild	2	2	1

Entropy of age
(Average Information Entropy)

$$I = \sum \frac{P_i + N_i}{P+N} \text{Entropy}(A)$$

$$= \frac{0+3}{5+5} (0) + \frac{3+0}{5+5} (0) + \frac{2+2}{5+5} (1) = 0.4$$

$$\text{Gain} = \text{Entropy}(S) - \text{Entropy}(\text{age})$$

$$= 1 - 0.4 = 0.6$$

Competition	P	N	Entropy	
Yes	1	3	0.8127	$\rightarrow -\frac{P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$
No	4	2	0.9182	$\rightarrow -\frac{1}{4} \log_2 \left(\frac{1}{4} \right) - \frac{3}{4} \log_2 \left(\frac{3}{4} \right)$
				$\rightarrow -\frac{4}{6} \log_2 \left(\frac{4}{6} \right) - \frac{2}{6} \log_2 \left(\frac{2}{6} \right)$

Entropy Competition

$$= \frac{(1+3)}{(5+5)} \times (0.8127) + \frac{(4+2)}{(5+5)} (0.918295)$$

$$= 0.8754$$

$$\text{Gain} = \text{Class Entropy} - \text{Entropy(Competition)}$$

$$= 1 - 0.8754$$

$$\boxed{\text{Gain} = 0.12455}$$

	P	N	Entropy
Software	3	3	1
Hardware	2	2	1

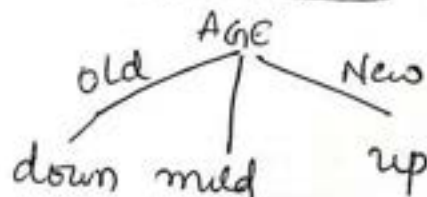
Entropy(Type)

$$\text{Avg Inf} = \frac{3+3}{5+5} (1) + \frac{2+2}{5+5} (1) = \frac{6}{10} + \frac{4}{10} = \frac{10}{10} = 1$$

$$\text{Gain} = \text{Class Entropy} - \text{Entropy(Type)} \\ = 1 - 1 = 0$$

Information Gain:-

Age	0.6	
Compete	0.12451	
Type	0	



Next Repeat

Age	competition	Type	Profit
mild	Yes	Software	Down
mild	Yes	Hardware	Down
mild	No	Hardware	up
mild	No	Software	up

} Current dataset

Entropy for the class

$$P=2 \quad N=2$$

$$\text{Entropy(Profit)} = 1$$

(6)

Competition

	P	N	Entropy
Yes	0	2	0
NO	2	0	0

$$\text{Entropy(Competition)} \\ \text{Average In} = \frac{2}{4}(0) + \frac{2}{4}(0) = 0$$

$$\text{Gain} = \text{Entropy}_{\text{class}} - \text{Entropy(Competition)} \\ = 1 - 0 = 1$$

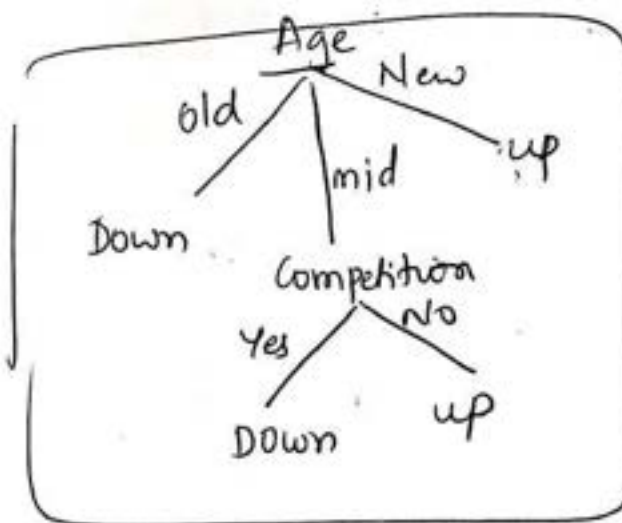
Type

	P	N	Entropy
Soft	1	1	1
Hard	1	1	1

$$\text{Average In Entropy(Type)} \\ = \frac{2}{4}(1) + \frac{2}{4}(1) = \frac{2}{4} + \frac{2}{4} = 1$$

Gain

Competition	1
Type	0



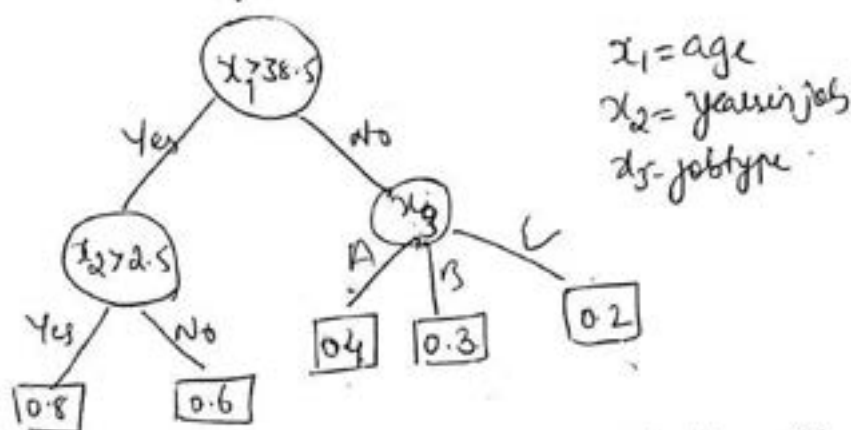


Rule Extraction from a Decision tree

(1)

One rule is created for each path from the root to a leaf.

- Each attribute-value pair along a path forms conjunction. The leaf holds the class prediction.



Rules

R_1 : If (age > 38.5) and (years-in-job > 2.5) then $y = 0.8$

R_2 : If (age > 38.5) and (years-in-job ≤ 2.5) then $y = 0.6$

R_3 : If (age ≤ 38.5) and (job-type = 'A') then $y = 0.4$

R_4 : If (age ≤ 38.5) and (jobtype = 'B') then $y = 0.3$

R_5 : If (age ≤ 38.5) and (jobtype = 'C') then $y = 0.2$

Interpretability: Another main advantage of decision trees is interpretability.

Learning rules from Data

Sequential covering Algm.

→ Rules are learned sequentially, where each rule for a given class will ideally cover many of the class tuples.

Pruning

- Stop splitting a node if a decision based on too few instances causes variance and thus generalization error. Stopping tree construction early on before is called Prepruning the tree.
- ④ Eg: If you have 5 instances and split into 2 & 3 instances is not a significant split.

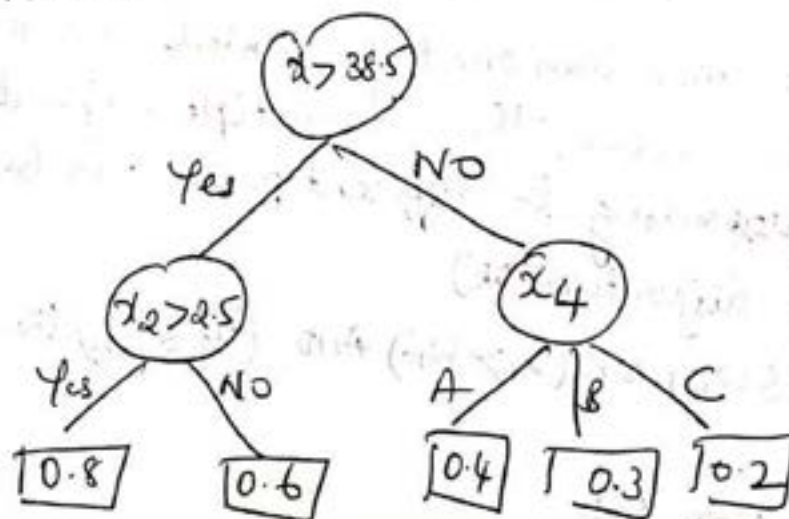
Post Pruning: Split the entire dataset into training set and testing set.

Train the training set for decision tree and run the testing set on the decision tree, and accuracy is not 100%, because the tree is overfit to the training data.

- Remove the node one by one from tree
- Measure the performance on test set.
- Remove the node that results in least improvement.
- Repeat until further pruning is harmful.

Rule Extraction from a Decision tree

- One rule is created for each path from the root to a leaf
- Each attribute-value pair along a path forms a conjunction; the leaf holds the class prediction.
- Rules are mutually exclusive



$x_1 = \text{age}$
 $x_2 = \text{Years in job}$
 $x_3 = \text{job type}$

- R_1 : If (age > 38.5) and (Years in job > 2.5) then $y = 0.8$
 R_2 : If (age > 38.5) AND (Years in job \leq 2.5) then $y = 0.6$
 R_3 : If (age \leq 38.5) AND (job-type = 'A') then $y = 0.4$
 R_4 : If (age \leq 38.5) AND (job-type = 'B') then $y = 0.3$
 R_5 : If (age \leq 38.5) AND (job-type = 'C') then $y = 0.2$

Interpretability: - Another main advantage of decision trees is interpretability. The decision nodes carry conditions that are simple to understand. Each path from root to a leaf corresponds to one conjunction of tests, as all those conditions should be satisfied to reach to the leaf. These paths together can be written down as a

Set of IF-THEN rules, called a rule base.
Rule base is a knowledge extraction, which is the model learned from data.

- Rule support indicates the percentage of training data covered by the rule.
- Rule reflects the main characteristics of the dataset. They show the important features and split positions.
- There may be more than one leaf labeled with the same class. In such case, these multiple conjunction expression corresponding to different paths can be combined as disjunction (OR).
$$\text{If } (x \leq w_{10}) \text{ OR } (x_1 > w_{10}) \text{ AND } (x_2 \leq w_{22}) \text{ then } c_1$$

Learning rules from Data

Sequential Covering Algm:- Rules are learned sequentially, where each rule for a given class will ideally cover many of the class tuples. A covering algm is an algm that develops a cover taking into account only set of positive examples but none of negative examples. It learns one rule at a time and repeat this process to gradually cover the full set of positive examples.

Steps

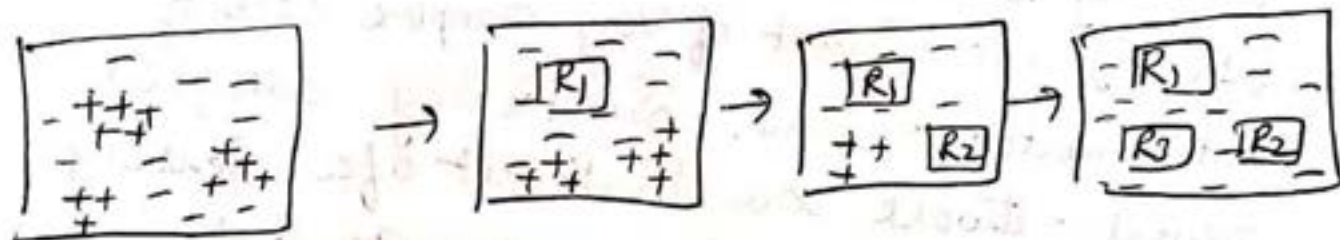
- 1) Starts with empty cover (no rules generated).
- 2) Rules are learned one at a time. Rule is said to cover an example if the example satisfies all the conditions of the rule.
- 3) Once a rule is obtained, all the training

examples covered by the rule are removed from the training set.

4) Goto step 2 until enough rules are added.

Example

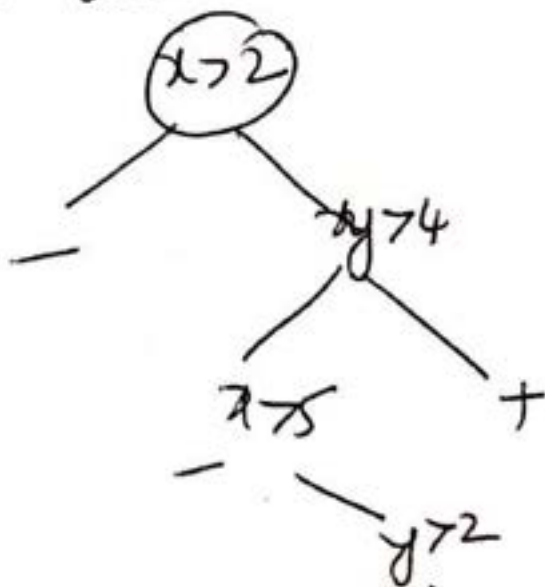
Rippee :- One example of rule induction is Rippee. We start with the case of positive and negative examples. Rules are added to explain positive examples. Other instances which are not covered by the rule is negative examples.



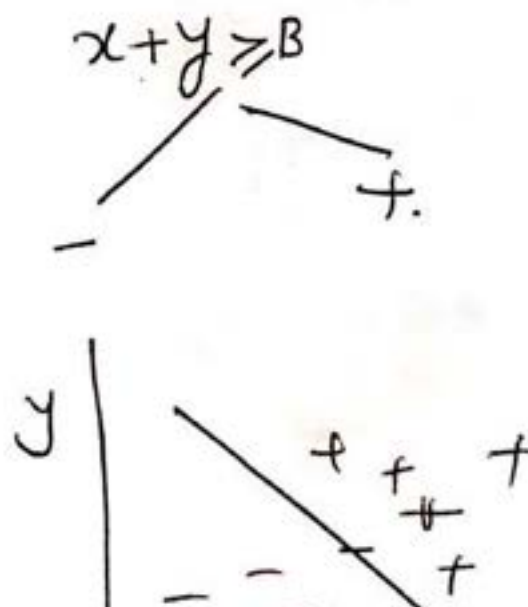
Multivariate trees :-

• Multivariate Decision trees are able to generalize well when dealing with attributes correlation different from univariate.

Ex :- univariate



Multivariate



Artificial neuron is called perceptron

(3)

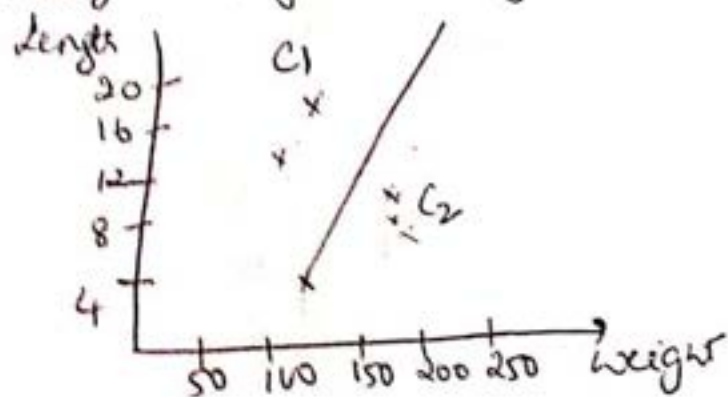
(1) Single Layer perceptron

	Weight (grams)	Length (cm)
Fruit 1 (class 1)	121	16.8
	114	15.2
Fruit 2 (class 2)	210	9.4
	195	8.2

Sample data obtained from two different fruit.

→ Train a single layer perceptron model using the above parameters to classify the two fruits.

Using the model parameters we can classify the fruit with weight 140 gm & length 17.9 cm.



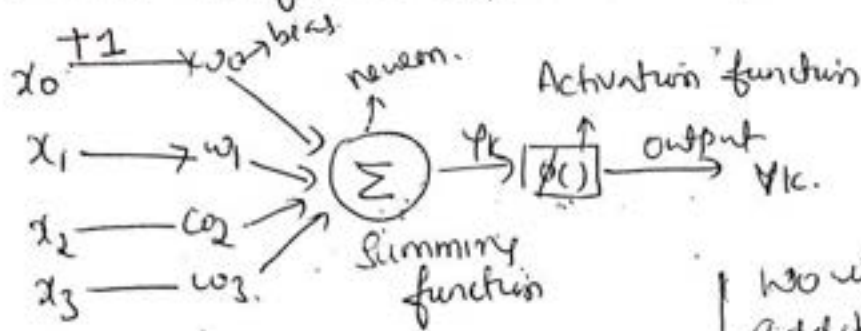
We found a decision boundary which classifies the two classes.

If new weight and height is given, this model is able to predict C₁ and C₂.

How will implement neural network for the above problem is shown in the next page.

Perceptrons [Perceptrons training] Single layer perceptron

One type of ANN system, is based on a unit called perceptron. A perceptron takes a vector of real valued inputs, calculates the linear combination of these inputs and outputs.



$$V_k = \sum_{j=0}^n w_j x_j + w_0 \quad \text{--- (1)}$$

w_0 is the bias is additional parameter which is used to adjust the output. Help the model to best fit the given data.

Activation is a threshold function.

$$\phi(x) = \phi V = \begin{cases} 1 & \text{if } V \geq 0 \\ 0 & \text{if } V < 0 \end{cases} \quad \text{or} \quad \begin{cases} 1 & \text{if } V \geq 0 \\ -1 & \text{if } V < 0 \end{cases}$$

where each w_i in eqn (1) is a real valued constant or weight, that determines the contribution of input x_i to perceptron output.

→ Perceptron is the simplest form of a neural network used for the classification of patterns said to be linear separable.

Training a perceptron

5

(b) $x(n) = [+1, x_1(n), x_2(n), \dots, x_m(n)]$

$w(n) = [b(n), w_1(n), w_2(n), \dots, w_m(n)]$

$b(n) \rightarrow$ bias, $y(n) \rightarrow$ actual response, $d(n) \rightarrow$ desired response

$\eta \rightarrow$ learning rate \rightarrow It is degree to which weights are changed at each step

Steps:-

(1) Initialization: Initialize weights & inputs.

(2) Activation: At time step n , activate the perceptron by applying input vector $x(n)$ and desired response $d(n)$

(3) Computation of actual response of the perceptron.

$$y(n) = \text{sgn}[w^T(n) \cdot x(n)] \quad \text{--- (1)}$$

where $\text{sgn}(\cdot)$ is the signum function (Activation function)

$$\text{sgn}(x) = \begin{cases} +1 & \text{if } x \geq 0 \\ -1 & \text{if } x < 0 \end{cases} \quad \text{--- (2)}$$

(4) update the weight vector of the perceptron.

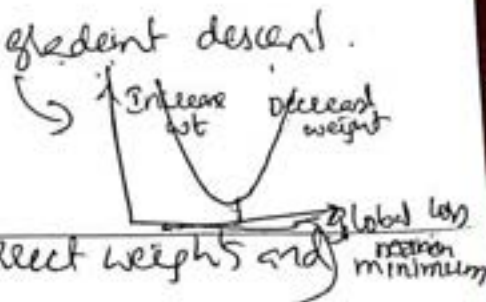
$$w(n+1) = w(n) + \eta [d(n) - y(n)] x(n) \quad \text{--- (3)}$$

(5) Go back to step 2 [retraining]

Note - (i) One way to learn acceptable weight vector is to begin with random weights then iteratively apply the perceptron to each training examples, modify the perceptron weights whenever it misclassifies an example.

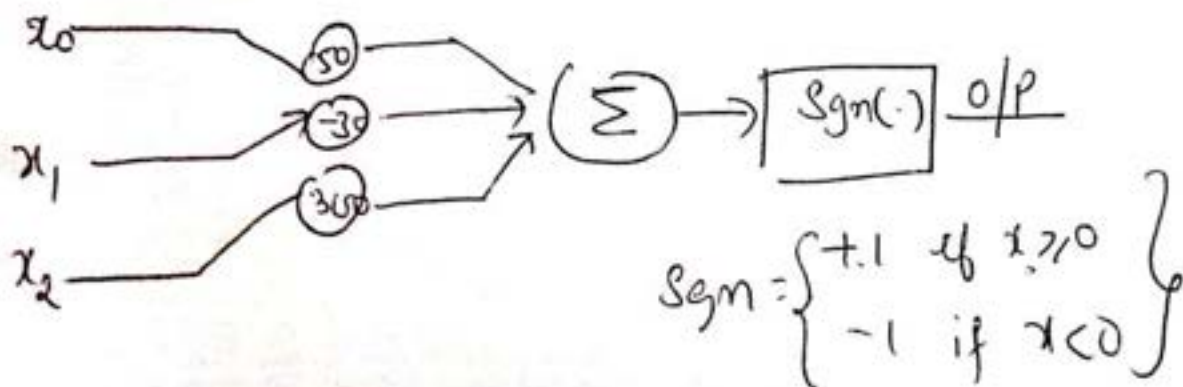
(6)

- (b) This process is repeated iteratively through the training examples correctly.
- (c) Weights are modified at each step according to the perceptron training rule, which reuses the weight associated with input x_i .
- (d) Suppose training example is correctly classified ($d(n) - y(n) = 0$) then no weight are updated.
- (e) Suppose perceptron outputs -1 , when target $d(n)$ is $+1$. To make perceptron output $+1$ instead of -1 , weights must be altered. Increasing w_i will bring the perceptron closer to correctly classifying this example.
- (f) If $d(n) = -1$ and $y(n) = +1$, then weight will be decreased rather than increased.
- This is also known as stochastic gradient descent.



Example 1

Initially Let us assume we know the correct weights and bias and see how we get the desired classification



$$w_1(0) = -30, w_2(0) = 300, b(0) = 50, \eta = 0.01$$

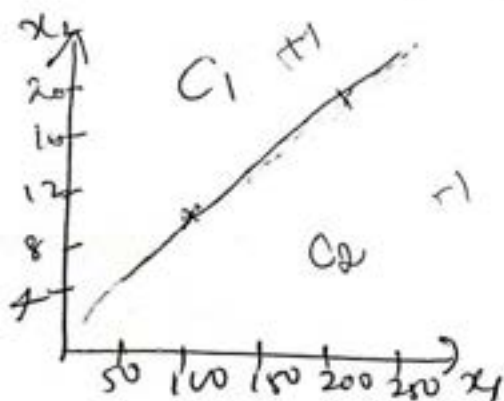
Decision boundary ~~data~~ for this example is

$$w_1 x_1 + w_2 x_2 + b = 0$$

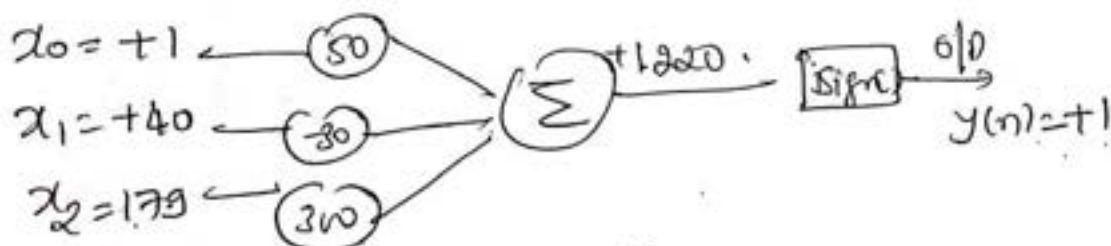
$$-30 x_1 + 300 x_2 + 50 = 0$$

$$x_1 = 100 \rightarrow x_2 = \frac{30 \times 100 - 50}{300} = 9.83$$

$$x_1 = 200 \rightarrow x_2 = \frac{30 \times 200 - 50}{300} = 19.83$$



This decision boundary clearly separates two classes. So the weights are fixed and decision boundary is obtained. Let us give test input (new unknown fruit) & determine the class.



$$x(n) = [+1, +40, +17.9]$$

$$w(3) = [50, -30, 300]$$

$$d(n) = +1$$

$$y(n) = \text{Sgn}(w^T(3) \cdot x(n)) = \text{Sgn}(50 \times 1 - 30 \times 40 + 300 \times 17.9)$$

$$\text{Sgn}(1220) = +1 \rightarrow \text{This unknown fruit belongs to class } C_1$$

Example 2

Initially let us assume an random weight and bias and train the neural network with the training example.

(8)

Train with known fruit

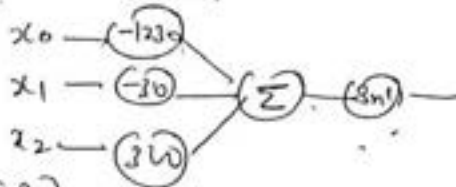
$$x(n) = [+1, 114, 15.2]^T \text{ and } d(n) = +1 \rightarrow \text{class 1}$$

$$w(n) = [-1230, -30, 300]$$

$$y(n) = \text{sgn}(w(n) \cdot x(n))$$

$$= \text{sgn}(-1230 \times 1 - 30 \times 114 + 300 \times 15.2)$$

$$= \text{sgn}(-90) = -1 \neq d(n) \rightarrow \text{class 2}$$



The target class is class 1 but the output class is class 2. So new weight should be recalculated.

$$w(n) = [-1230, -30, 300], x(n) = [+1, 114, 15.2]$$

$$d(n) = +1, y(n) = -1, \eta = 0.01$$

$$w(n+1) = w(n) + \eta [d(n) - y(n)] x(n)$$

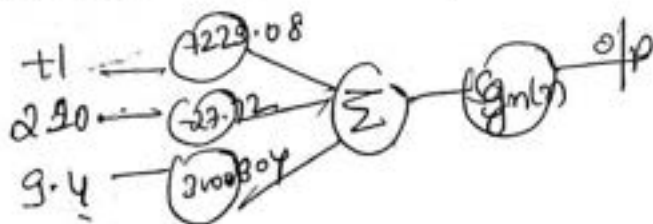
In general

$$w(n+1) = w(n) + \eta [d(n) - y(n)] x(n)$$

$$w(2) = [-1230, -30, 300]^T + [0.02, 2.28, 0.304]^T$$

$$w(n) = [-1229.08, -27.72, 300.304]^T$$

Train with known fruit



$$w(2) = [-1229.08, -27.72, 300.304]$$

$$x(2) = [+1, 210, 9.4] \text{ and } d(2) = -1 \rightarrow \text{class 2}$$

$$y(2) = \text{sgn}(w(2) \cdot x(2)) =$$

$$= \text{sgn}(-1229.08 \times 1 - 27.72 \times 210 + 300.304 \times 9.4)$$

$$= \text{sgn}(-4227.424) = -1 = d(2)$$

Hence no need to recalculate the weights.

$$d(n) - y(n) = 0$$

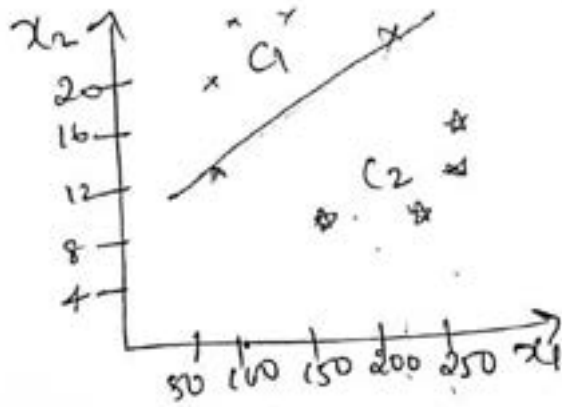
(9)

Therefore decision boundary after neural network training is.

$$-24.72x_1 + 300.304x_2 - 1229.08 = 0$$

$$x_1 = 100 \quad x_2 = \frac{24.72 \times 100 + 1229.08}{300.304} = 13.32$$

$$x_2 = 200 \quad x_2 = \frac{24.72 \times 200 + 1229.08}{300.304} = 22.55$$



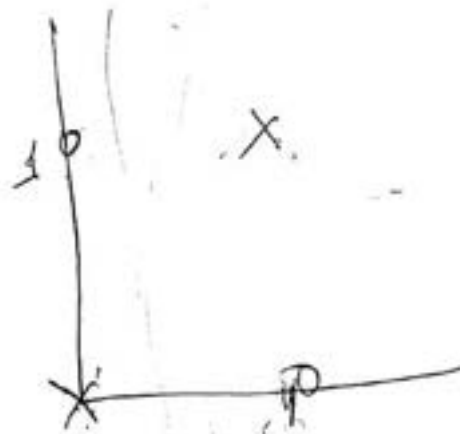
Once the decision boundary is built, New test examples can be applied to get to see which class it belongs.

Limitations of Perceptron (Single layer Perceptron)

- It can separate classes that lie either side of straight line easily.
- But it is reality division b/w classes are more complex. For example Classical Exclusive (XOR) Problem

XOR

x_1	x_2	o/p
0	0	0
0	1	1
1	0	1
1	1	0



XOR Problem is not linearly separable. We cannot draw a line that separates X class and circle class.

An example of a perception that implements AND and its geometric interpretation in two is given below

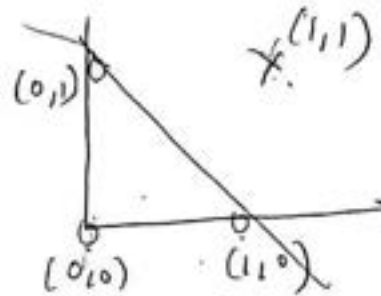
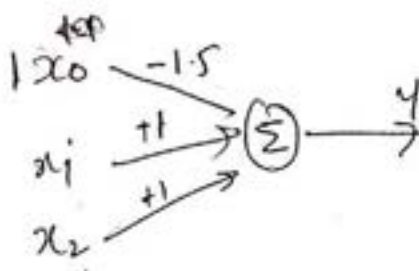
(10)

AND

$$x = [1, x_1, x_2] \text{ and } w = [-1.5, 1, 1]$$

x_1	x_2	o/p
0	0	0
0	1	0
1	0	0
1	1	1

3D



Discrete function
 $w_0x_0 + w_1x_1 + w_2x_2 + b = 0$

$$y = \text{sgn}(x_1 + x_2 - 1.5)$$

$$x_1 = 1, x_2 = 0$$

$$y = \text{sgn}(1 + 0 - 1.5)$$

$$y = \text{sgn}(-0.5)$$

$$y = 0$$

$$x_1 = 1, x_2 = 1$$

$$y = \text{sgn}(1 + 1 - 1.5)$$

$$y = \text{sgn}(0.5)$$

$$y = 1$$

Activation function

$$\begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

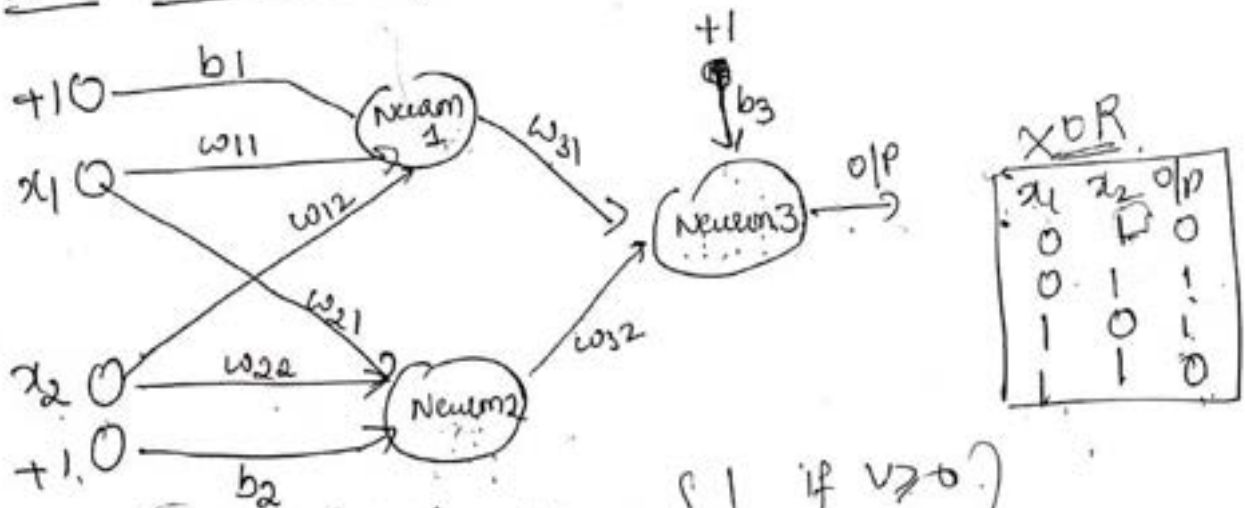
Multilayer Perceptron

(11)

A perceptron that has a single layer of weights can only approximate linear function of inputs and cannot solve problems like XOR, where the discriminant to be estimated is nonlinear. Similarly a perceptron cannot be used for nonlinear regression.

Multilayer perceptron can be used for classification which can implement nonlinear discriminants.

XOR in multilayer perceptron



$$\phi(v) = \text{Activation function} = \begin{cases} 1 & \text{if } v \geq 0 \\ 0 & \text{if } v < 0 \end{cases}$$

(i) $w_{11} = w_{12} = w_{21} = w_{22} = w_{31} = w_{32} = +1$
 $w_{31} = -2, b_1 = -1.5, b_2 = b_3 = -0.5$

For $[0, 1]$ input $w_H(n) = \begin{bmatrix} -1.5 & -0.5 \\ 1 & 1 \\ 1 & 1 \end{bmatrix}^T$ $x(n) = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$

$$y^H(n) = \phi[w^T(n) x(n)] = \phi \left[\begin{bmatrix} -1.5 & -0.5 \\ 1 & 1 \\ 1 & 1 \end{bmatrix}^T \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right]$$

$$= \phi \left[\begin{bmatrix} -1.5 & 1 & 1 \\ -0.5 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \right] = \phi \begin{bmatrix} -0.5 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} y_1^H \\ y_2^H \end{bmatrix}$$

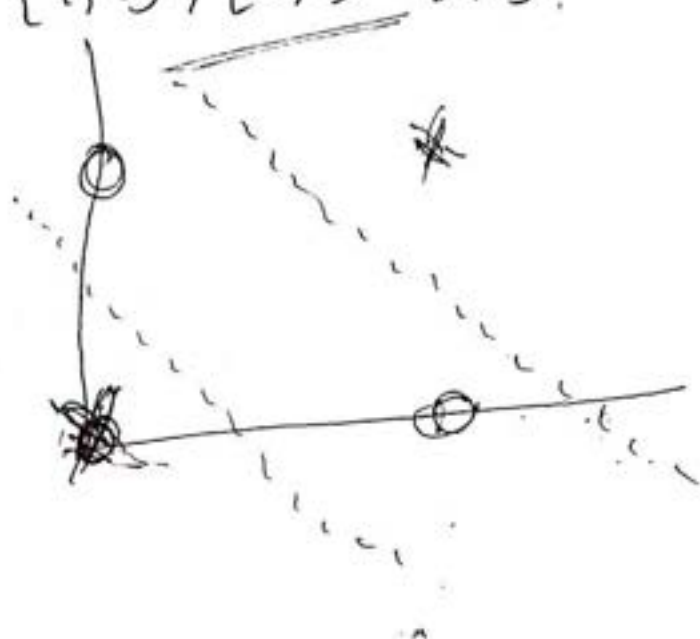
$$y^H(n) = \phi(w^T(n)x(n)) = \phi\left[\begin{bmatrix} b_3 \\ w_{31} \\ w_{32} \end{bmatrix}^T \begin{bmatrix} 1 \\ y_{1H} \\ y_{2H} \end{bmatrix}\right]$$

See the activation function

$$= \phi\left[\begin{bmatrix} -0.5 & -2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}\right] = \phi[0.5] = \underline{\underline{\begin{pmatrix} 1 \\ 0 \end{pmatrix}}}$$

For input $[0, 1] \rightarrow$ we got $= 0$ as 1

Similarly, it is done for input $[1, 0], [0, 0] \& [1, 1]$.



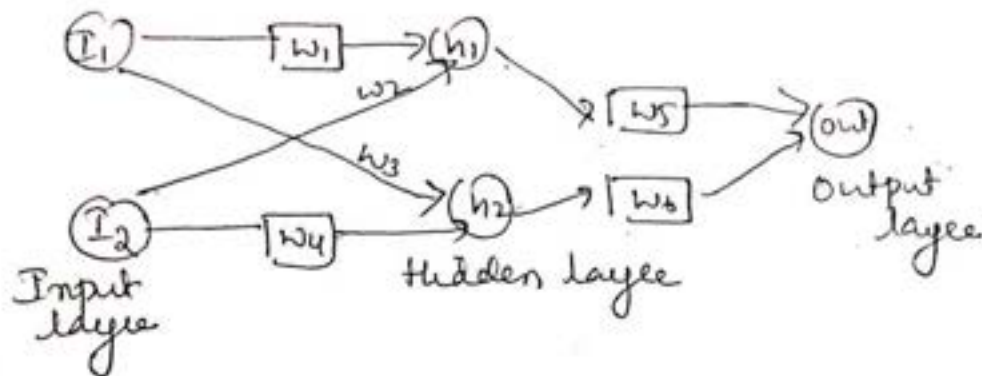
Backpropagation Algm

(13)

What is Backpropagation?

Backpropagation is a supervised learning algm, for training neural networks.

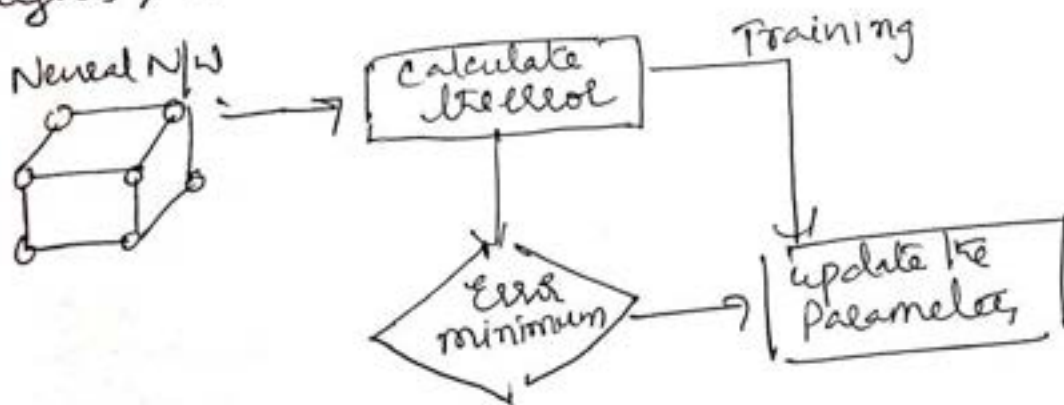
Why we need Backpropagation?

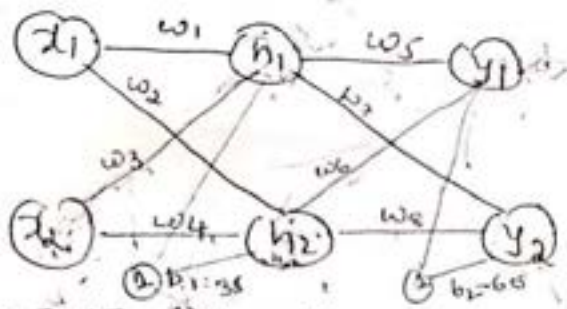


When we need to design neural network, we need to initialize the weights with random value. There is a least probability with whatever weights assigned will lead to correct output.

If $error = 0$, if Prediction output = actual o/p .

If error is high we need to change the parameters (weights) so that error becomes minimum.





$$x_1 = 0.05, x_2 = 0.10$$

$$w_1 = 0.15, w_2 = 0.20, w_3 = 0.25, w_4 = 0.30, b_1 = 0.3$$

$$w_5 = 0.4, w_6 = 0.45, w_7 = 0.5, w_8 = 0.55, b_2 = 0.6$$

$$y_1 = 0.01, y_2 = 0.99, \eta = 0.5$$

Forward pass :- Using the initialized random weights, we first compute the predicted output. Then find the error between predicted output and actual output. To reduce the error we update the weights from Backward using gradient descent, so that error is minimum.

Again we go forward pass, calculate the difference between predicted and actual output. To minimize the error update the weights through backward pass. This is repeated until the error is minimum.

Note : A complete pass over all the patterns in the training set is called an epoch.

$$net h_1 = w_1 x_1 + w_3 x_2 + b_2 x_1$$

$$net h_1 = 0.15 \times 0.05 + 0.2 \times 0.1 + .35 \times 1 = .0.3775$$

Apply the activation function = $\frac{1}{1+e^{-x}}$

$$(a) \quad out h_1 = \frac{1}{1+e^{-x}} = \frac{1}{1+e^{-0.3775}} = 0.593269$$

Similarly.

$$out h_2 = 0.59688$$

Calculating y_1

$$y_1 = out h_1 \times w_5 + out h_2 \times w_6 + b_2$$

$$y_1 = 0.5932 \times .4 + 0.45 \times .5968 + 0.6$$

$$y_1 = 1.1059059$$

$$(b) \quad out y_1 = \frac{1}{1+e^{-y_1}} = \frac{1}{1+e^{-1.1059059}} = 0.751365$$

In the same way,

$$out y_2 = 0.71292$$

Calculating total error

$$E_{total} = \sum \frac{1}{2} (Prediction - actual)^2$$

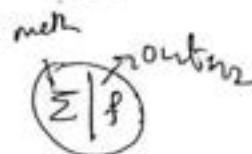
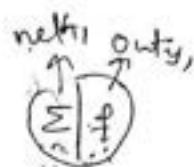
$$= \frac{1}{2} (out y_1 - T_1)^2 + \frac{1}{2} (out y_2 - T_2)^2$$

$$= \frac{1}{2} (.07513 - 0.01)^2 + \frac{1}{2} (.075136 - .99)^2$$

$$E_{total} = .298$$

11.1

(15)

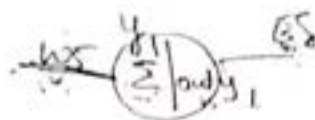


To reduce the error we will backpropagate and update the weights.

(16)

Consider w_5 .

$$\text{Error at } w_5 = \frac{\delta E_{\text{total}}}{\delta w_5}$$



$$\frac{\delta E_{\text{total}}}{\delta w_5} = \frac{\delta E_{\text{total}}}{\delta \text{out}_{y1}} \times \frac{\delta \text{out}_{y1}}{\delta y_1} \times \frac{\delta y_1}{\delta w_5}$$

$$\frac{\delta E_{\text{total}}}{\delta \text{out}_{y1}} = \frac{1}{2} (T_1 - \text{out}_{y1})^2 + \frac{1}{2} (T_2 - \text{out}_{y2})^2$$

$$\frac{\delta E_{\text{total}}}{\delta \text{out}_{y1}} = 2 \times \frac{1}{2} (T_1 - \text{out}_{y1}) \times -1 + 0$$

$$= -(T_1 - \text{out}_{y1})$$

$$= -(0.01 - 0.75136) = 0.74136$$

$$\begin{aligned} \frac{\delta \text{out}_{y1}}{\delta y_1} &= \text{out}_{y1} (1 - \text{out}_{y1}) \\ &= 0.75136 (1 - 0.75136) \\ &= 0.186815 \end{aligned}$$

$$\text{Note: } \text{out}_{y1} = \frac{1}{1 + e^{-y_1}}$$

$$\frac{\delta y_1}{\delta w_5} = \frac{1 \times \text{out}_{y1} \times w_5}{\text{out}_{y1}}$$

$$[y_1 = \text{out}_{y1} \times w_5 + \text{out}_{y2} \times w_6 + b_1]$$

$$\frac{\delta y_1}{\delta w_5} = 0.59326$$

$$\begin{aligned} \frac{\delta E_{\text{total}}}{\delta w_5} &= \frac{\delta E_{\text{total}}}{\delta \text{out}_{y1}} \times \frac{\delta \text{out}_{y1}}{\delta y_1} \times \frac{\delta y_1}{\delta w_5} \\ &= 0.74136 \times 0.1868 \times 0.59326 \\ &= 0.082167 \end{aligned}$$

Updating w_5

$$w_5 = w_5 - \eta \times \frac{\delta E_{\text{total}}}{\delta w_5}$$

$$w_5 = 0.4 - 0.5 \times 0.082167$$

$$= 0.258916$$