

1. Consider the following multi-layer perceptron network. Assume that the neurons have a sigmoid activation function, perform a forward pass and a backward pass on the network. Assume that the actual output of y is 1 and learning rate is 0.9. Perform another forward pass.

→ Activation function : sigmoid

$$\text{Actual Output } y = 1$$

$$\text{Learning Rate } \eta = 0.9$$

→ Net input to H_4 ,

$$H_4 = 1(0.2) + 0(0.4) + 1(-0.5) + (-0.4)$$

$$= -0.7$$

$$H_4 = \frac{1}{1 + e^{-(0.7)}} = 0.3318122$$

→ Net input to H_5

$$H_5 = 1(-0.3) + 0(0.1) + 1(0.2) + 0.2 \\ = 0.1$$

$$H_5 = \frac{1}{1 + e^{(0.1)}} = 0.5249792$$

→ Net input to O_6

$$O_6 = 0.3318122(-0.3) + 0.5249792(-0.2) + 0.1 \\ = -0.1045395$$

$$O_6 = \frac{1}{1 + e^{(-0.1045395)}} = 0.473889$$

$$\delta_6 = O_6(1 - O_6)(1 - O_6) = 0.473889 = 0.1311690$$

$$\Delta w_{46} = \eta \cdot \delta_6 \cdot H_4 = 0.9 \times 0.474 \times 0.332 = 0.1416$$

$$\Delta w_{56} = \eta \cdot \delta_6 \cdot H_5 = 0.9 \times 0.474 \times 0.525 = 0.223965$$

$$\Delta \theta_6 = \eta \cdot \delta_8 = 0.9 \times 0.1312 = 0.1181$$

$$\Delta w_{46} = \eta \cdot \delta_8 \cdot H_4 = 0.9 \times 0.1312 \times 0.332 = 0.0392$$

$$\Delta w_{56} = \eta \cdot \delta_8 \cdot H_5 = 0.9 \times 0.1312 \times 0.522 = 0.061992$$

$$\begin{aligned}\delta_4 &= H_4(1 - H_4)(\delta_6 \cdot w_{46}) \\ &= 0.3318(1 - 0.3318)(0.1312 \cdot 0.0392) \\ &= 0.00114\end{aligned}$$

$$\begin{aligned}\delta_5 &= H_5(1 - H_5)(\delta_6 \cdot w_{56}) \\ &= 0.525(1 - 0.525)(0.1312 \times 0.062) \\ &= 0.002025\end{aligned}$$

$$w_{46} = -0.3 + 0.0392 = -0.2608$$

$$w_{56} = -0.2 + 0.061992 = -0.138$$

$$\theta_6 = 0.1 + 0.1181 = 0.2181$$

$$\Delta w_{14} = \eta \cdot \delta_4 \cdot \alpha_1 = 0.9(0.00114)(1) = 0.001026$$

$$\Delta w_{24} = \eta \cdot \delta_4 \cdot \alpha_2 = 0$$

$$\Delta w_{34} = \eta \cdot \delta_4 \cdot \alpha_3 = 0.9(0.00114)(1) = 0.001026$$

$$\Delta \theta_5 = \eta \cdot \delta_5 \cdot \alpha_1 = 0.9 \times 0.002025 \times 1 = 0.0018225$$

$$\Delta w_{25} = \eta \cdot \delta_5 \cdot \alpha_2 = 0$$

$$\Delta w_{35} = \eta \cdot \delta_5 \cdot \alpha_3 = 0.9 \times 0.002025 = 0.0018225$$

$$w_{154} = 0.2 + 0.001026 = 0.201026$$

$$w_{24} = 0.4 + 0 = 0.4$$

$$w_{34} = -0.5 + 0.001026 = -0.498974$$

$$w_{25}^{\frac{15}{15}} = -0.3 + 0.0018225 = -0.2981775$$

$$w_{25} = 0.1$$

$$w_{35} = 0.2 + 0.0018225 = 0.2018225$$

$$H_4 = \alpha_1 \times w_{14} + \alpha_2 \times w_{24} + \alpha_3 \times w_{34} = \\ = 1(0.201026) + 0(0.4) + 1(-0.498934) \\ = -0.297948$$

$$H_4 = \frac{1}{1 + e^{-(0.297948)}} = 0.4260592$$

$$H_5 = \alpha_1 \times w_{15} + \alpha_2 \times w_{25} + \alpha_3 \times w_{35} \\ = 1(-0.2981775) + 0 + 1(0.2018225) \\ = -0.096355$$

$$H_5 = \frac{1}{1 + e^{-(0.096355)}} = 0.4759299$$

$$O_6 = H_4 \times w_{46} + H_5 \times w_{56} + O_6 \\ = 0.4260592 \times (-0.2608) + 0.4759299 \times (-0.138) + 0.22 \\ = 0.04320$$

$$O_6 = \frac{1}{1 + e^{-(0.0432)}} = 0.5107997$$

2. Apply K-nearest neighbour classifier on below mentioned fitness dataset and classify the test record based on the values of input feature.

Row	Height	weight	Class	Distance
1	167	51	underweight	$\sqrt{(170-167)^2 + (57-51)^2} = 6.71$
2	182	62	Normal	$\sqrt{(144+25)} = 13.0$
3	176	69	Normal	$\sqrt{(36+144)} = 13.42$
4	173	64	Normal	$\sqrt{9+49} = 7.62$
5	172	65	Normal	$\sqrt{4+64} = 8.25$
6	174	56	Underweight	$\sqrt{16+1} = 4.12$
7	169	58	Normal	$\sqrt{1+1} = 1.41$
8	173	67	Normal	$\sqrt{9+0} = 3$
9	170	55	Normal	$\sqrt{0+4} = 2$

K = 3

- Pick three nearest neighbour
- Row 7 → Distance = 1.41 → Normal
- Row 8 → Distance = 2.00 → Normal
- Row 9 → Distance = 3.00 → Normal
- ∴ (770, 57) → Normal

3. Define support vector machine. Discuss its advantages and application in brief

→ A support vector machine (SVM) is supervised machine learning algorithm used for classification and regression task. The main idea behind SVM is to find the optional hyperplane that best separates the data points of different classes in a high-dimensional space.

⇒ Advantages of SVM:

→ Effective in high dimensional space.

→ Memory efficient

→ Versatile

→ Robust to overfitting

⇒ Applications of SVM:

→ Text and hypertext classification

→ Image classification

→ Bioinformatics

→ Financial applications.

→ Medical diagnosis

4 Explain SVM classifier algorithm for linearly separable data. Use diagram to support your discussion.

⇒ steps of algorithm :

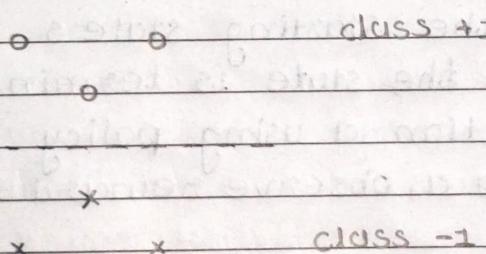
- Input label training data (x_i, y_i) where $y_i \in \{-1, 1\}$
- Find the hyperplane $w \cdot x + b = 0$ that separates the classes.
- Maximise the margin between the support vectors.
- Use the optimal w and b to classify new data point

→ Decision Rule :

If $w \cdot x + b \geq 0 \Rightarrow$ class +1

Else \Rightarrow class -1

⇒ Diagram for SVM with linearly separable data:



5 Discuss the working of Q-learning algorithm in reinforcement learning.

→ Q-learning is a model-free reinforcement learning algorithm that helps an agent learn the optimal action-selection policy by learning the Q-values of state-action pairs. It does not require a model of the environment.

\Rightarrow Q-Learning Intuition :

- \rightarrow The agent interacts with the environment in discrete time steps.
- \rightarrow At each time step :
 - \rightarrow The agent observes the current state (s).
 - \rightarrow Choose an action (a) based on some policy
 - \rightarrow Receives a reward (r) and observes the next state (s').
- \rightarrow Updates the Q-values of the pair using the Q-Learning formula.

\Rightarrow Formula :

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

\Rightarrow Steps of the Q-Learning Algorithm :

1. Initialize Q-table with zeros for all state-action pairs.
2. For each episode :
 1. Initialize the starting state s
 2. Repeat until the state is terminal :
 - \rightarrow choose action a using policy
 - \rightarrow Take action a , observe reward r and new state s'
 - \rightarrow Update $Q(s, a)$ using the update rule.
 - \rightarrow set $s = s'$
3. Discuss the working .
4. State and explain the inductive Learning hypothesis.
- \rightarrow Any hypothesis found to approximate the target function well over a sufficiently large training

dataset will also perform well over unseen data from the same distribution.

⇒ Working :

- You collect a training dataset
- You train a model using dataset.
- You evaluate the model's performance on new data.
- The inductive learning hypothesis allows you to expect that good training performance indicate good performance on similar, unseen data.

⇒ Example :

- Train a spam classifier on 10,000 emails labeled as "spam" or "not spam".
- Formal signature significance.
- Evaluating model performance.
- Avoiding overfitting.
- Developing generalization strategies.

7. what is Occam's Razor with respect to hypothesis search?

- Among competing hypothesis that explain the observed data equally well, the simplest one should be preferred.
- When multiple hypotheses fit the training data well, Occam's Razor suggests that the simplest hypothesis is more likely to generalize better to unseen data.

⇒ Benefits :

- Improve generalization
- Reduce overfitting.
- Encourages interpretability and efficiency.

8. Explain the following terms: (i) restriction bias
(ii) performance bias. What is the preference bias in decision tree?

- ⇒ (i) Restriction Bias (also called language bias)
 - Restriction bias refers to the limitations or constraints on the set of hypotheses that a learning algorithm can consider.
- ⇒ (ii) Preference Bias (also called search bias).
 - Preference bias is a soft bias where the algorithm has a preference for some hypotheses over others, even if multiple hypotheses fit the data.

⇒ Preference Bias in decision trees:

1. Information gain or gini index:
 - The algorithm prefers attributes that split the data in a way that maximize information gain.
2. Tree simplicity:
 - Some implementations may prefer shallower trees or fewer splits to avoid overfitting.
 - Pruning is used to remove parts of the tree that do not provide much predictive power, reflecting a bias toward simpler trees.

9. What is overfitting off an ML model? How can you detect whether overfitting has occurred in a decision tree?

- Overfitting occurs when a machine learning model learns the training data too well,

including noise and irrelevant patterns, resulting in poor generalization to new, unseen data.

⇒ How to detect overfitting in decision Trees:

1. Performance crap:

→ Train Accuracy > Test Accuracy.

→ Large difference between training and validation/test accuracy.

2. Complex Tree structure:

→ The tree has too many nodes or depth.

→ Leaf nodes with very few samples.

3. cross-validation Results:

→ Significant variance in accuracy across fold may indicate overfitting

4. visual Inspection:

→ The tree is too specific to the training data

10 Explain the advantage and disadvantage of using the following attribute selection criteria for decision trees: (i) information gain (ii) Gini Ratio.

⇒ (i) Information gain:

→ Information gain measures the reduction in entropy achieved by splitting the data based on a particular attribute.

Information Gain = Entropy (parent)

$$= \sum \left[\frac{|\text{Subset}_i|}{|\text{Total}|} \times \text{Entropy}(\text{subset}_i) \right]$$

⇒ Advantage:

1. simple and effective.

→ Easy to compute and interpret.

→ works well with categorical attribute

is used in ID3 Algorithm:

- A core component of the ID3 decision tree algorithm.
- Quick identifies attributes that provide the most immediate "purity" improvement.
- ⇒ Disadvantage:
 1. Bias toward many-valued Attributes
 2. overfitting on training data.

Ciui Cratio Ratio:

- Cratio Ratio is normalized version of information gain that penalize attributes with many values.

$$\text{Cratio Ratio} = \frac{\text{Information Gain}}{\text{split Information}}$$

- split information measures how broadly and uniformly the data is split.

⇒ Advantages:

- Reduce Bias
- Better generalization
- Used in C4.5 Algorithm.

⇒ Disadvantages:

- can be unstable
- computationally more expensive.

1) consider the following set of training example

Instance	Classification	c_1	c_2
1	+	T	T
2	+	T	T
3	-	T	F
4	+	F	F
5	-	F	T
6	-	F	T

a. what is the entropy of this collection of training samples with respect to the target function classification ?

→ Positive (+) : 3

Negative (-) : 3

$$\text{Entropy}(S) = -P_+ \log_2(P_+) - P_- \log_2(P_-)$$

$$P_+ = \frac{3}{6} = 0.5 \quad P_- = \frac{3}{6} = 0.5$$

$$\text{Entropy}(S) = -0.5 \log_2(0.5) - 0.5 \log_2(0.5) = 1$$

(b) what is the information gain of c_2 relative to these training example.

→ split data based on c_2

$c_2 = T$: Instance 1,2,5,6 → classes : +, +, -, -

$c_2 = F$: Instance 3,4 → classes : -, +

→ Entropy for $c_2 = T$

$$\text{Entropy}_T = -0.5 \log_2(0.5) - 0.5 \log_2(0.5) = 1$$

→ Entropy for $c_2 = F$

$$\text{Entropy}_F = -0.5 \log_2(0.5) - 0.5 \log_2(0.5) = 1$$

→ weighted Entropy after split

$$\text{Entropy}_{c_2} = \frac{4}{6} \cdot 1 + \frac{2}{6} \cdot 1 = \frac{4}{6} + \frac{2}{6} = 1$$

→ Information gain
 $\text{gain}(u_2) = \text{Entropy}(s) - \text{Entropy}_{u_2}$

$$= I - I$$

$$= 0$$

→ Information gain of $u_2 = 0$