

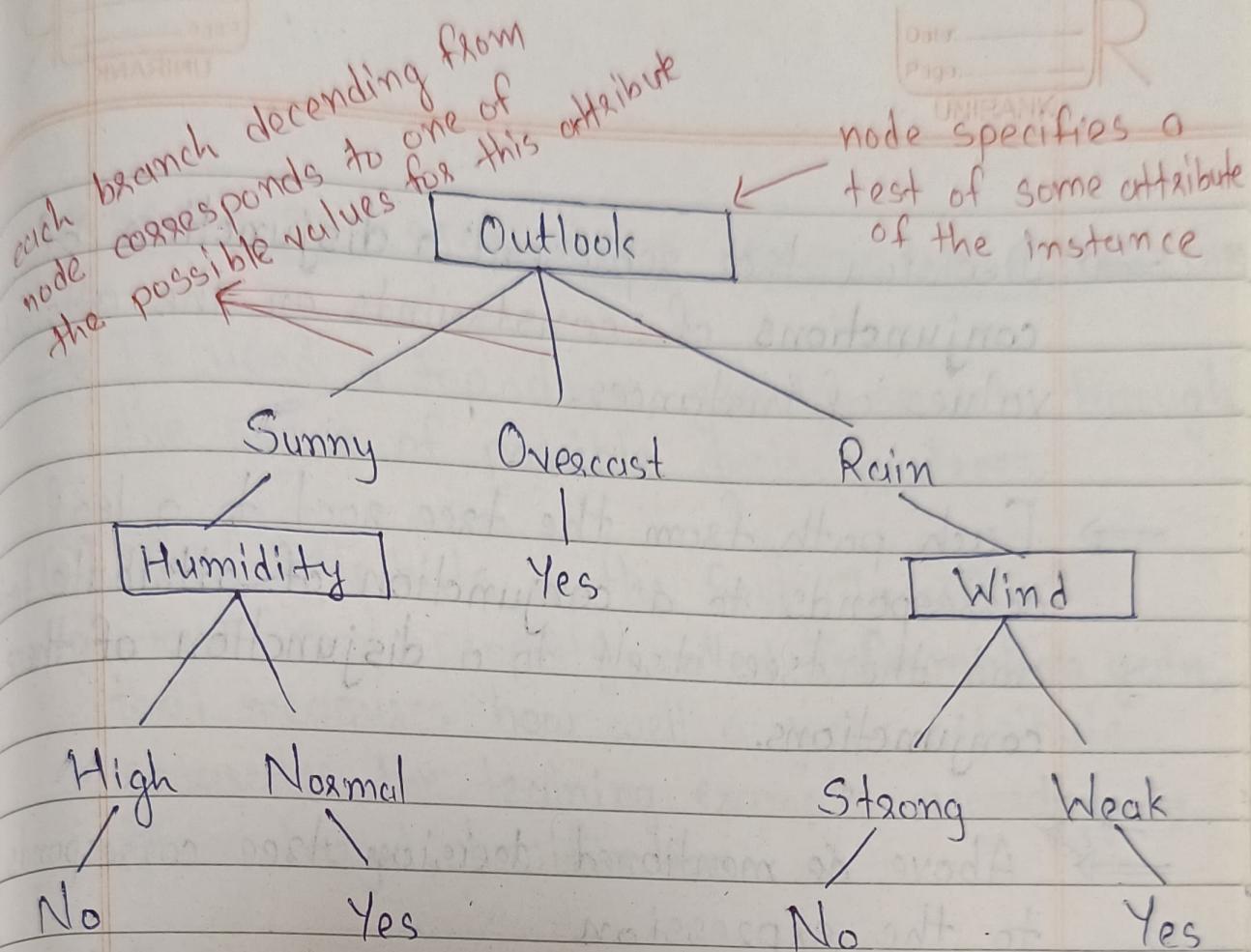
Advantages of Inductive bias :-

- ⇒ Inductive bias provides a nonprocedural means of characterizing their policy for generalizing beyond the observed data.
- ⇒ It allows comparison of different learners according to the strength of the inductive bias they employ.

Module 2: Supervised Learning

★ Decision Tree Learning:

- ⇒ Decision Tree learning is one of the most widely used and practical methods for - inductive inference.
- ⇒ It is a method for approximating discrete valued functions that is robust to noisy data and capable of learning disjunctive expressions.
- ⇒ Decision Tree Learning is a method for - approximating discrete-valued target functions, in which the learned function is represented by a decision tree.



⇒ Above tree classifies Saturday mornings - according to whether or not they are suitable for playing tennis.

~~Ex/~~ Instance = <Outlook = Sunny, Temperature = Hot,
Humidity = High, Wind = Strong>

⇒ above instance would be sorted down the leftmost branch of this decision tree and would therefore be classified as a negative instance.

- Decision tree represent a disjunction of conjunctions of constraints on the attribute values of instances.
- Each path from the tree root to a leaf corresponds to a conjunction of attribute test and the tree itself to a disjunction of these conjunctions.
- Above mentioned decision tree corresponds to the expression

$$(\text{Outlook} = \text{Sunny} \wedge \text{Humidity} = \text{Normal})$$
$$\vee (\text{Outlook} = \text{Overcast})$$
$$\vee$$
$$(\text{Outlook} = \text{Rain} \wedge \text{Wind} = \text{Weak})$$

* When to use Decision Tree Learning?

1. Instances are represented by attribute-value pairs.
2. The target function has discrete output values.
3. Disjunctive descriptions may be required.
4. The training data may contain errors.
5. The training data may contain missing attribute values.

* BASIC DECISION TREE LEARNING ALGORITHM

→ It used a top-down, greedy search through the space of possible decision trees.

(a) Which attribute is the Best Classifier?

→ Statistical property called Information gain, that measures how well a given attribute separates the training examples according to their target classification.

(i) Entropy measures homogeneity of Examples:

<sup>boolean
classification</sup> \rightarrow Entropy(S) = $-P_+ \log_2 P_+ - P_- \log_2 P_-$

<sup>multi-class
classification</sup> \rightarrow Entropy(S) = $\sum_{i=1}^c -P_i \log_2 P_i$

S = collection of examples containing positive and negative examples of some target concept.

~~Ex~~ S = collection of 14 examples of some boolean concept, including 9 positive and 5 negative examples.

Note:-

Entropy is considered as a measure of the impurity in a collection of training examples.

Note: \rightarrow Entropy is 0 if all members of S belong to the same class.
 \rightarrow Entropy is 1 if collection contains an equal number of positive and negative examples.

$$\text{Entropy}(S) = -(9/14)\log_2(9/14) - (5/14)\log_2(5/14)$$
$$= 0.940$$

(ii) Information gain measures the expected reduction in Entropy:-

\Rightarrow Information gain for this is simply the expected reduction in entropy caused by partitioning the examples according to this attribute.

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

where;

S = collection of examples

A = attribute

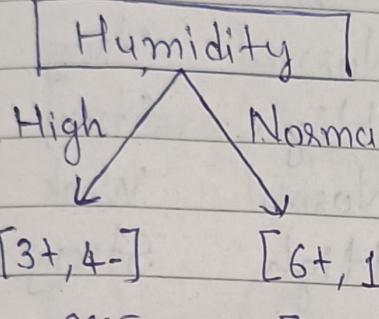
$\text{Values}(A)$ = set of all possible values for attribute A

S_v = subset of S for which attribute A has value v . (i.e. $S_v = \{s \in S | A(s)=v\}$)

Which attribute is the best classifier?

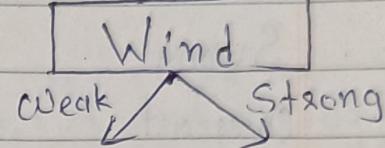
$$S: [9+, 5-]$$

$$E = 0.940$$



$$S: [9+, 5-]$$

$$E = 0.940$$



$$E = 0.985$$

$$E = 0.592$$

$$\text{Gain}(S, \text{Wind}) = 0.940 - \left(\frac{7}{14}\right)0.985 - \left(\frac{6}{14}\right)0.592$$

$$= 0.048$$

$$= 0.151$$

\Rightarrow Humidity provides greater information gain than Wind, relative to the target classification.

Example to illustrate the operation of ID3 :-

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
01	Sunny	Hot	High	Weak	No
02	Sunny	Hot	High	Strong	No
03	Overcast	Hot	High	Weak	Yes
04	Rain	Mild	High	Weak	Yes
05	Rain	Cool	Normal	Weak	Yes
06	Rain	Cool	Normal	Strong	No
07	Overcast	Cool	Normal	Strong	Yes

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Sol:-

Step 1:- Select attribute which will be the topmost node of the decision tree.

→ ID3 determines the information gain for each candidate attribute, then selects the one with highest information gain.

$$\text{Gain}(S, \text{Outlook}) = 0.246$$

$$\text{Gain}(S, \text{Humidity}) = 0.151$$

$$\text{Gain}(S, \text{Wind}) = 0.048$$

$$\text{Gain}(S, \text{Temperature}) = 0.029$$

$$\text{Gain}(S, \text{Outlook}) = \text{Entropy}(S) - \sum_{v \in \{\text{Sunny, Overcast, Rainy}\}} \frac{|\text{S}_v|}{|S|} \text{Entropy}(\text{S}_v)$$

$Gain(S, \text{Outlook})$

$$\begin{aligned} \text{Entropy}(S) &= -\left(\frac{9}{14}\right) \log_2 \left(\frac{9}{14}\right) - \left(\frac{5}{14}\right) \log_2 \left(\frac{5}{14}\right) \\ &\approx 0.940 \end{aligned}$$

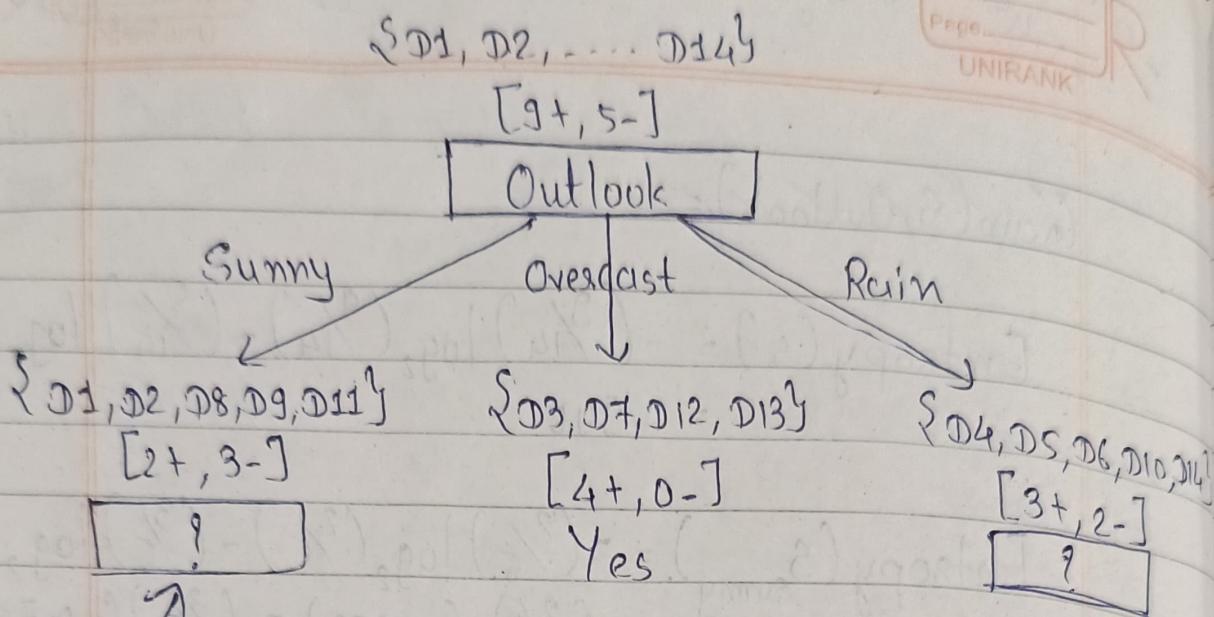
$$\begin{aligned} \text{Entropy}(S_{\text{sunny}}) &= -\left(\frac{2}{5}\right) \log_2 \left(\frac{2}{5}\right) - \left(\frac{3}{5}\right) \log_2 \left(\frac{3}{5}\right) \\ &= -\left(\frac{2}{5}\right)(-1.322) - \left(\frac{3}{5}\right)(-0.737) \\ &= 0.5288 + 0.4422 = 0.971 \end{aligned}$$

$$\text{Entropy}(S_{\text{overcast}}) = 0$$

$$\begin{aligned} \text{Entropy}(S_{\text{Rainy}}) &= -\left(\frac{3}{5}\right) \log_2 \left(\frac{3}{5}\right) - \left(\frac{2}{5}\right) \log_2 \left(\frac{2}{5}\right) \\ &= 0.971 \end{aligned}$$

$$\begin{aligned} \text{Gain}(S, \text{Outlook}) &= 0.940 - \left[\frac{\frac{5}{14} \cdot 0.971 + \frac{5}{14} \cdot 0.971}{14} \right] \\ &= 0.940 - [0.3468 + 0.3468] \\ &= 0.2464 \end{aligned}$$

⇒ According to the information gain measure, the Outlook attribute provides the best prediction of the target attribute, PlayTennis, over the training examples. Therefore, Outlook selected as the root node for decision tree and branches are created below the root for each of its possible value.



Step 2:

Which attribute should be tested here?

$$S_{\text{sunny}} = \{D1, D2, D8, D9, D11\}$$

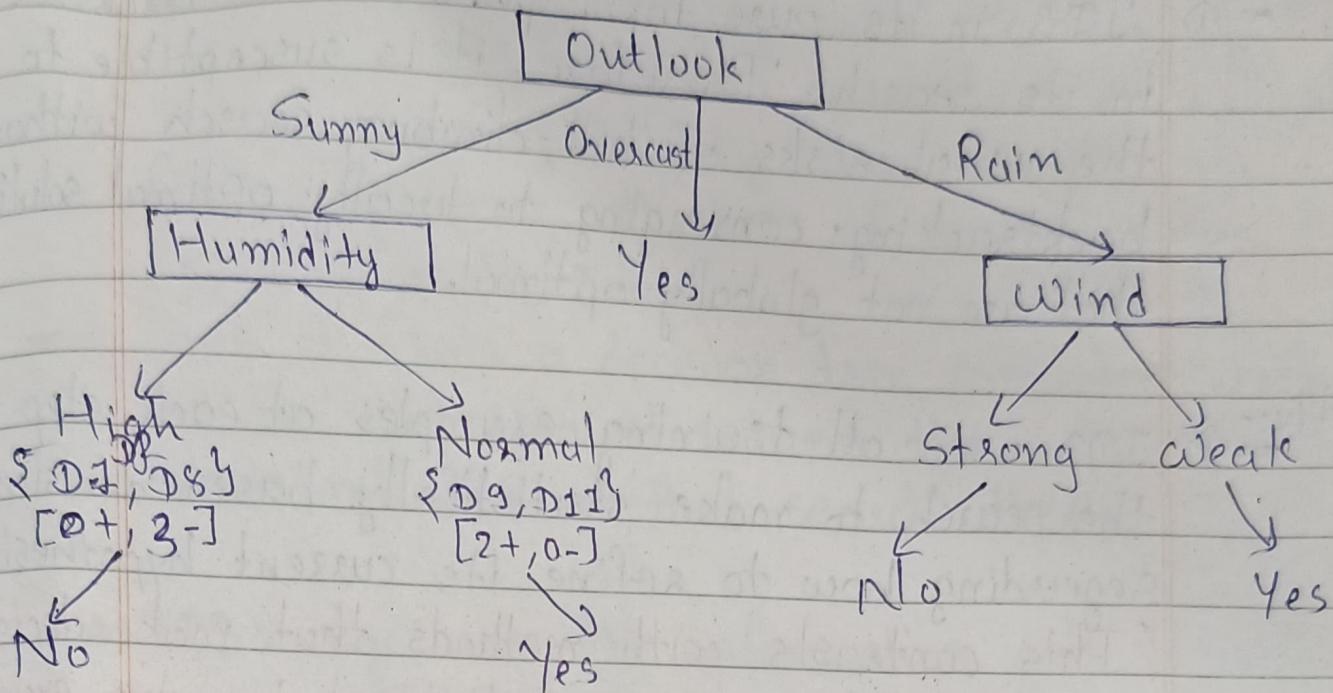
$$G_{\text{rain}}(S_{\text{sunny}}, \text{Humidity}) = 0.970 - \left(\frac{3}{5}\right)0 - \left(\frac{2}{5}\right)0 = 0.970$$

$$G_{\text{rain}}(S_{\text{sunny}}, \text{Temperature}) = 0.970 - \left(\frac{2}{5}\right)0 - \left(\frac{2}{5}\right)1 - \left(\frac{1}{5}\right)0 = 0.570$$

$$G_{\text{rain}}(S_{\text{sunny}}, \text{Wind}) = 0.970 - \left(\frac{2}{5}\right)1 - \left(\frac{3}{5}\right)0.918 = 0.19$$

\Rightarrow Humidity attribute is selected at this step.

* Final constructed decision tree:-



* Capabilities and Limitations of ID3:-

* Capabilities:

→ ID3's hypothesis space of all decision trees is a complete space of finite discrete-valued functions, relative to the available attributes.

Limitations:

→ ID3 maintains only a single current hypothesis as it searches through the space of decision trees. So ID3 loses the capabilities that follow from explicitly representing all consistent hypotheses.

limitation:-

→ ID3 in its pure form performs no backtracking in its search. Therefore, it is susceptible to the usual risks of hill-climbing search without backtracking: converging to locally optimal solutions that are not globally optimal.

Adv:-

→ ID3 uses all training examples at each step in the search to make statistically based decisions regarding how to refine its current hypothesis.

(This contrasts with methods that make decision incrementally, based on individual training example

~~or FIND-S or CANDIDATE-ELIMINATION~~). One advantage of using statistical properties of all the examples (~~information gain~~) is that the resulting search is much less sensitive to ~~error~~ in individual training examples.

→ ID3 can be easily extended to handle noisy training data by modifying its termination criterion to accept hypotheses that imperfectly fit the training data.

* Inductive bias in Decision Tree learning (ID3):
 Shorter trees are preferred over larger trees. Can

• BFS-ID3: algo. begins with the empty tree and searches breadth first through progressively more complex trees, first considering all trees of depth 1, then all trees of depth 2 etc.

→ Once it finds a decision tree consistent with the training data, it returns the smallest consistent tree at that search depth.

★ ID3: algo. uses the information gain heuristic and a hill climbing strategy, it exhibits a more complex bias than BFS-ID3.

★ A closer approximation to the inductive bias of ID3

- Shorter trees are preferred over longer trees.
Trees that place high information gain attributes close to the root are preferred over those that do not.

Difference between the hypothesis space search in these following two approaches

CANDIDATE-ELIMINATION

- The version-space of this algo. searches an incomplete hypothesis space. It searches this space completely finding every hypothesis fitting data.

ID3 searches a complete hypothesis space but it searches incompletely through this space (from simple to complex)