

mood-book



Introduction to Machine Learning

Def:-

- * Machine Learning is a growing technology which enables computers to learn automatically from past data. Machine Learning uses various algorithms for building mathematical models and making predictions using historical data (or) information.

Applications of Machine Learning:-

- * Image recognition
- * Speech recognition
- * Email filtering (Spam detector)
- * Virtual personal Assistants
- * Online customer Support
- * Online fraud detection.
- * Stock market trading
- * Traffic prediction.
- * Self driving cars.
- * Medical diagnostics etc.

What is machine Learning?

- * The term machine learning was first introduced by **Arthur Samuel** in 1959.
- * It enables a machine to automatically learn from data improve performance from experience, and predict

things without being explicitly programmed.

* Training data :- past / historical data which is used to train the program (or) algorithm.

Features of ML :-

- * ML uses data to detect various patterns in a given dataset.
- * It can learn from past data and improve automatically.
- * It's a data-driven technology.
- * Machine learning is much similar to data mining as it also deals with the huge amount of the data.

Well posed Learning problem :-

An agent solves a problem (or) task T , performance P and gain some experience E .

If P is measured at T it can improve E ..

----- Learning by Experience -----

Examples :-

1. playing checkers problem
2. Handwritten recognition problem
3. Robo driving learning problem.

problem	Task (T)	performance (P)	Experience (E)
1. playing checker learning	playing against opponents to win game	Make perfect moves to win game	It plays itself to improve
2. Handwriting recognition learning	classifying the images & text	Better classification	A database of homework text
3. Robot driving learning problem	Drive the car in a 4 lane highway.	Source to dest. the avg distance travelled (long & safe)	Images, vehicles on Road.

Designing a learning System:-

- * To get a Successful learning System, it should be designed
- * For a proper design, several steps should be followed.
- * These steps can create perfect & efficient System.

5 * Important steps :-

1. choosing the training experience.
2. choosing the target function
3. choosing a representation for target function

4. Choosing a learning algorithm for approximating the target function.

5. Final design:

* By using these 4 steps we will get a final design.

* 1st choosing The training Experience :-

Below are the attributes which will impact on success and failure of data.

* 1st attribute :-

The training Experience will be able to provide "direct (or) indirect feedback" regarding choices.

Ex:- while playing chess the training data will provide feedback to itself like instead of this move if this is chosen the chances of success increases.

* 2nd attribute :-

It is "the degree" to which the learner will control the sequence of training examples.

Ex:- when training data is fed to the machine then at that time accuracy is very less but when it gains experience while playing again and again with itself (or) opponent the machine algorithm will get feedback and control the chess game accordingly.

* 3rd attribute :-

It is how it will represent "the distribution of examples" over which performance will be measured.

2. Choosing the target function!:-

What type of knowledge is learnt and how it is used by the performance system.

Ex:- Checkers game

* While moving diagonally set of all possible moves is called legal moves.

legal moves



one move



target move

* Travel only in forward direction

* Only one move per chance

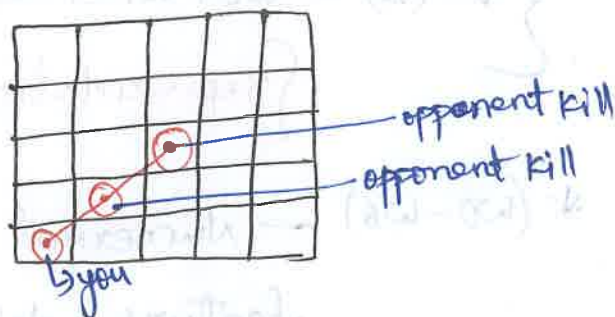
* only in diagonal direction.

* Jump over opponent.

* Target function = $v(b)$

* Board state = b .

* Legal moves set = b .



* b is final board state that is won, then $v(b) = 100$

* b is final board state that is lost, then $v(b) = -100$

* b is final board state that is drawn, then $v(b) = 0$

* If b is not final state then $v(b) = v(b')$.

$b' \rightarrow$ best final state.

3. Choosing A Representation for Target Function :-

For any board state, we calculate function "C" as linear combination of following board features i.e.,
"C(b)"

Features :-

- * x_1 — No. of black pieces on board
- * x_2 — No. of Red pieces on board
- * x_3 — No. of black kings on board
- * x_4 — No. of red kings on board.
- * x_5 — No. of black pieces threatened by red (black which can be beaten by red)
- * x_6 — No. of red pieces threatened by black.

$$V^{\wedge}(b) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + w_5 x_5 + w_6 x_6$$

{ representation of C(b) }

- * $(w_0 - w_6)$ — Numerical coefficients (or) weights of each feature determine the importance / weightage of features
- * $w_0 \rightarrow$ additive constant.

4. Choosing a Learning Algorithm for approximating the target function :-

- * To learn a target function (f) we need a set of

* describes a particular board state and training value

Ordered pair = $(b, V_{\text{train}}(b))$
 (training Example representation)

Ex:-

Black won the game

(ie, $x_2=0$, which means no red

$$V_{\text{train}}(b) = +100$$

$$b = (x_1=3, x_2=0, x_3=1, x_4=0, x_5=0, x_6=0)$$

$$\langle (x_1=3, x_2=0, x_3=1, x_4=0, x_5=0, x_6=0) + 100 \rangle$$

We need to do 2 steps in this phase.

1. Estimating training values:-

In every step, we consider Successor.

(depending on the next step of opponent).

$$V_{\text{train}}(b) \leftarrow V^{\wedge}(\text{Successor}(b))$$

↓
it represents the next board state.

{ Estimating that, this move will help/destroy opponent }

V^{\wedge} → Represents approximation

2. Adjusting the weights:- There are some algorithms

to find weights of linear functions.

Here we are using LMS (Least Mean Square).

(Used to minimize the error)

$$\text{Error} = (V(\text{train}(b)) - V^n(b))^2$$

If error = 0, No need to change weights

error = +ve, Each weight is increased \uparrow in proportion

error = -ve, Each weight is decreased \downarrow in proportion

* Final design :-

It has 4 different modules.

1. performance System

2. Critic

3. Generalizer

4. Experiment Generator.

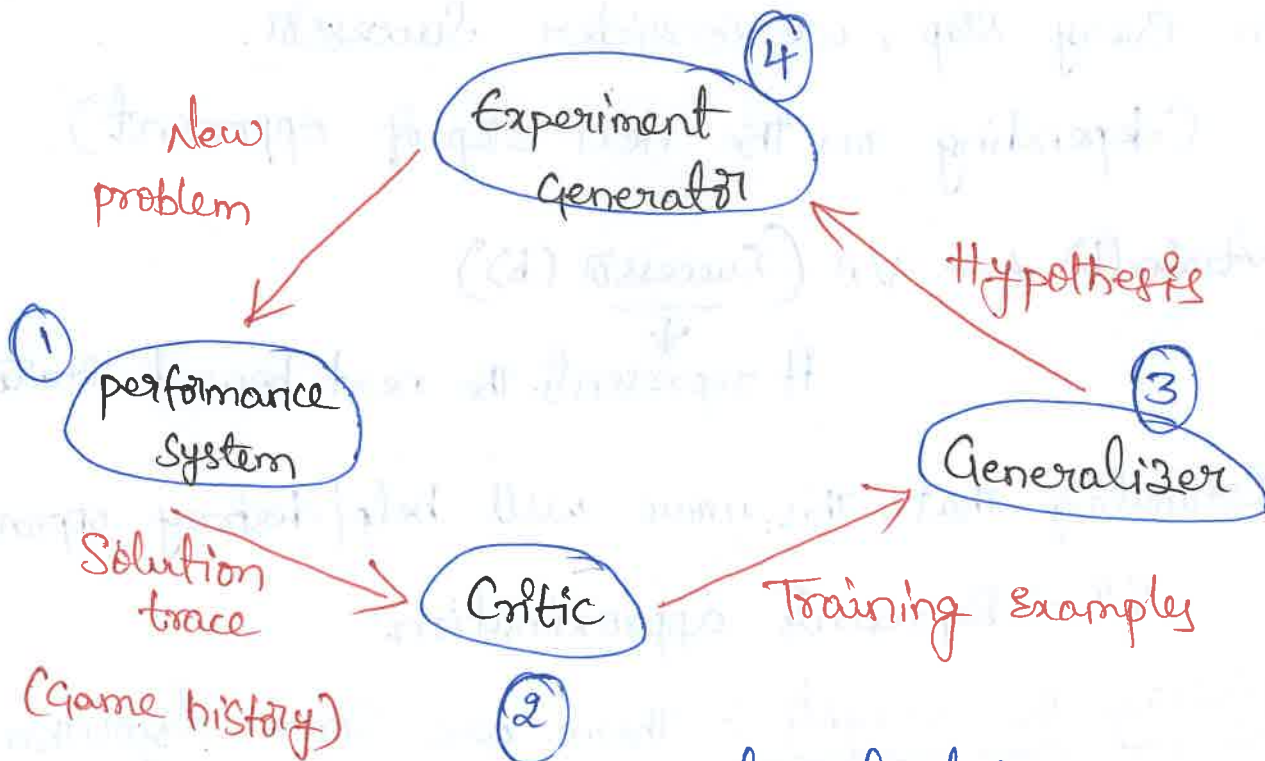


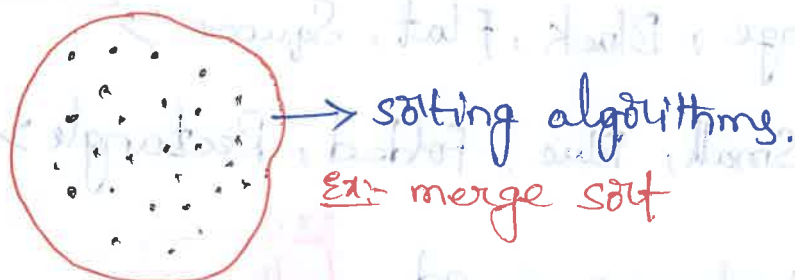
fig: final design

Concept Learning and the general to Specific ordering ⁽⁹⁾

perspectives of machine Learning :-

perspective of machine learning involves searching very large space of possible hypotheses to determine one that best fits the observed data and any prior knowledge held by learner.

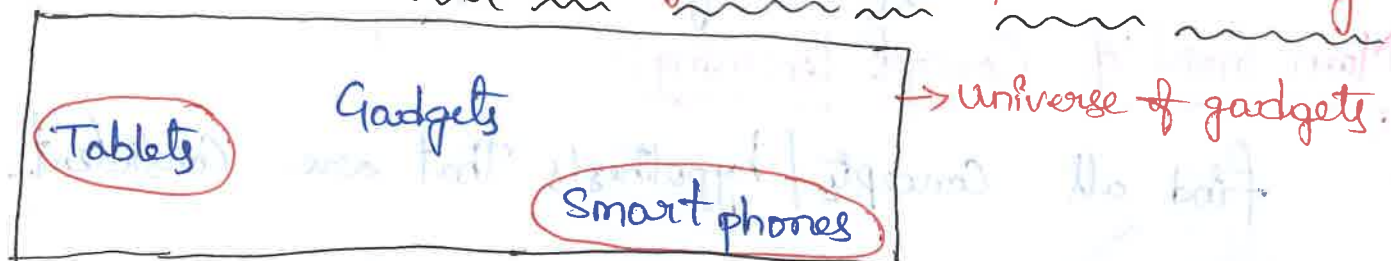
Ex:-



Issues in Machine Learning :-

1. What algorithm should be used?
2. Which algorithms perform best for which types of problems?
3. How much training data is sufficient? and testing data.
4. What kind of methods should be used?
5. What methods should be used to reduce learning overhead?
6. For which type of data which methods should be used?

Concept Learning and the general to Specific ordering :-



Features (Binary valued attributes) :-

Size :- Tablet Large, Smart phone Small (x₁)

Color :- Black, Blue (x₂)

Screen type :- Flat, Folded (x₃)

Shape :- Square, Rectangle (x₄)

Concept Learning :- Inferring/Assuming Boolean-valued function from Concept = $\langle x_1, x_2, x_3, x_4 \rangle$ { training example of it's i/p's and o/p's.

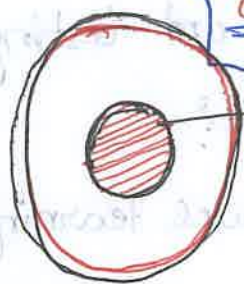
Tablet = $\langle \text{large, Black, Flat, Square} \rangle$

Smart phone = $\langle \text{Small, Blue, Folded, Rectangle} \rangle$

No. of possible instances = 2^d ex:- $2^4 = 16$

Total no. of possible Concepts = 2^{2^d} 2^{16}

? - it represents that any value is acceptable for this attribute.
φ - " " " " " " " " " " " "



Target Concept | hypothesis Space

ϕ (it represents Empty set of instances ^{-ve}).

$\langle \phi, \phi, \phi, \phi \rangle \rightarrow$ Reject all \rightarrow (Most Specific hypothesis).
 $\langle ?, ?, ?, ? \rangle \rightarrow$ Accept all

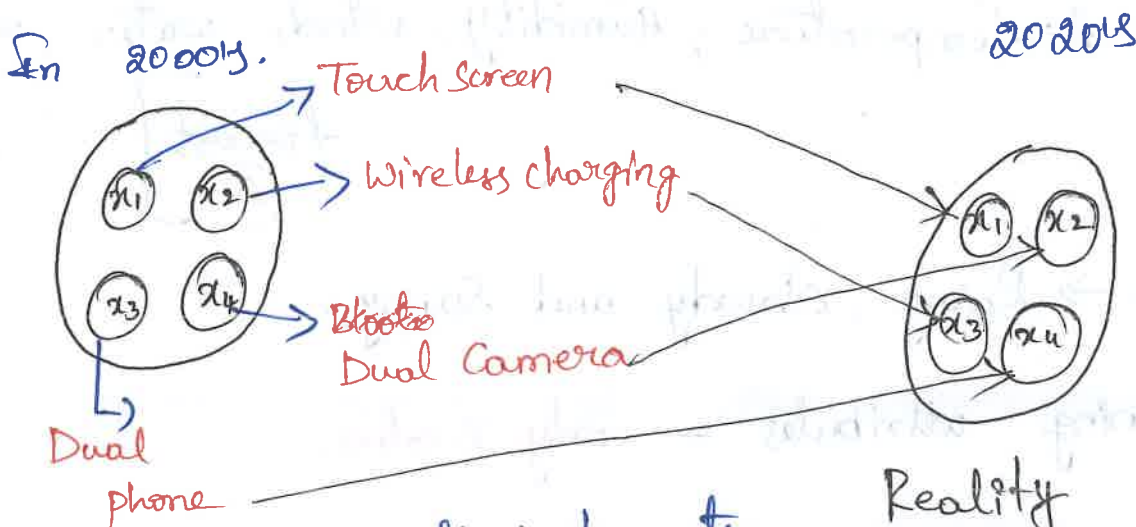
\hookrightarrow (Most general hypothesis).

∴ ϕ = phi ^{called}

Main Goal of Concept Learning :-

find all concepts/hypothesis that are consistent.

Example:-



Hypothesis

it implements which are consistent.

Reality

* Concept learning defines to find out the concepts, hypothesis are the always consistent.

Concept Learning as Search:-

Main goal of this search is to find the hypothesis that best fits the training examples.

Example:- Enjoy Spot learning task.

6 attributes = (Sky, Air-temperature, humidity, wind, water, and forecast). Attributes (forecast) concept

Example	Sky	Airtemp	Humidity	wind	water	forecast	Enjoy spot
1	Sunny	warm	Normal	Strong	warm	same	yes
2	Sunny	warm	high	Strong	warm	change	yes
3	Rainy	cold	high	strong	warm	change	No
4	Sunny	warm	high	strong	cool	change	yes

General-to-Specific Ordering of Hypothesis :-

(13)

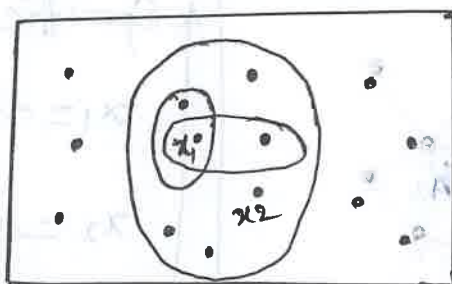
(More general than) \times

Def:- Let h_j and h_k be Boolean-valued functions defined over X . Then h_j is more general than or equal to h_k (written, $h_j \geq h_k$) if and only if.

① More general than:-

Let h_j and h_k be Boolean-valued functions defined over X . For any instance 'x' in 'X', hypothesis h in H , we can say 'x' satisfies 'h' iff $h(x) = 1$.

i.e., h_j is "more general than" h_k ($h_j \geq h_k$).

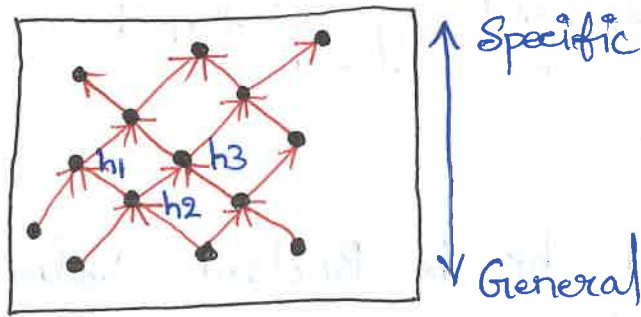


Instances X

② More general than or equal to:-

When set of training examples satisfy the 2 hypotheses h_j and h_k , h_j is "more general than or equal to" h_k ($h_j \geq h_k$) iff any instance x that satisfies h_k also satisfies h_j .

* \rightarrow we can also define the more-specific-than ordering.



$$(\forall x \in X) [(h_k(x)=1) \rightarrow (h_j(x)=1)]$$

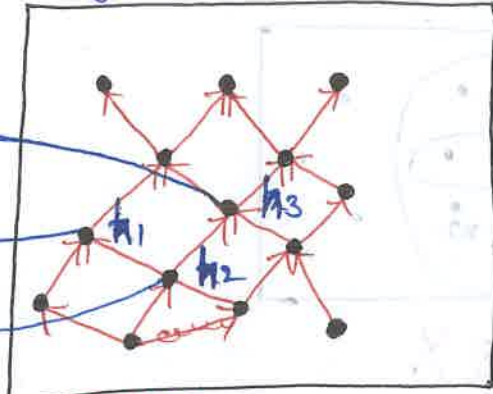
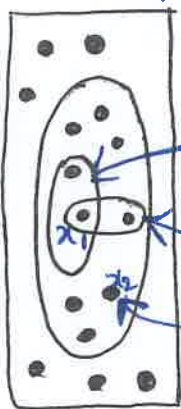
\uparrow
 x satisfies h_k .

③ General-to-Specific Ordering:-

Many Concept learning algorithms organize the search through the hypothesis space by taking advantage of a "Naturally occurring structure" over it.

Instance X

Hypothesis H



\uparrow Specific
 \downarrow General

$x_1 = \langle \text{Sunny, warm, High, Strong, Cool, Same} \rangle$
 $x_2 = \langle \text{Sunny, warm, High, Light, warm, Same} \rangle$

Suppose that h_1 & h_2 classify the examples.

eg:- let's take two hypotheses, and two training ex's.

$$h_1 = \langle \text{Sunny, ?, ?, Strong, ?, ?} \rangle$$

$$h_2 = \langle \text{Sunny, ?, ?, ?, ?, ?} \rangle$$

h_2 - is more general than h_1 .

* h_2 - it imposes fewer constraints on instances.

* h_2 - classify more the instances than h_1 .

Find-S: Finding a maximally Specific Hypothesis

- * The find-S algorithm is a basic concept learning algorithm in machine learning.
- * "The find-S algorithm finds the most specific hypothesis that fits all the positive examples"

Note:- The algorithm considers only the examples (training examples).

- * The find-S algorithm starts with the most specific hypothesis and generalizes this hypothesis until it cannot classify further positive training data (i.e., further classification will not be possible as it reaches the end.)

Algorithm:-

1. Initialize h to the most specific hypothesis
2. For every positive training instance x for each attribute constraint a in h , If the constraint a is satisfied by x ,
Then do nothing

Else

Replace a , in h by the next more general
Constraint that is satisfied by the x .

3. Output Hypothesis. h

[0x]

Algorithm!:-

Step 1! Initiate with most specific hypothesis,

$h_0 = \langle \phi, \phi, \phi, \phi, \phi \rangle$

Step 2! for each positive sample/example,

for each attribute

if (value = hypothesis value) = ignore.

else

Replace with most general hypothesis.

<u>Example</u>	<u>SKY</u>	<u>AirTemp</u>	<u>Humidity</u>	<u>wind</u>	<u>water</u>	<u>forecast</u>	<u>Enjoy Spot.</u>
1	Sunny	warm	Normal	Strong	warm	same	yes
2	Sunny	warm	High	Strong	warm	same	yes
3	Rainy	cold	High	Strong	warm	change	no
4	Sunny	warm	High	Strong	cool	change	yes

$$h_0 = \langle \phi, \phi, \phi, \phi, \phi, \phi \rangle$$

$$h_1 = \langle \text{Sunny}, \text{warm}, \text{Normal}, \text{Strong}, \text{warm same} \rangle$$

$$h_2 = \langle \text{Sunny}, \text{warm}, ?, \text{Strong}, \text{warm same} \rangle$$

$$h_3 = h_2 \quad [\because \text{it doesn't allow -ve hypothesis}]$$

$$h_4 = \langle \text{Sunny}, \text{warm}, ?, ?, ?, ? \rangle$$

Drawbacks :-

- * Considers only +ve values
- * h_4 may not be sole hypothesis that fits complete data

Version Space :-

Subset of hypothesis H consistent with the training examples D .

$$VS_{H,D} = \{ h \in H \mid \text{Consistent}(h, D) \}$$

where $H =$ Hypothesis space

$h =$ Subset of hypothesis space H .

$D =$ training examples.

Here Consistent :-

[where $x =$ instance

$$h(x) = C(x).$$

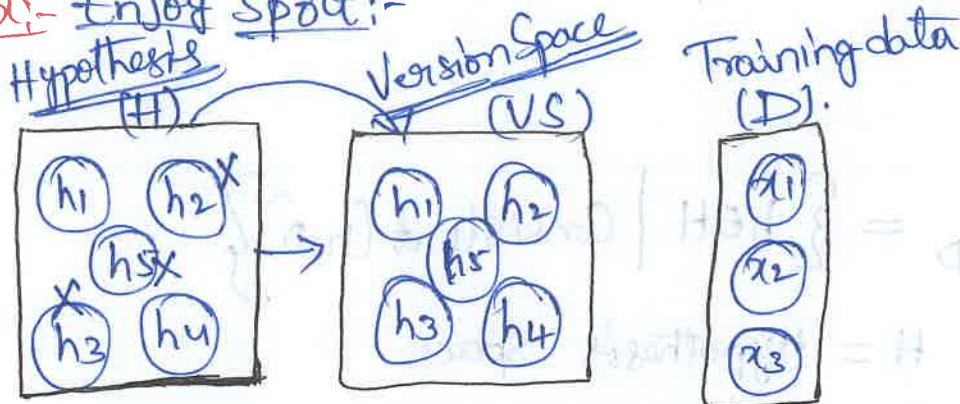
$h(x) =$ hypothesis of instance x
 $C(x) =$ target concept of instance x

Algorithm to obtain Version Space :-

(List then eliminate algorithm)

1. All the hypothesis H is kept in version space
(irrespective of the fact that it's consistent or not).
2. Now we keep on removing inconsistent hypotheses from version space.
→ for each training example $\langle x, c(x) \rangle$ remove any hypothesis that is $h(x) \neq c(x)$.
3. O/p the list of hypothesis into version space after checking for all training examples.

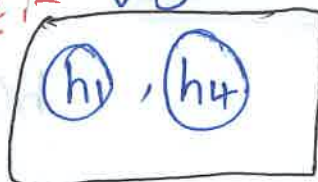
Ex:- Enjoy Sport:-



$h_1(x_1) = c(x_1)$		$h_2(x_1) \neq c(x_1)$		$h_4(x_1) = c(x_1)$
$h_1(x_2) \neq c(x_2)$		$h_2(x_2) = c(x_2)$		$h_4(x_2) = c(x_2)$
		$h_2(x_3) = c(x_3)$		$h_4(x_3) = c(x_3)$

$h_5(x_1) = c(x_1)$
 $h_5(x_2) = c(x_2)$
 $h_5(x_3) = c(x_3)$

final output:- V.S



Candidate Elimination Algorithm:-

- * It's an extended form of find-S algorithm.
- * It considers both +ve and -ve examples.
- * Here +ve examples are used as find-S algorithm (i.e., specific to general).
- * While -ve examples from general to specific.

$$S = \{ \phi \phi \phi \phi \phi \} \quad + \downarrow$$

$$G = \{ ? ? ? ? ? \} \quad - \uparrow$$

Algorithm:-

1. Initialize both general and specific hypotheses (S and G).

$$S = \{ \phi, \phi, \phi, \phi \dots \phi \} \quad \text{depends on}$$

$$G = \{ ?, ?, ?, ? \dots ? \} \quad \text{no. of attributes}$$

2. For each example,

if example is positive

make specific to general.

else example is -ve.

make general to specific.

ex:- enjoy sport. as a example.

$$S_0 = \{\phi, \phi, \phi, \phi, \phi, \phi\}, G_0 = \{?, ?, ?, ?, ?, ?\} \textcircled{20}$$

1. tue,

$$S_1 = \{\text{'Sunny'}, \text{'warm'}, \text{'normal'}, \text{'Strong'}, \text{'warm'}, \text{'same'}\}$$

$$G_1 = \{?, ?, ?, ?, ?, ?\}$$

2. tue (Specific to general).

$$S_2 = \{\text{'Sunny'}, \text{'warm'}, ?, \text{'Strong'}, \text{'warm'}, \text{'same'}\}$$

$$G_2 = \{?, ?, ?, ?, ?, ?\}$$

$$3. S_3 = \{\text{'Sunny'}, \text{'warm'}, ?, \text{'Strong'}, \text{'warm'}, \text{'same'}\}$$

$$G_3 = \{\langle \text{'Sunny'}, ? ? ? ? ? \rangle, \langle ?, \text{'warm'}, ? ? ? ? ? \rangle, \langle ? ? ? ? ? \text{'same'} \rangle\}$$

$$4. S_4 = \{\text{'Sunny'}, \text{'warm'}, ?, \text{'Strong'}, ? ? \}$$

$$G_4 = \{\langle \text{'Sunny'}, ? ? ? ? ? \rangle, \langle ? \text{'warm'} ? ? ? ? ? \rangle\}$$

Final output:-

$$S_4 = \{\text{Sunny}, \text{warm}, ?, \text{Strong}, ?, ?\}$$

$$G_4 = \{\langle \text{'Sunny'}, ? ? ? ? ? \rangle, \langle ? \text{'warm'} ? ? ? ? ? \rangle\}$$

Data Set:-

Example "Enjoy Sport" Table.

Remarks on Version Space and Candidate-Elimination

① Will the CANDIDATE-ELEMINATION Algorithm Converge to the correct hypothesis?

* The learned version space correctly describes the target concept, provided:

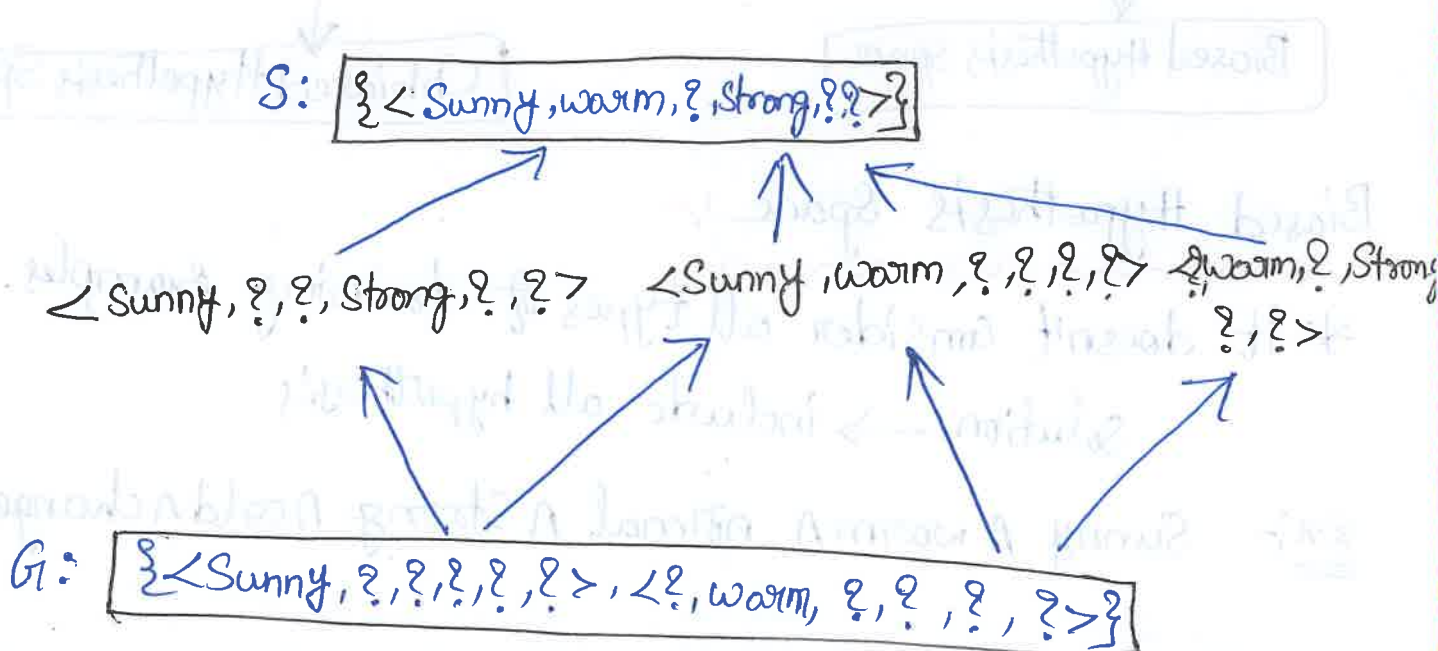
1. There are no errors in the training Examples
2. There is some hypothesis that correctly describes the target concept.

* If S and G converge to a single hypothesis the concept is exactly learned.

* An empty version space means no hypothesis in H is consistent with training examples.

② How can partially learned concepts be used?

* The learner is required to classify new instances that it has not yet observed.



③ What training example should the learner request next?

* Consider the version space learned from the four training examples of the enjoy spot concept.

* What would be a good query for the learner to pose at this point?

* What is a good query strategy in general?

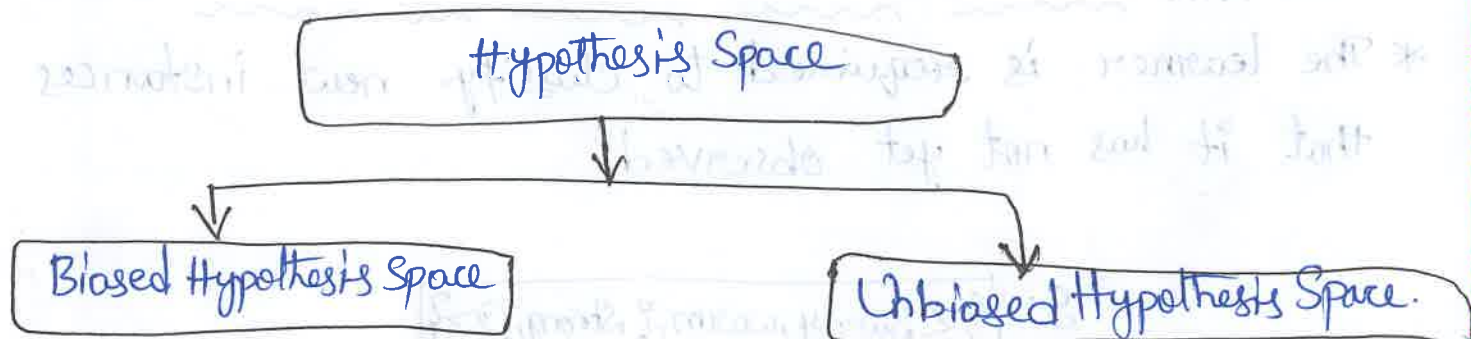
Inductive Bias!:- It refers to a set of assumptions that the learner uses to predict o/p's given unseen i/p's.

Remarks on CE and VS Algorithms!:-

1. Will the CE algorithm give us correct hypotheses?
2. What training examples should the learner request next?

Inductive Learning!:-

Derive rules from examples



Biased Hypothesis Space!:-

* It doesn't consider all types of training examples.

Solution → include all hypothesis's

Ex:- Sunny \wedge warm \wedge normal \wedge Strong \wedge cold \wedge change = Yes.

Hypothesis Space!

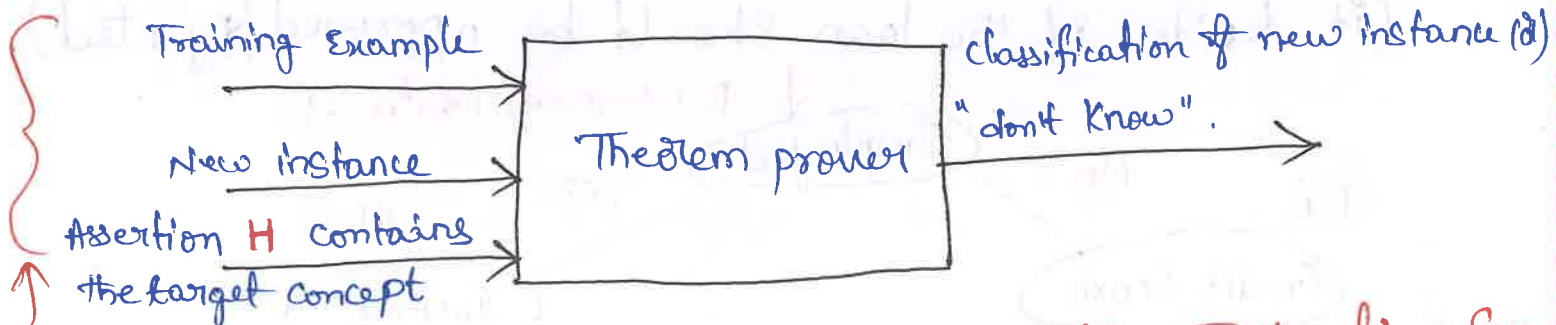
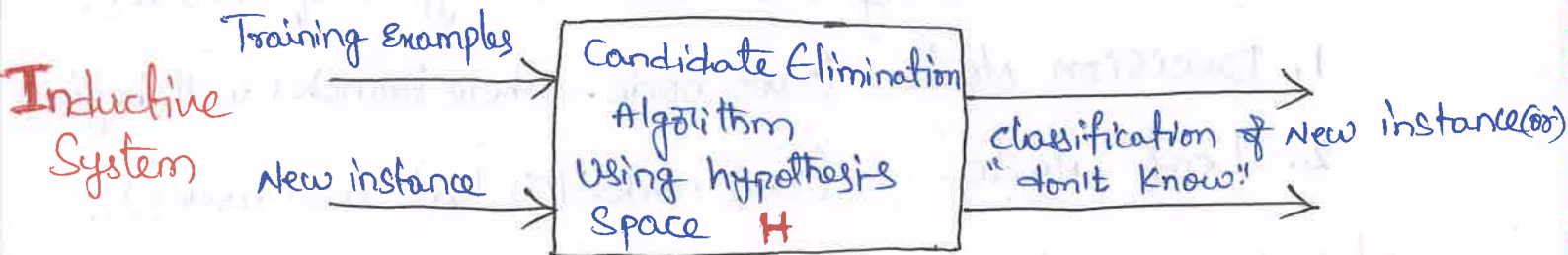
* providing a hypothesis capable of representing set of all examples.

possible Instances :- $3 \times 2 \times 2 \times 2 \times 2 \times 2 = 96$.

Target Concepts :- 2^96 (huge).

(practically not possible),
it takes a lot of time.

<u>Instance X</u>	→	<u>Attributes</u>
Sky	<u>3</u> →	Sunny, cloudy, Rainy
AirTemp	<u>2</u> →	warm, cold
Humidity	<u>2</u> →	Normal, high
Wind	<u>2</u> →	Strong, weak
water	<u>2</u> →	warm, cold
forecast	<u>2</u> →	same, change.



Equivalent Deductive System.

Fig:- Modeling Inductive Sys.

Decision Tree Learning (ID3):

Decision Tree Learning is a method that uses inductive inference to approximate a target function, which will produce discrete values. And the following features;

- * widely used
- * Robust to noisy data.
- * Considers a practical method for learning disjunctive expressions.

* It's a tree structured classifier, where

① → Internal nodes represent the features;

② → Branches represent the decision rules;

③ → Each leaf node represents the Outcomes;

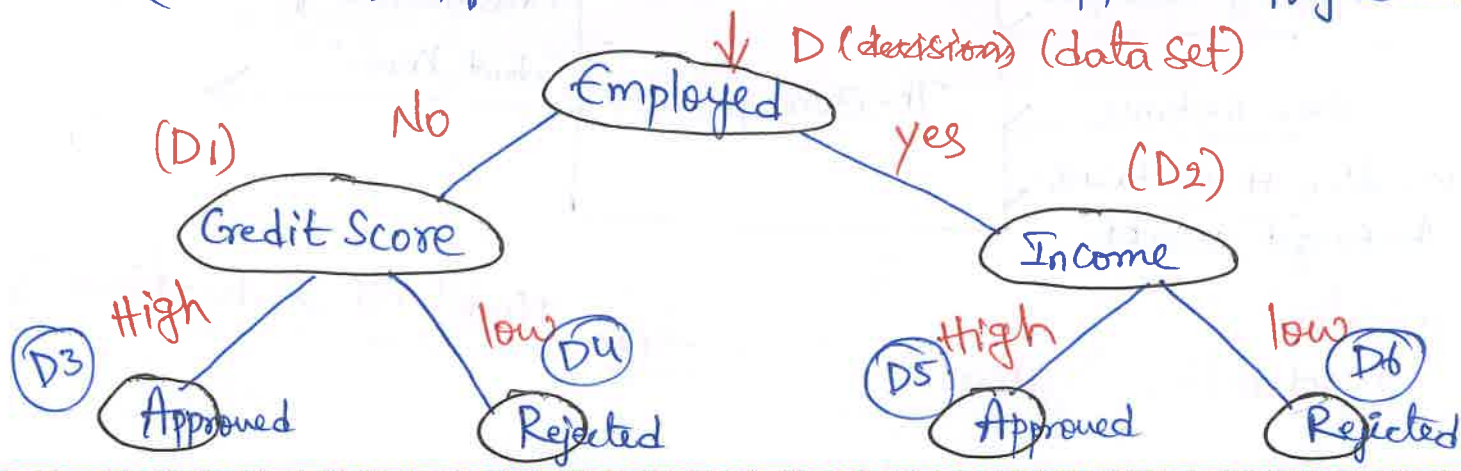
In decision tree mainly we have 2 types of Nodes.

1. Decision Node (root node - where branches will begin)

2. Leaf Node. (End nodes or last row nodes).

Example: Loan System

(It decides if the loan should be approved/rejected).



Algorithm Decision Tree Terminologies :-

Root/Decision Node :- It is the node where the decision tree starts. It represents entire dataset.

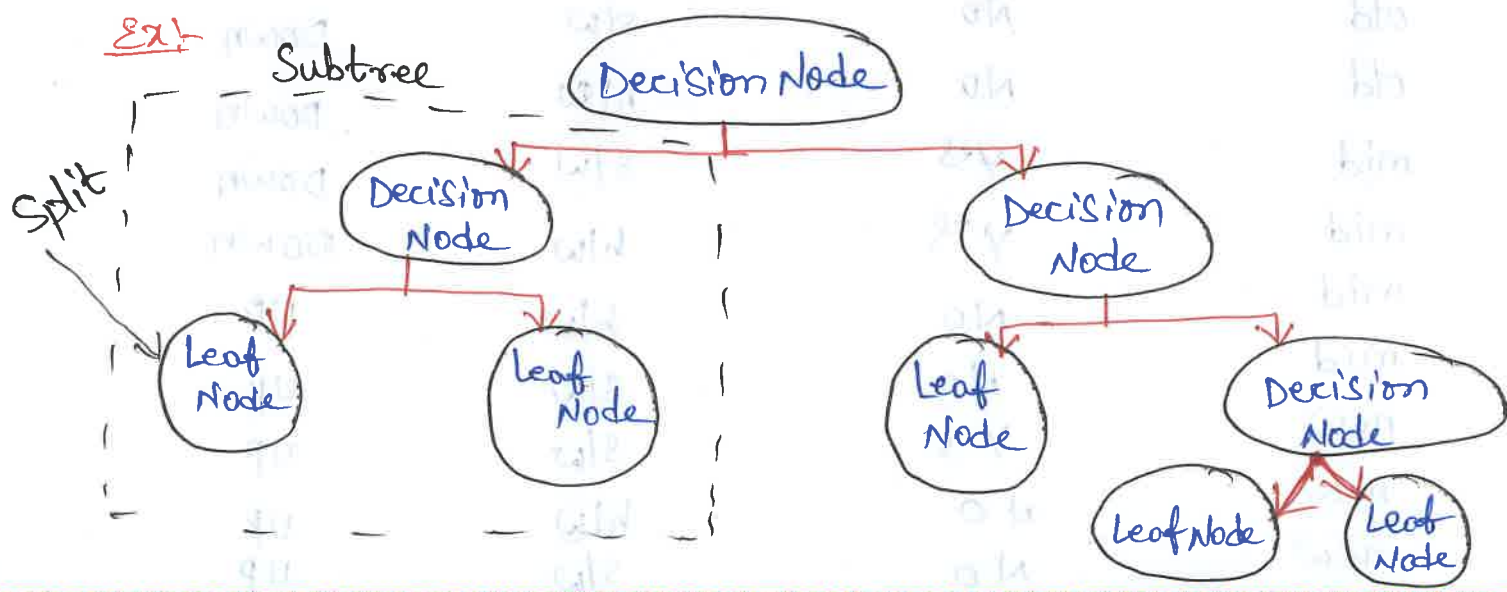
Leaf Node :- Leaf nodes are the final output node, and tree can't be segregated further after getting a leaf node.

Splitting :- Splitting is the process of dividing the decision node into subnodes according to given conditions.

Branch/
Sub-tree :- A tree formed by splitting the tree

Pruning :- It is a process of removing unwanted branches from the tree

Depth :- The depth of a node is the minimum no. of steps required to reach the node from the root.



Algorithm of Decision Tree Learning :-

Step 1 : In the given dataset, choose a target attribute.

Step 2 : Calculate information gain of target attribute.

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

Step 3 : For remaining attributes, find entropy.

$$[\text{Entropy} = IG \times \text{probability}]$$

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

Step 4 : calculate gain = $IG - E(A)$

Example :-

Age	Competition	Type	profit
old	yes	s/w	Down
old	No	s/w	Down
old	No	h/w	Down
mid	yes	s/w	Down
mid	yes	h/w	Down
mid	No	h/w	UP
mid	No	s/w	UP
new	yes	s/w	UP
new	No	h/w	UP
new	No	s/w	UP

Step - 1

Target Attribute = "profit".

Step - 2

Information gain

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

$$\rightarrow P = \text{Count (down)} = 5$$

$$\rightarrow N = \text{Count (up)} = 5$$

$$= \frac{-5}{5+5} \log_2 \left(\frac{5}{5+5} \right) - \frac{5}{5+5} \log_2 \left(\frac{5}{5+5} \right)$$

$$= \frac{-5}{10} \log_2 \left(\frac{5}{10} \right) - \frac{5}{10} \log_2 \left(\frac{5}{10} \right)$$

$$= - \left[\frac{1}{2} \log_2 (2^{-1}) + \frac{1}{2} \log_2 (2^{-1}) \right] \quad [\because \log_x m^n = n \log_x m]$$

$$= - \left[\frac{1}{2} \times -1 \log_2 2 + \frac{1}{2} \times -1 \log_2 2 \right] \quad [\because \log_2 2 = 1]$$

$$= - \left[\frac{1}{2} \times -1 + \frac{1}{2} \times -1 \right]$$

$$= - \left[-\frac{1}{2} - \frac{1}{2} \right] = - \left[\frac{-1-1}{2} \right] = - \left[\frac{-2}{2} \right] = -[-1] = 1$$

$$\therefore \text{IG} = 1.$$

Step - 3

calculate entropy for remaining attributes.

$$E(A) = \sum \frac{P_i + N_i}{P+N} I(P_i, N_i) \quad [\text{IG} \times \text{probability}]$$

Age:

(1) prepare a table for each attribute.

rows - Undertaken attribute values
(old, mid, new).

columns - target attribute values.
(down, up).

Age values	(P) Down	(N) UP
old	3	0
mid	2	2
new	0	3

$$\text{Entropy} = \text{IG} \times \text{Probability}$$

$$P = \text{down count}$$

$$N = \text{up count}$$

$$1. \quad \text{IG}(\text{old}) = - \left(\frac{3}{3} \log_2 \left(\frac{3}{3} \right) + \left(\frac{0}{3} \right) \log_2 \left(\frac{0}{3} \right) \right) = 0.$$

$$\text{probability} = 3/10$$

$$[\because \log 1 = 0]$$

$$\text{Entropy}(\text{old}) = 0 \times \frac{3}{10} = 0.$$

$$2. \quad \text{IG}(\text{mid}) = - \left(\frac{2}{4} \log_2 \left(\frac{2}{4} \right) + \frac{2}{4} \log_2 \left(\frac{2}{4} \right) \right)$$

$$= - \left(\frac{1}{2} \log_2 \left(\frac{1}{2} \right) + \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right)$$

$$= - \left(\frac{1}{2} \log_2 2^{-1} + \frac{1}{2} \log_2 2^{-1} \right)$$

$$= - \left(-\frac{1}{2} + \left(-\frac{1}{2} \right) \right) = -(-1) = 1.$$

$$\text{probability} = 4/10$$

$$\text{Entropy}(\text{mid}) = 1 \times \frac{4}{10} = 0.4.$$

$$3. \text{IG (new)} = - \left[\frac{0}{3} \log_2 \left(\frac{0}{3} \right) + \frac{3}{3} \log_2 \left(\frac{3}{3} \right) \right] = 0$$

$$\text{probability} = 3/10$$

$$\text{Entropy} = 0 \times \frac{3}{10} = 0.$$

$$\text{Entropy (Age)} = E(O) + E(M) + E(N)$$

$$= 0 + 0.4 + 0$$

$$= 0.4$$

Step 4 :-

$$\text{Gain} = \text{IG} - E(A) = 1 - 0.4$$

$$= 0.6.$$

In the same way, calculate Gain of other attributes.

$$\text{Gain (Competition)} = 0.124$$

$$\text{Gain (Tyre)} = 0$$

$$\text{Gain (Age)} = 0.6$$

Highest gain \rightarrow root node **Age**.



Here

old - All down

mid - Some down & Some up

new - All up

Why Competition \rightarrow next highest gain.

Why not other attribute (Type).

Type \rightarrow gain = 0 and even in DT also, no requirement with it.

\therefore It can be ignored.

* Appropriate problems for Decision Tree Learning:-

1. Instances are represented by attribute value pairs.

\rightarrow when a ϕ we have a finite list of attributes & each attribute contains small no. of distinct values, then, it's easier for the decision tree to reach the solution.

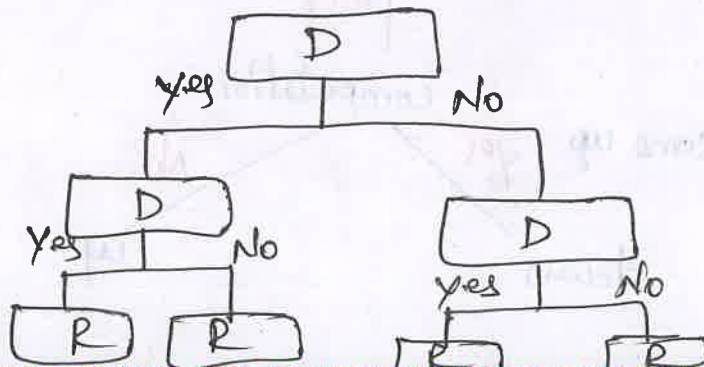
ex:-

Student Branch — (CST, ECT, IT etc). finite list.

2. The target function has discrete output.

\rightarrow it works with simplest way, i.e., when there are 2 possible ways (Boolean classification).

\rightarrow However it's easy to extend the decision tree to with more than 2 possible op values.



③ Disjunctive descriptions may be required.

↳ Decision tree naturally represents disjunctive expressions i.e., it reduces large consideration set to a more manageable no. of alternatives.

④ The training data may contain errors.

↳ errors in the classification of examples are handled well by decision trees & making them a robust machine learning method.

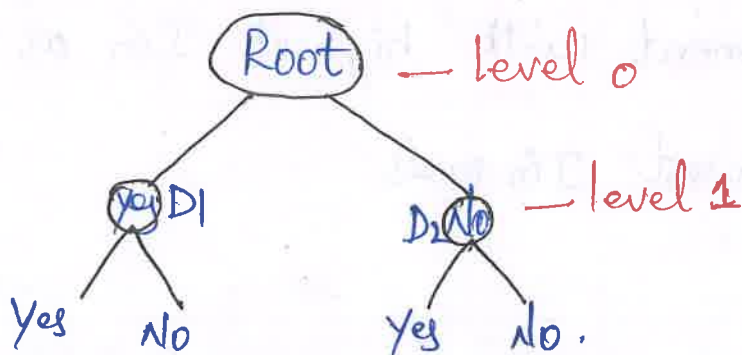
⑤ The training data may contain missing attribute values.

↳ Decision tree methods can be used even when some training examples have unknown values.

Hypothesis Space Search in Decision Tree Learning:-

- * ID3 can be characterized as searching a space of hypothesis for one that fits the training examples.
- * ID3 will search set of possible decision trees from available hypothesis.

ID3 performs simple to complex searching



- * first, start with empty tree and keep on adding.
- * Every discrete valued (finite) function can be described by some decision tree.
- * Avoids major risk of searching in incomplete hypothesis.
- * it has only single current hypothesis.
- * it can't determine alternative decision trees.
- * Back tracking is not possible.
- * it can also handle noisy data.

Inductive Bias in Decision Tree Algorithm:-

↳ set of assumptions.

Inductive bias of ID3 consists of describing the basis by which ID3 chooses one consistent decision tree over all the possible DT's.

ID3 Search Strategy:-

- * Selects in favour of shorter trees over longer ones.
- * Selects element with highest IG_i as root attribute over lowest IG_i ones.

Types of Inductive Bias:-

1. Restrictive Bias (Based on conditions).

2. preference " (Based on priorities).

Restrictive Biases and preference Biases:

There is difference b/w the inductive bias exhibited by ID3 and by the candidate elimination algorithm.

* ID3 Searches incompletely through this space, from simple to complex hypotheses, until its termination condition is met. Its inductive bias is solely a consequence of the ordering of hypotheses by its search strategy.

* The inductive bias of ID3 is thus a preference for certain hypotheses over others, with no hard restriction on the hypotheses that can be eventually enumerated. This form of bias is typically called a preferential bias.

* Candidate-Algorithm searches this space completely, finding every hypothesis consistent with the training data. Its inductive bias is solely a consequence of the expressive power of its hypothesis representation.

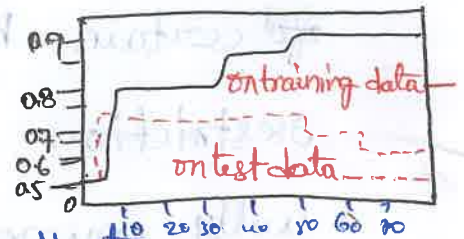
- * A preference bias is more desirable than a restriction bias because it allows the learner to work within a complete hypothesis space that is assured to contain the unknown target function.
- * A restriction bias that strictly limits the set of potential hypothesis is generally less desirable, because it introduces the possibility of excluding the unknown target function altogether.

Why short hypotheses?

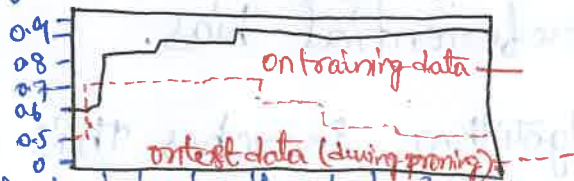
Occam's razor: prefer the simplest hypothesis that fits the data.

Issues in Decision Tree Learning!

1. Overfitting the data



2. Incorporating continuous valued attributes



3. Handling training examples with missing attribute values.

4. Handling attributes with different costs.

5. Alternative measures for selecting attributes!