

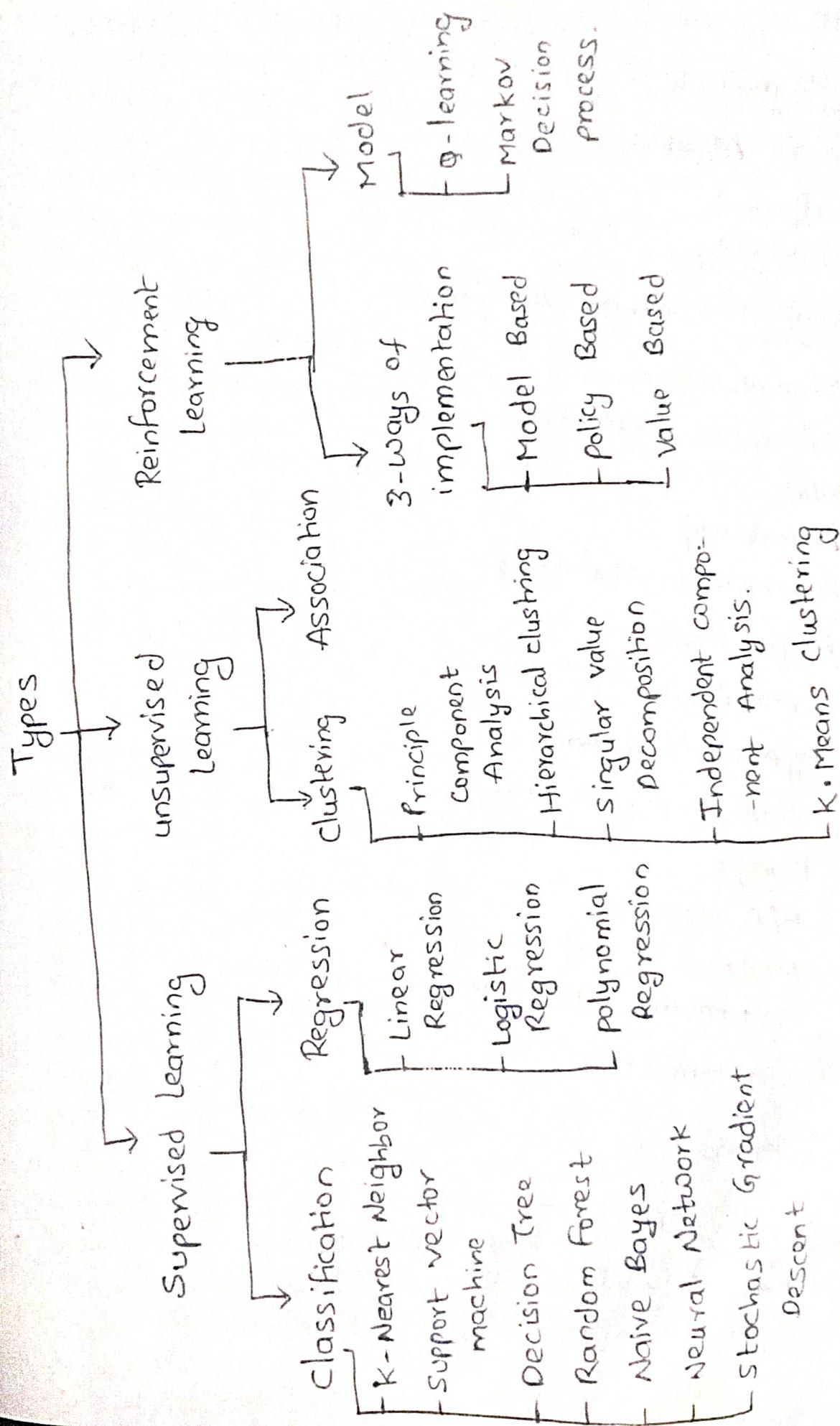
Machine
Learning

N.Likhith Reddy

Definition:-

* It Learn from Data & Solves problems.

Types:-

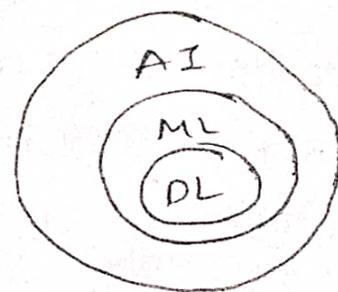


Machine Learning

It is a field of study that gives computers a capability to learn without being explicitly programmed.

Example:- online shopping.

→ Machine adapts to the user based on data.



Well posed learning Problem:-

An agent solves a problem or task τ ,
Performance P and gain some experience E .
If P is measured at τ it can improve E .

Example:- 1. Handwritten recognition Problem.
2. Robo driving learning Problem.

<u>Problem</u>	<u>task(r)</u>	<u>Performance(P)</u>	<u>experience(e)</u>
1. Hand writing recognition learning	classifying the Images & Text	Better classification	A database of home work by
2. 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.
3. playing chess learning	Playins against opponen to win game	make perfect moves to win Game	It plays itself to improve

Perspectives of machine learning:

perspective of machine learning involves
Searchins very large Space of possible
hypotheser to determine one that best fit
the observed data and any prior knowledge
held by learner

* It is to enable computers to learn from
data and make predictions without
being explicitly programmed to perform task

Ex :-



Sorting Algorithms
(merge sort)

Issues in Machine Learning:-

1. what algorithms should be used?
2. which algorithms perform best for which types of problems?
3. How much training data is sufficient?
and types 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?

Designing A learning System:-

To get a successful learning system,
it should be designed for a proper
design, several steps should be followed

* To get perfect & efficient system.

6 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.

Step-1:- choosing a training experience.

In choosing a training experience,
3 attributes are taken.

1. type of feedback.
2. degree.
3. distribution of examples.

i) whether the training experience provides direct or indirect feedback regarding the choices made by performance system.

Direct feedback

Indirect feedback

Ex:- Checkers game, learning driving.

2) Degree to which learner will control the sequence of training.

Ex:- learning driving.

/
with trainer,
help

\
with trainers
partial help

Completely on
your own.

3) How will it represent the distribution of examples over which the performance of final system is measured.

- more possible combinations.
- more situations, more examples.
- learning over distributed range of examples.

Ex:- learning driving.

Step-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
 - ↓
 - one move
 - ↓
 - target move
 - travel only in forward dir
 - only one move per chance.
 - only in diagonal direction.
 - jump over opponent
- target function $\Rightarrow v(b)$
Board style $\Rightarrow b$
legal moves $\Rightarrow B$

1. b is final board state that is won,
then $v(b) = 100$.
2. b is final board state that is lost,
then ($v(b) = -100$).
3. b is final board state that is draw,

then $v(b) = 0$

4. If b is not final state then $v(b) = v(b')$
 $b' \rightarrow$ best final state

Step-3 :- choosing A Representation
for target function:-

for any board state, we calculate
function ' c ' as linear combination
of following board features $c(b)$

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's which can be beaten
by red)

x_6 - No. of red pieces threatened by
black

$$v(b) = w_0 + w_1x_1 + w_2x_2 + w_3x_3$$

(representation of $v(b)$)

(w_0, w_1) = numerical coefficients (or)
weight of each feature determine
the importance / weightage of features
 $w_0 \rightarrow$ additive constant.

Step-4: choosing a learning algorithm
for approximating the target function

To learn a target function (f) we
need a set of training examples.

(describe a particular board state (b)
and training value $v_{\text{train}}(b)$).

ordered pair = $(b, v_{\text{train}}(b))$

(training example representation)

example: black won the game

(i.e. $x_2 = 0$, 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)

$$v_{\text{train}}(b) \leftarrow \overline{v^*(\text{successor}(b))}$$

\downarrow
represent the next board state.

(estimating that this move will help/destroy opponent)

v^* → represents approximation.

2 Adjusting the weight

There are some algorithms to find weights of linear functions

Here, we are using LMS (least mean square)
(used to minimise the error)

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

If $\text{error} = 0$, no need to change weight.

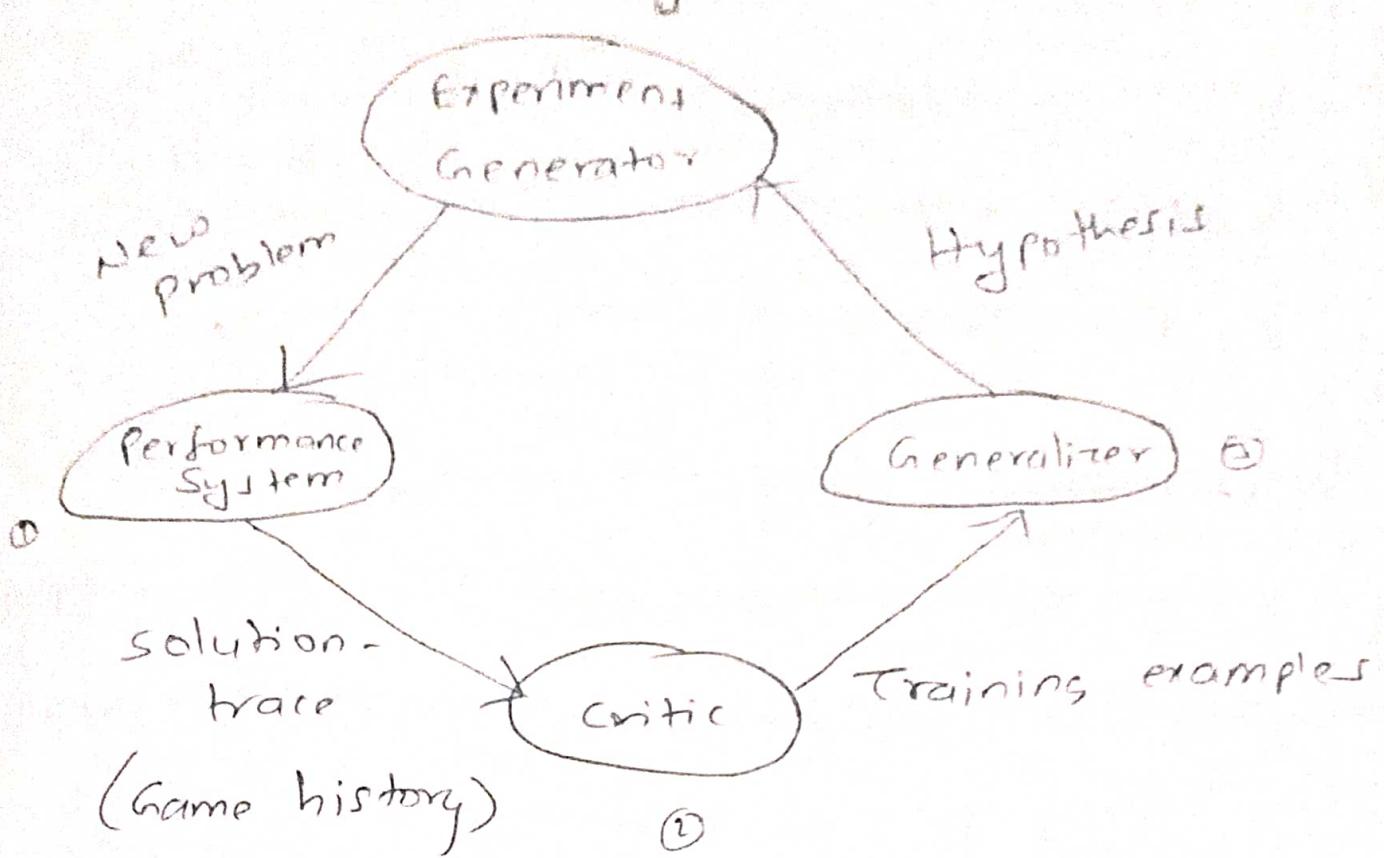
If error is positive, each weight is increased in proportion.

If error is negative, each weight is decreased in proportion.

Final Design:-

has 4 different modules.

1. performance system.
2. critic
3. Generalizer
4. Experiment generator.



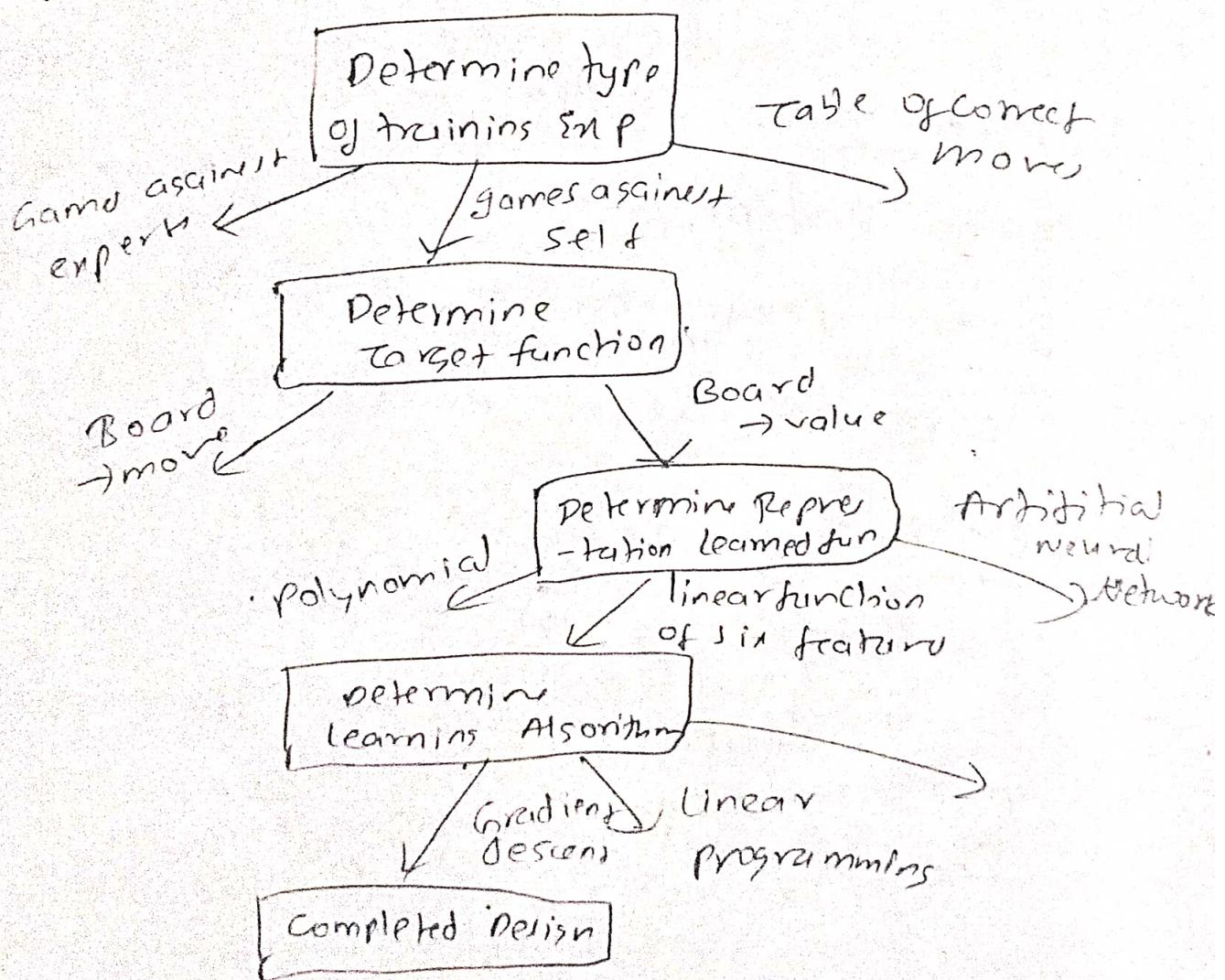
1. Performance System:- It is the module

that performs and must solve the given performance task by using the learned target functions.

2. Critic:- takes as input the history

or trace of the game and produces as output a set of training examples of a target function.

3. Generalizer :- takes as input the training examples and produces an output hypothesis that is its estimate of the target function.
4. Experiment Generator :- It takes as input the current hypothesis and outputs a new problem for the performance system to explore.



Concept Learning :-

Acquiring general concepts from specific training examples.

1. Learning General Concepts :-

* The learning process involves understanding general concepts from specific examples.

* like "bird", "car", "study more for the exam".

2. Concept as subset Definitions :-

* General concepts represent subsets of objects or events within a larger set.

* Example: "bird" is a subset of all animals

3. Concepts as Boolean Functions :-

* Alternatively, concept can be seen as boolean-valued functions.

* Example: "A function could return "true" for birds and "false" for other animals."

4. Automatic interface of concepts:-

- * The chapter shows automatically deducing the general definition of a concept.

5. Concept learning:-

- * It involves inferring boolean-valued functions.
- * Input :- Specific objects.
Output :- whether they belongs to the concept.
- * The process is about approximating a boolean-valued function using these examples.

Concept Learning Task :-

1. Example Task :- "Enjoy Sport" concept
* Target concept :- "Days on which my friend Aldo enjoys his favorite water sport."

* Table provide example days, each with attributes and the "Enjoy Sport" value.

Example	Sky	Air Temp	Humidity	Wind	Water
1	Sunny	warm	Normal	strong	warm
2	Sunny	warm	High	strong	warm
3	Rainy	cold	High	strong	warm
4	Sunny	warm	High	strong	cool

Forecast	Enjoy sport
Same	Yes
Same	Yes
Change	No
Change	Yes

Table 2.1

2. Hypothesis Representation

Simple representation: Each hypothesis consists of a conjunction of constraints on six attributes: sky, airtemp, humidity, wind, water, forecast.

- * Constraints can be "?", a specific value, or "0"

3. General Hypotheses

- Most general hypothesis: (?, ?, ?, ?, ?, ?)
- Most specific hypothesis: (0, 0, 0, 0, 0, 0)

4. Learnings Task Summary

- * Concept Learnings Task: - Learnings: the set of days for which EnjoySport = yes by describing this set through attribute constraints.
- * Components of the task: - Instances, Target concept, Hypothesis, Training examples.

5. Notation

- * instances x_i : Possible days with attributes (sky, Airtemp, Humidity, Wind, water, forecast)

- * Target concept c : $f(x) = \text{enjoy_park} : x + f_0, y$.
- * Hypotheses H : Conjunction of attribute constraints.
- * Training examples D : Positive and negative examples of the target function.

6. Inductive Learning Hypothesis

- * Assumption: The best hypothesis for unseen instances is the one that best fits observed training data.

- * Fundamental assumption of inductive learning.

Concept Learning as Search (Section 2.3)

- * Concept learning as a search problem
- * Concept learning can be seen as a search task in a space of hypotheses defined by the hypothesis representation.

- * The goal is to find the hypothesis that best fits the provided training examples.
- * The choice of hypothesis representation defines the space of hypotheses available to the learning algorithm.
- * The example of the "Enjoy Sport" learning task involves a finite hypothesis space but real-world tasks can have larger or infinite spaces.

2. Strategies for Searching Hypothesis

Space for

- * Learning efficiently algorithms need to search through the hypothesis space.
- * The text introduces the concept of a "general-to-specific ordering of hypothesis."

General to - Specific ordering of Hypotheses

- * Many Concept learning algorithms reverse the structure of hypotheses in a general-to-specific order.
- * Hypotheses can be compared in terms of being more general or more specific.
- * A hypothesis is considered "more general" if it classifies at least as many instances as another hypothesis.
- * The formal definition of the "more-general-than" relation is provided.
- * The relation is reflexive, anti-symmetric, and transitive, making it a partial order over the hypothesis space.

4. Illustrating General-to-Specific

Ordering

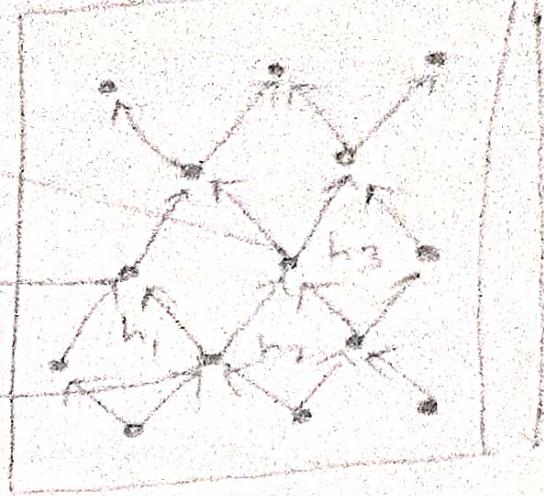
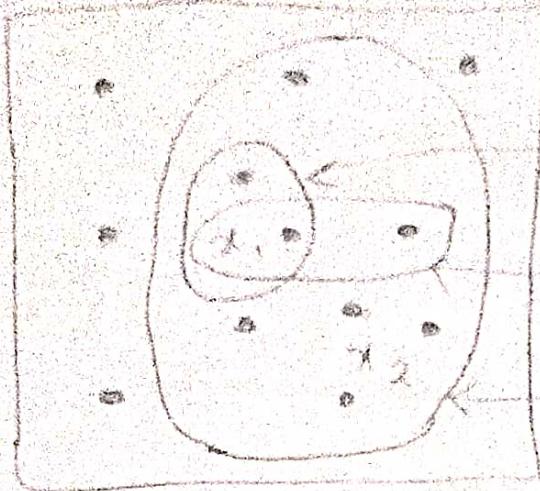
- * Example hypotheses h_1 , h_2 , and h_3 from the "Enjoy sport" task are shown in figure 2-1
- * h_2 is more general than h_1 and h_3 .
- * h_2 is more general than h_3 .
- * Neither h_1 nor h_3 is more general than the other, they don't subsume each other.

5. Importance of General-to-Specific ordering

- * The partial order (\leq) provides structure in the hypothesis space for concept learning.
- * Concept learning algorithms can be w.e. this structure for more efficient hypothesis search.

To start with

Hypothesis H



General

$x_1 = \langle \text{Sunny}, \text{warm}, \text{High}, \text{Strong},$
 $\text{cool}, \text{Same} \rangle$

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

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

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

$h_3 = \langle \text{Sunny}, ?, ?, ?,$

$\text{cool}, ? \rangle$

Fig 2.1

FIND-S Algorithm and its Application:

I. Using the More-General-Than-partial
orderings:-

* Concept learning involves searching
for a hypothesis consistent with
training examples.

- * The search starts with most specific hypothesis and generalizes when necessary to cover positive samples.
- * The more-geared-than partial ordering is employed to structure the search.
- * The "find-s" algorithm (Table 2.3) introduced for this purpose.

1. Initialize h to the most specific hypothesis in H .

2. for each positive training instance x .

- for each attribute constraint a_i in h
 - if the constraint a_i is satisfied by x
Then do nothing
 - Else replace a_i in h by the next
~~more general constraint~~
more general constraint that is satisfied by x .

3. output hypothesis h .

Table 2.3
find-s Algorithm.

Mysteries: the Pupil's Algorithm

Problem 8: Sequence of Training

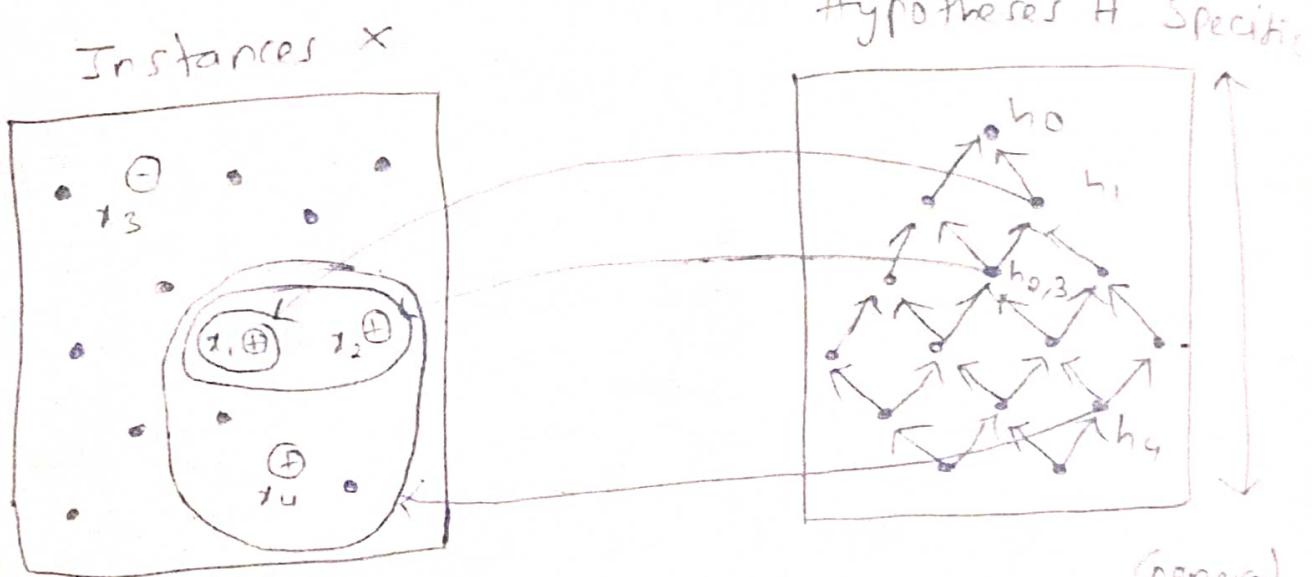
Example: draw the "Baby spot" rule
Table 2.1

- * Start with the most specific hypothesis and generalize as needed to cover positive examples.
- * The algorithm ignores negative examples as long as the hypothesis is consistent with them.

3. Guarantees and Questions Raised -

- * The finds algorithm guarantees the most specific hypothesis consistent with the training data in conjunction based hypothesis spaces
- * important questions remain:
 - a. Has the learner converged to the correct target concept?
 - b. Why prefer the most specific hypothesis? is it always the best choice

- c. Handling inconsistent training examples
(error or noise).
- d. what if there are several maximally specific, consistent hypotheses in non-conjunction-based spaces?
- e. defining hypothesis spaces without a maximally specific consistent hypothesis.



$x_1 = \langle \text{sunny, warm, normal, strong, warm, same} \rangle, +$

$x_2 = \langle \text{sunny, warm, high, strong, warm, same} \rangle, +$

$x_3 = \langle \text{rainy, cold, high, strong, warm, change} \rangle, -$

$x_4 = \langle \text{sunny, warm, high, strong, cool, change} \rangle, +$

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

$$h_1 = \langle \text{sunny, warm, normal, strong, warm, same} \rangle$$

$$h_2 = \langle \text{sunny, warm? strong, warm, same} \rangle$$

$$h_3 = \langle \text{sunny, warm? strong, warm, same} \rangle$$

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

The hypothesis space performed by finds the search begins (to) with the most specific hypothesis in H_0 then considers increasingly general hypotheses (H_1 , through H_n) as mandated by the training examples. In the instance space diagram, positive training examples are denoted by "+", negative by "-" and instances that have not been presented as training examples are denoted by a solid circle.

Version spaces and the candidate-Elimination Algorithm:-

Concept Learning and the candidate-Elimination Algorithm :-

1. Introduction :-

- The candidate-Elimination algorithm is an approach to concept learning.
- It addresses limitations of the AND-algorithm by maintaining a set of consistent hypotheses.

2. Representation

- A hypothesis h is consistent with training examples if it correctly classifies them.
- A hypothesis h is consistent with a set of training examples D if and only if $h(x) = c(x)$ for each example $(x, c(x))$ in D .
- The version space, denoted $V \subseteq H_D$, contains all plausible versions of the target concept.

3. The List-Then-Eliminate Algorithm:

- * It is a simple learning algorithm that lists all hypotheses in H and eliminates inconsistent ones when processing training examples.
- * However, it requires enumerating all hypotheses in H , making it unrealistic for non-trivial hypothesis space.

4. A more compact Representation for version spaces

- * The Candidate-Elimination algorithm represents the version space with the most general (h) and most specific (s) hypothesis.
- * h and s delimit the version space within the partially ordered hypothesis space.
- * The version space includes hypotheses that are more general or more specific than h and s , forming a compact representation.

5. Version Space Representation Theorem

- * The version space is the set of hypotheses contained in h , plus those contained in s , plus those that lie between h and s in the hypothesis space.

6. The candidate-elimination learning algorithm

- * It initializes h and s with the most general and most specific hypotheses, respectively.
- * It processes training examples and updates h and s based on the consistency of hypotheses.
- * Positive examples lead to the generalization of s , while negative examples lead to the specialization of h .
- * The algorithm terminates with a version space containing consistent hypotheses.

7. An Illustrative Example :-

- * Example trace of the candidate - Elimination algorithm applied to training examples.
- * Positive examples lead to generalization of the s boundary, and negative examples lead to the specialization of the

6 boundary.

- * The S boundary summarizes part Positive examples, and the G boundary summarizes part Negative examples.
- * The algorithm converges to a version space containing all consistent hypotheses.

Example - 1 :- Training with positive Examples :-

Initial State:

• Most general hypothesis: $h = (?, ?, ?, ?, ?)$

• Most specific hypothesis: $h = (<\text{sunny}, \text{warm}, \text{normal}, \text{strong}, \text{warm}, \text{same}>)$

• Training Example: $<\text{sunny}, \text{warm}, \text{normal}, \text{strong}, \text{warm}, \text{same}>$

Action:-

• The training example is positive.

• The S boundary is revised to be more

general: $S = \{ <\text{sunny}, \text{warm}, ?, \text{strong}, \text{warm}, \text{same}> \}$

• No change in the G boundary

Example-2 :- Training with a Negative Example.

Example.

- Current state:-
- S boundary: ' $s = \{ <\text{sunny}, \text{warm}, ?, \text{strong}, ?, ? > \}$ '
- G boundary: ' $g = \{ < ?, ?, ?, ?, ?, ? > \}$ '
- Training Example: ' $\langle \text{Rainy}, \text{cold}, \text{high}, \text{strong}, \text{warm}, \text{change} \rangle$ '
- Action:
 - The training example is negative.
 - The G boundary is revised to be more specific:-
- New hypothesis in G: ' $g = \{ < \text{sunny}, \text{warm}, ?, \text{strong}, ?, ?, ? >, < ?, \text{warm}, ?, ?, ?, ? > \}$ '
- No change in the S boundary.

Remarks on Version Spaces and Candidate-Elimination

1. Will the Candidate-Elimination Algorithm converge to the correct hypothesis?

- The Candidate-Elimination algorithm will converge to the correct hypothesis if there are no errors in the training examples and if there is at least one hypothesis in the hypothesis space (H) that correctly describes the target concept.

• The version space can be monitored to determine when sufficient training examples have been observed to unambiguously identify the target concept:

• Convergence to the correct hypothesis is achieved when the S and G boundary sets converge to a single hypothesis that is identical

- The version space can be monitored
2. What Training Example should the learner Request Next?

- The learner can choose to request new training examples from an external oracle.
- An optimal query strategy is to generate instances that satisfy exactly half of the hypotheses in the current version space.
- The goal to evenly split the candidate hypotheses into sets that predict "yes" and "no" to determine the correct hypothesis.
- This strategy minimizes the number of required experiments to identify the target concept.

How can partially learned concepts be used?

- Even if the version space contains multiple hypotheses, the learner can classify certain instances with confidence.
- If an instance is classified as positive by every hypothesis in the version space, it can be classified a positive with confidence.
- If an instances is classified as negative by every hypothesis in the version space, it can be classified as negative with confidence.
- If an instance is classified differently by different hypothesis in the version space, further training examples are needed to make a confident classification.

Inductive Bias

1. A Biased Hypothesis Space :-

- if the target concept is not contained in the initial hypothesis space (H), the candidate-elimination algorithm cannot converge to the correct hypothesis.

- one way to address this is to enrich the hypothesis space to include every possible hypothesis, but this may not be practical.

2. An Unbiased Learner :-

- An unbiased learner would have a hypothesis space (H') capable of representing every possible subset of instances.

- The power set of X , containing all possible subsets, provides an unbiased

by hypothesis space

In this highly expressive hypothesis space, the learner struggles to generalize beyond observed training examples, requiring all instances to be presented as training data.

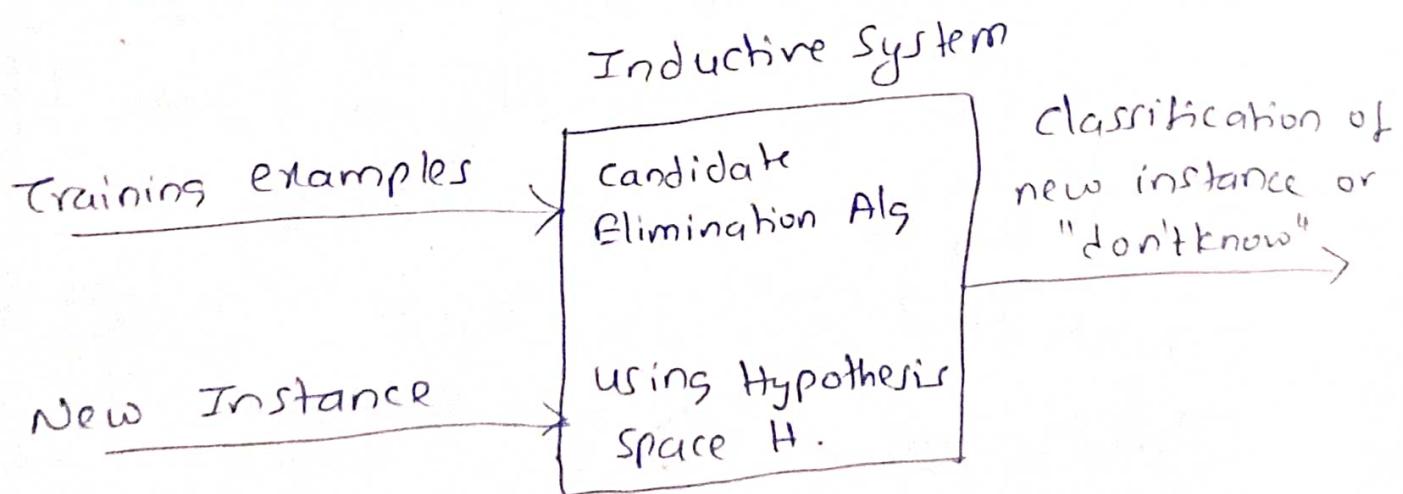
3. The Futility of Bias-free Learning -

- An unbiased learner makes no a priori assumptions about the target concept and cannot generalize to unseen instances without extensive training data.
- The Candidate-Elimination algorithm, as originally defined, has a bias that the target concept can be represented in its hypothesis space (inductive bias).
- The concept learning algorithms' inductive bias is the set of additional assumptions required for justifying

it's inductive inference as deduction.

Interface:

- By characterizing inductive systems by their inductive bias, we can improve and compare different learners based on the strengths of their biases.
- Examples of learners with varying levels of bias include ROFE-LEARNER, CANDIDATE-ELIMINATION Algorithm, and FIND-S, each making different inductive leaps and classifications



Equivalent deductive system

Training Examples

New Instance

Assertion "H contains
the target concept"

Theorem Prover

Classification of
new instance or
"don't know"



Inductive Bias
made explicit

Decision Tree Learning :-

Introduction:-

- Decision tree learning is a method for approximating discrete-valued target functions.
- g_t represents the learned function as a decision tree.
- These methods are widely used for various tasks, from medical diagnosis to credit risk assessment.

Decision Tree Representation:-

- Decision trees classify instances by traversing from the root to a leaf node.
- Each node tests an attribute of the instance.
- Branches correspond to possible attribute values.
- Instances are classified by following the path based on attribute tests.
- Decision trees represent a disjunction of conjunctions of attribute constraints
- Each path from root to leaf represent a conjunction of attribute tests.
- The entire tree is a disjunction of these conjunctions.

Example:

- Figure 3.1 shows a decision tree for classifying saturday morning for playing tennis.
- An instance is classified by moving down the tree based on attribute tests.
- The tree represents disjunctions and conjunctions of attribute constraints, allowing it to capture complex decision boundaries in the data.

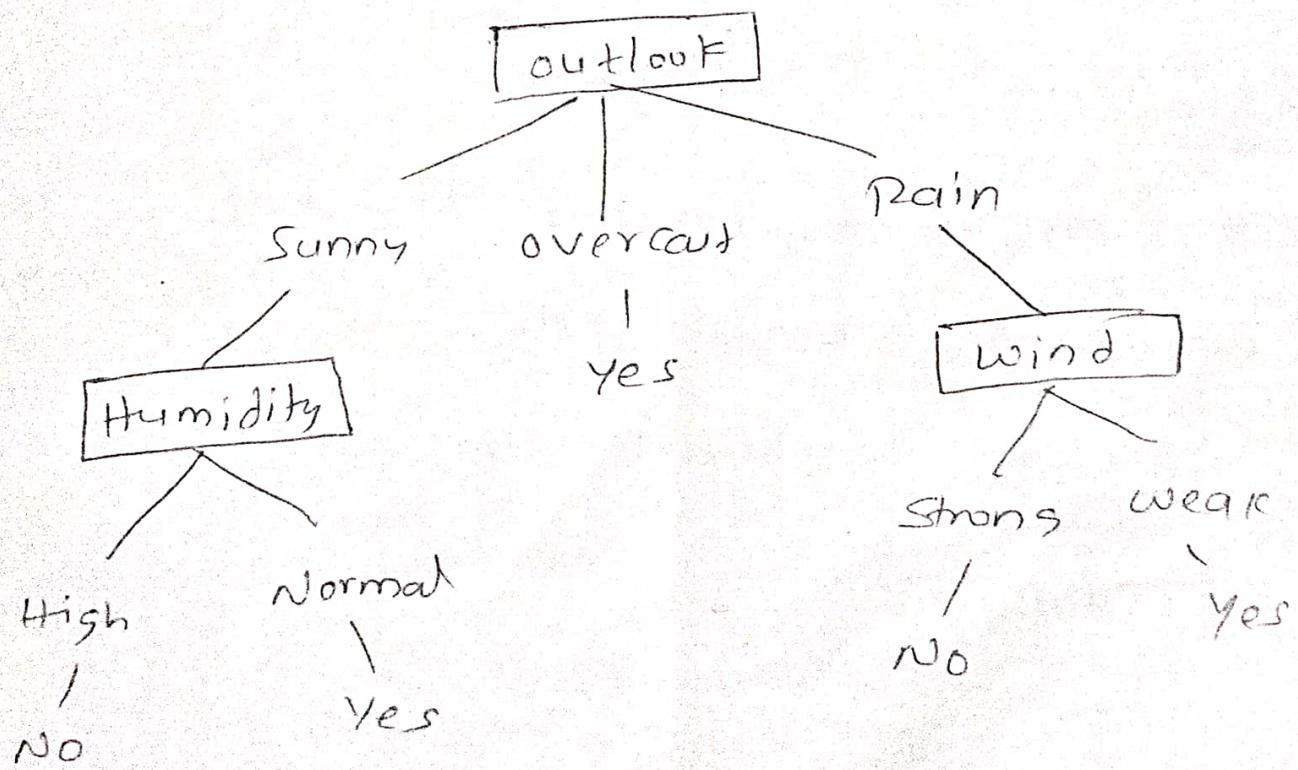


figure 3.1

How would we represent?

- \wedge, \vee, XOR
- $(A \wedge B) \vee (C \wedge \neg D \wedge E)$
- M of N

Decision Trees.

When to consider

- Instances describable by attribute value pairs.
- Target function is discrete valued
- Disjunctive hypothesis may be required
- Possibly noisy training data.

Examples:-

- Equipment or medical diagnosis.
- Credit risk analysis.
- Modeling calendar scheduling preferences.

for Figure 3.1 the representation rule
in DT's :-

if outlook = sunny AND humidity =
normal

OR

if outlook = overcast

OR

if outlook = rain AND wind = weak

then play tennis.

Rules can represent a decision tree:-

if item₁ then subtree 1

else item₂ then subtree 2

elseif ---

Advantages of rules over DT:-

Rules are easier to understand,

modify and combine.

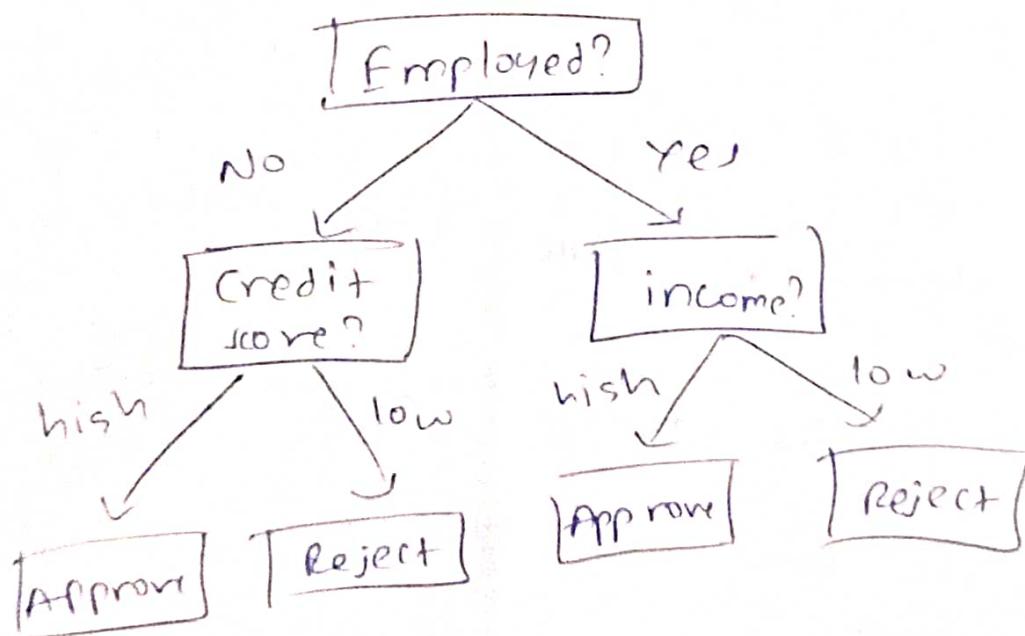
Rules can significantly improve
classification performance by eliminating
unnecessary tests.

- Rules make it possible to combine different decision trees for the same task.

Advantages of Decision Tree

- Decision trees are simple to understand.
- Decision trees have a clear evaluation strategy.
- Decision trees are able to handle both nominal and categorical data.

Example: whether to approve a loan

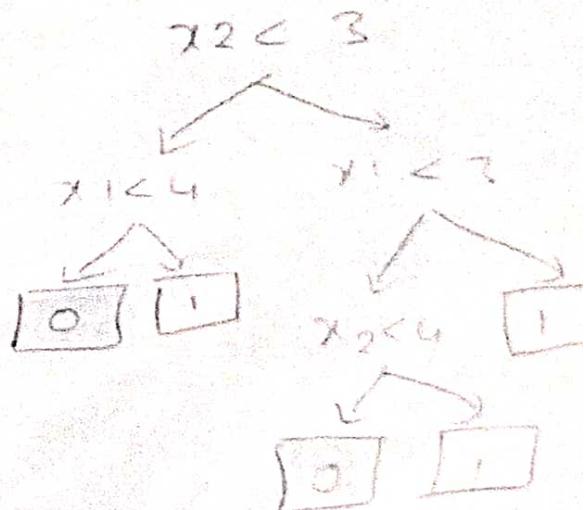
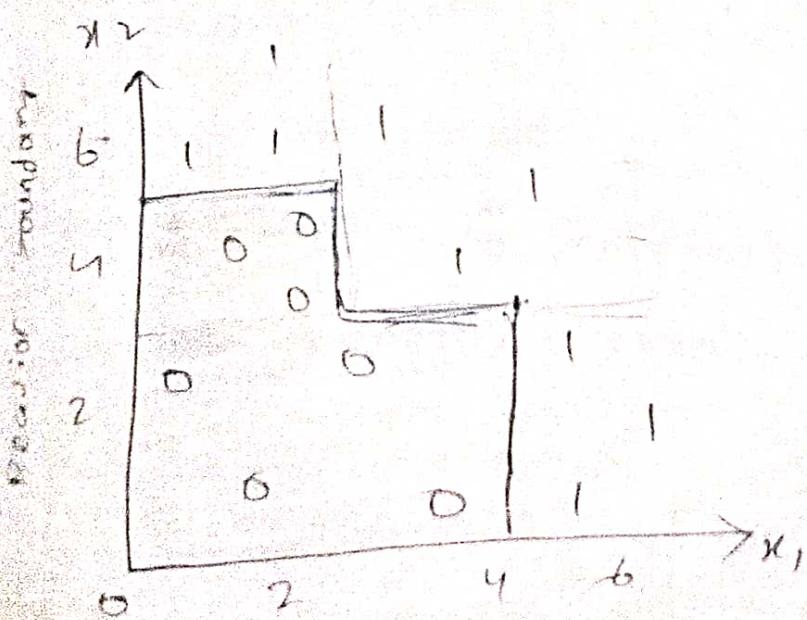


Applications

- Instances do not have binary attribute values.
- target function is discrete valued.
- Disjunctive hypothesis may be required.
- possibly noisy training data.

Decision Tree - Decision Boundary

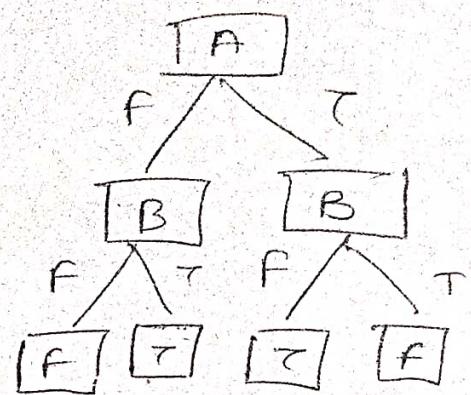
- Decision trees divide the features space into axis-parallel rectangles.
- Each rectangle region is labeled with one label - or a probability distribution over labels.



Expressiveness:-

- Decision trees can represent any function of the input attributes
- Boolean operations (and, or, xor, etc)
 - All Boolean functions.

A	B	$A \text{ xor } B$
F	F	F
F	T	T
T	F	T
T	T	F



Searching for a good tree :-

- The space of decision trees is too big for systematic search.
- Stop and
 - return the a value for the target feature or
 - a distribution over target feature values

- choose a test (e.g. an input feature) to split on.
- for each value of the test, build a subtree for those examples with this value for the test.

The Basic ID3 Algorithm :-

ID3 (Examples, Target-attribute, Attributes)

- * Examples are training examples.
 - * Target-attribute is the attribute whose value to be predicted by the tree.
 - * Attributes is the list of other attributes that may be tested by the learned decision tree.
- create a Root node for the tree.
- if all Examples are positive, Return the single-node tree Root, with label = +
- if all Examples are negative, Return the single-node tree Root with label = -

→ if Attribute is empty, return the single-node tree Root, with label = most common value of target-attribute in Examples.

→ otherwise begin

- $A \leftarrow$ the attribute from Attribute that best * classifies Examples.

- The decision attribute for Root $\leftarrow A$.

- for each possible value, v_i , of A

→ Add a new tree branch below

Root, corresponding to the test $A = v_i$.

→ The decision attribute for $\text{Root} \leftarrow A$.

→ for each

→ let Examples _{v_i} be the subset of Examples that have v_i for A .

→ if Examples _{v_i} is empty,

- Then below this new branch add a leaf node with label

= most common value of target-attribute in Examples

- Else below this new branch add the subtree

ID3 (Example v., Target-attribute,
Attributes = {A})

- End
- Return Root

Entropy :-

- A measure for
 - Uncertainty
 - Purity
 - information content

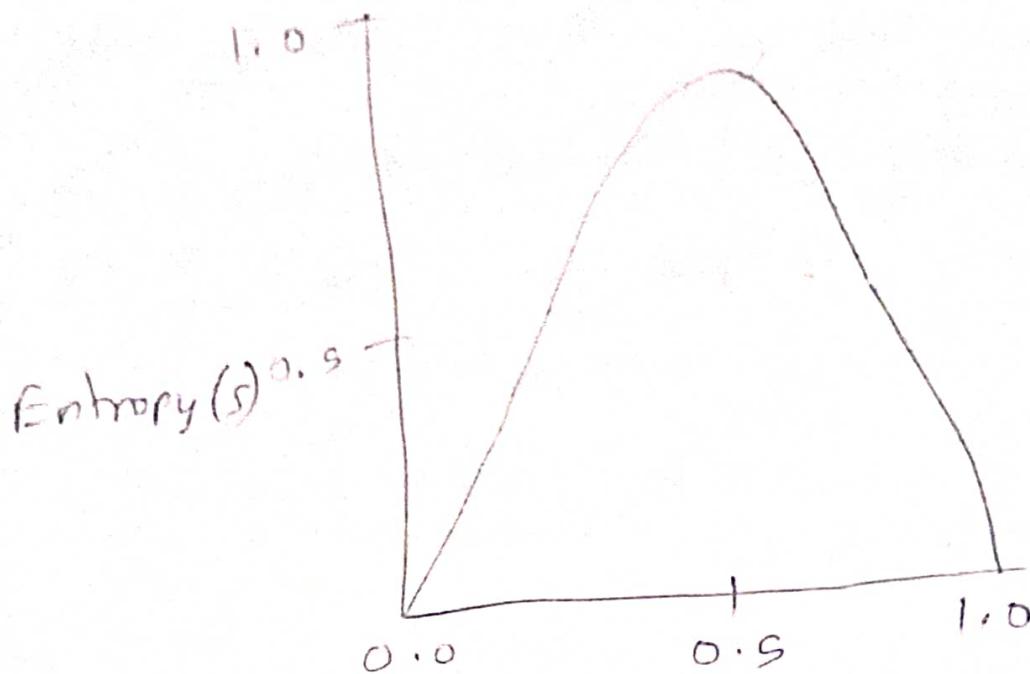
- The entropy is 0 if the outcome is "certain".
- The entropy is maximum if we have no knowledge of the system.

$$\text{Entropy} = \sum_{i=1}^c -p_i * \log_2(p_i)$$

$$\text{infoGain} = \frac{\text{Entropy}(\text{Parent}) - (\text{Weighted Avg})}{\text{Entropy}(\text{Children})}$$

$\text{Gain}(S, A) = \text{expected reduction in entropy due to sorting on } A$

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



$$P_+ \oplus \textcircled{+}$$

- Entropy characterizes then (im) purity of an arbitrary collection of examples.
- Given a collection S (containing positive and negative examples), the entropy of S relative to this boolean classification is.
- P_+/P_- is the proportion of positive/negative examples in S .
- Entropy is 0 if all members of S belongs to the same class; entropy is 1 when the collection

contains an equal number of positive and negative examples.

Information Gain :-

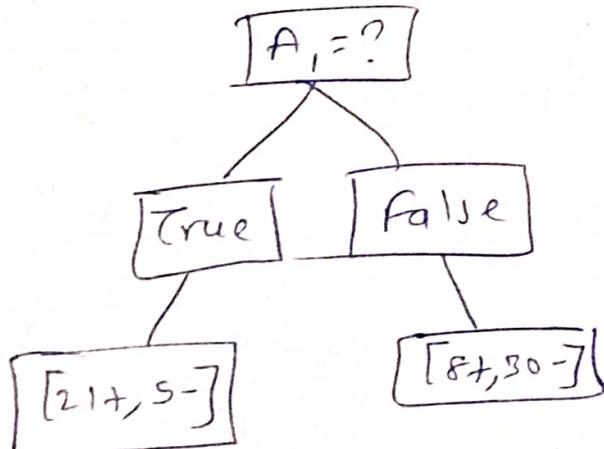
- Information gain measures the expected reduction in entropy caused by partitioning the examples according to an attribute.
- first term:- entropy of the original collection s ; Second term: expected value of entropy after s is partitioned using attribute A (s_A subset of s).
- $\text{Gain}(s, A)$: The expected reduction in entropy caused by knowing the value of attribute A.
- ID3 uses information gain to select the best attribute at each step in growing the tree.

$\text{Gain}(S, A)$: expected reduction in entropy due to partitioning S on attribute A .

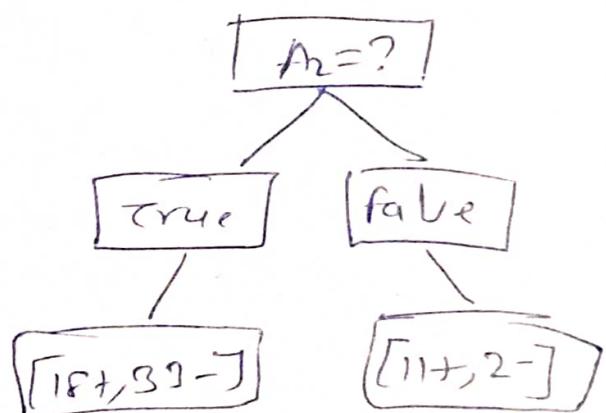
$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\begin{aligned} \text{Entropy}([29+, 35-]) &= -\frac{29}{64} \log_2 \frac{29}{64} \\ &\quad - \frac{35}{64} \log_2 \frac{35}{64} \\ &= 0.99 \end{aligned}$$

$[29+, 35-]$



$[29+, 35-]$



$$\text{Entropy}([21+, 5-]) = 0.71$$

$$\text{Entropy}([8+, 30-]) = 0.74$$

$$\text{Gain}(s, A_1) = \text{Entropy}(s)$$

$$-26/64 * \text{Entropy}([21+, 5-])$$

$$-38/64 * \text{Entropy}([8+, 30-])$$

$$= 0.27$$

$$\text{Entropy}([18+, 33-]) = 0.94$$

$$\text{Entropy}([8+, 30-]) = 0.62$$

$$\text{Gain}(s, A_2) = \text{Entropy}(s)$$

$$-51/64 * \text{Entropy}([18+, 33-])$$

$$-13/64 * \text{Entropy}([11+, 2-])$$

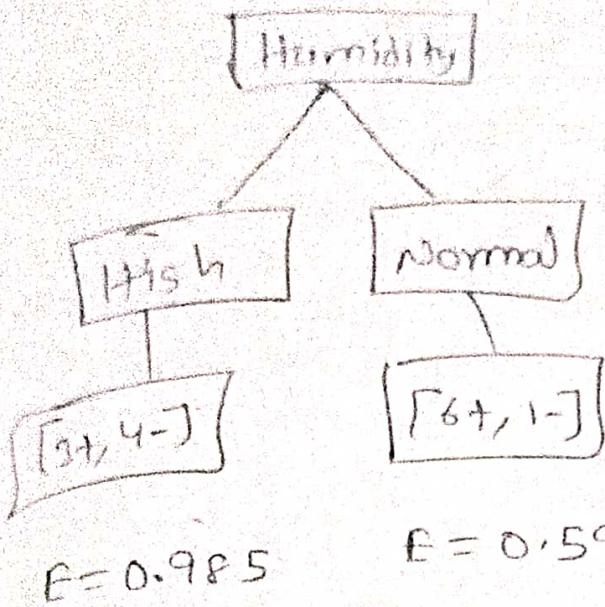
$$= 0.12$$

Training Examples

Day	outlook	Temp	Humidity	Wind	Terrain
D ₁	Sunny	Hot	High	weak	No
D ₂	Sunny	Hot	High	Strong	No
D ₃	overcast	Hot	High	Weak	Yes
D ₄	Rain	Mild	High	weak	Yes
D ₅	Rain	Cool	Normal	weak	Yes
D ₆	Rain	Cool	Normal	Strong	No
D ₇	overcast	Cool	Normal	Strong	Yes
D ₈	Sunny	Mild	High	weak	No
D ₉	Sunny	Cool	Normal	weak	Yes
D ₁₀	Rain	Mild	Normal	weak	Yes
D ₁₁	Sunny	Mild	High	Strong	Yes
D ₁₂	overcast	Mild	Normal	weak	Yes
D ₁₃	overcast	Hot	Normal	Strong	No
D ₁₄	Rain	Mild	High	Strong	No

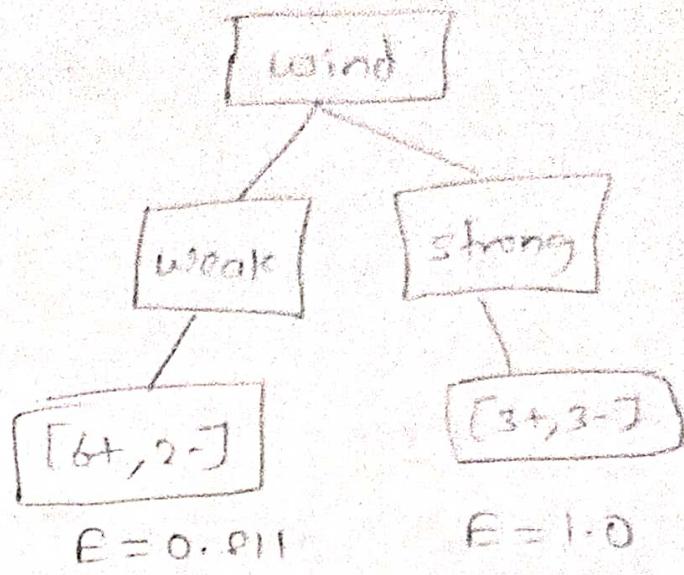
$$S = [9+, 5-]$$

$$E = 0.940$$



$$S = [4+, 5-]$$

$$E = 0.940$$



Gain (S, Humidity)

$$\begin{aligned} &= 0.940 - (7/14) * 0.985 \\ &\quad - (7/14) * 0.592 \\ &= 0.151 \end{aligned}$$

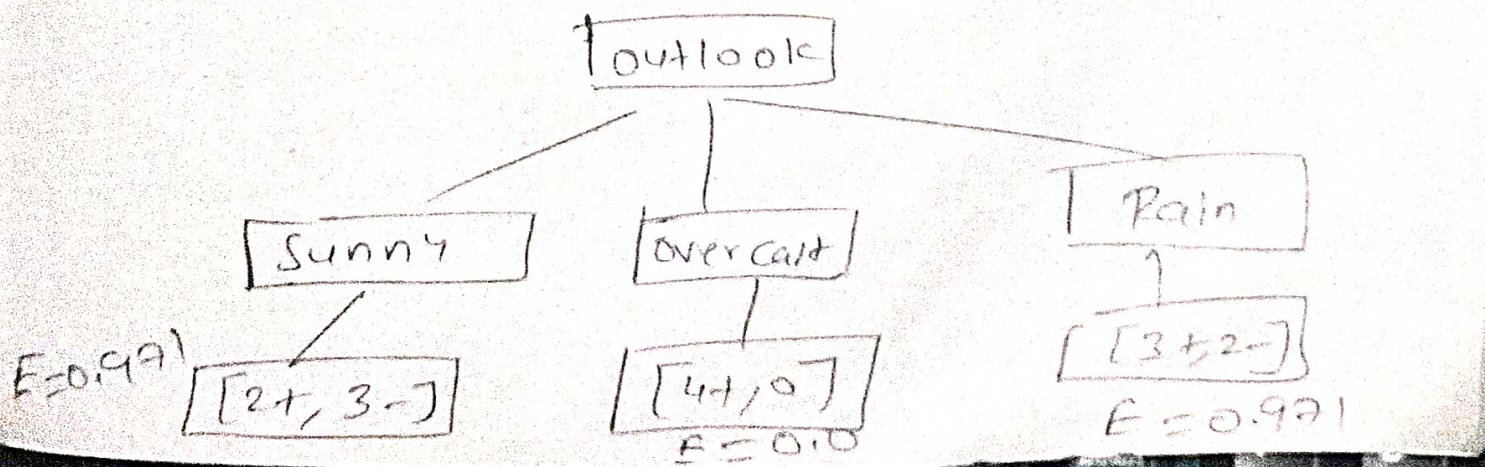
Gain (S, wind)

$$\begin{aligned} &= 0.940 - (8/14) * 0.811 \\ &\quad - (6/14) * 1.0 \\ &= 0.048. \end{aligned}$$

humidity provides greater info. gain
than wind, w.r.t target classification

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

$$E = 0.940$$



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

$$= (0.940 - (5/14)^* 0.971)$$

$$- (4/14)^* 0.0 - (5/14)^* 0.0971$$

$$= 0.247$$

The information gain value for
the 4 attributes are:

- $\text{Gain}(S, \text{outlook}) = 0.247$.

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

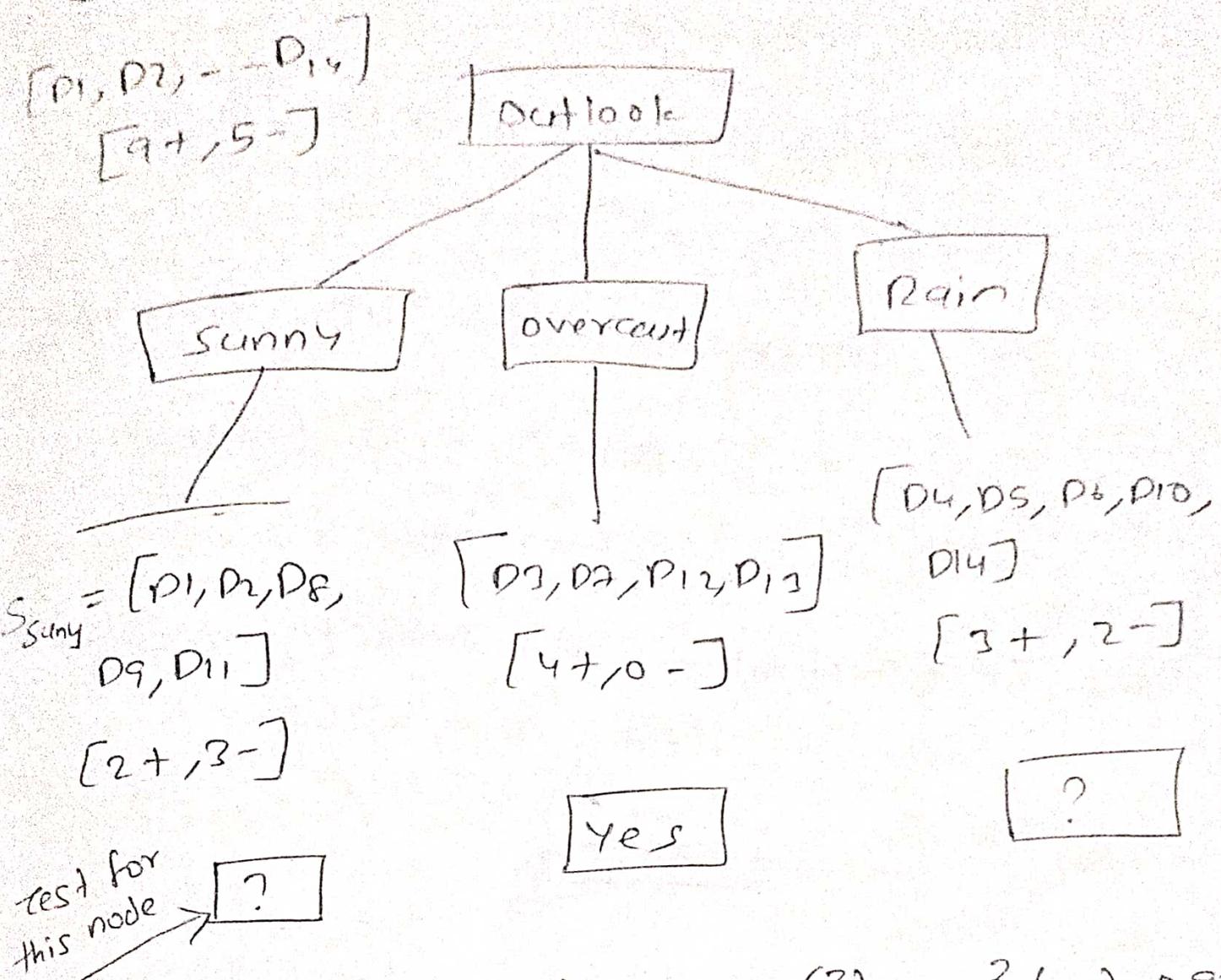
- $\text{Gain}(S, \text{wind}) = 0.048$

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

where S denotes collection of
trainings examples.

Note :- $0 \log_2 0 = 0$

ID3 - Algorithm



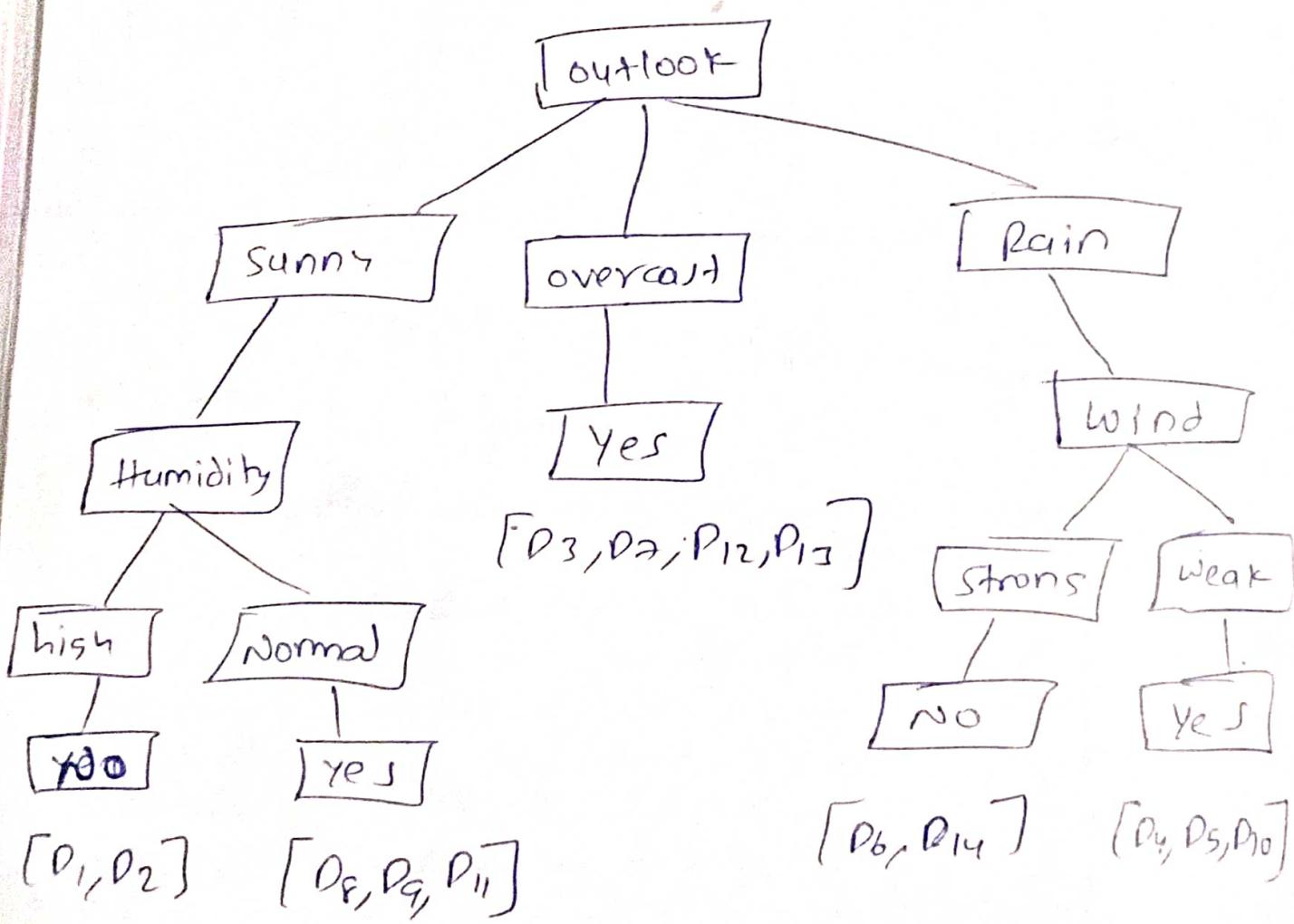
$$\text{Gain}(S_{\text{sunny}}, \text{Humidity}) = 0.970 - \left(\frac{3}{5}\right)0.0 - \frac{3}{5}(0.0) = 0.970$$

$$\text{Gain}(S_{\text{sunny}}, \text{temp}) = 0.970 - \left(\frac{2}{5}\right)0.0 - \frac{2}{5}(1.0)$$

$$- \left(\frac{1}{5}\right)0.0 = 0.570$$

$$\text{Gain}(S_{\text{sunny}}, \text{wind}) = 0.970 - \left(\frac{7}{5}\right)1.0 - \frac{7}{5}(0.918) \\ = 0.019$$

- $R_1: \text{IF}(\text{outlook} = \text{sunny}) \text{ AND } (\text{windy} = \text{false})$
 THEN play = Yes.
- $R_2: \text{IF}(\text{outlook} = \text{sunny}) \text{ AND } (\text{windy} = \text{true})$
 THEN play = No
- $R_3: \text{IF}(\text{outlook} = \text{overcast}) \text{ THEN play} = \text{Yes}$
- $R_4: \text{IF}(\text{outlook} = \text{rainy}) \text{ AND } (\text{humidity} = \text{high}) \text{ THEN play} = \text{No}$.
- $R_5: \text{IF}(\text{outlook} = \text{rain}) \text{ AND } (\text{humidity} = \text{normal}) \text{ THEN play} = \text{Yes}$.



Splittings Rule: GINI Index

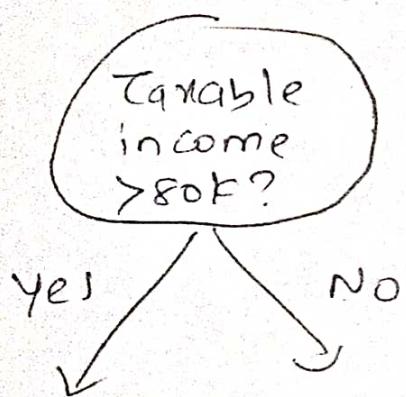
GINI index

- Measure of node impurity

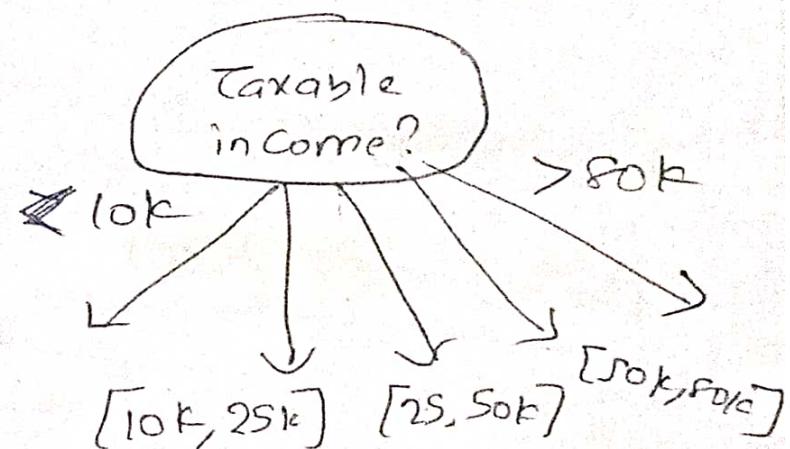
$$GINI_{node}(\text{Node}) = 1 - \sum_{c \in \text{classes}} [P(c)]^2$$

$$GINI_{split}(A) = \sum_{S \in \text{values}(A)} \frac{|S_v|}{|S|} GINI(N_v)$$

Splittings Based on continuous attributes



i) Binary split



ii) Multi-way split

Binary split

- For continuous attribute

- partition the continuous value of attribute A into a discrete set of intervals.

temperature = 82.5

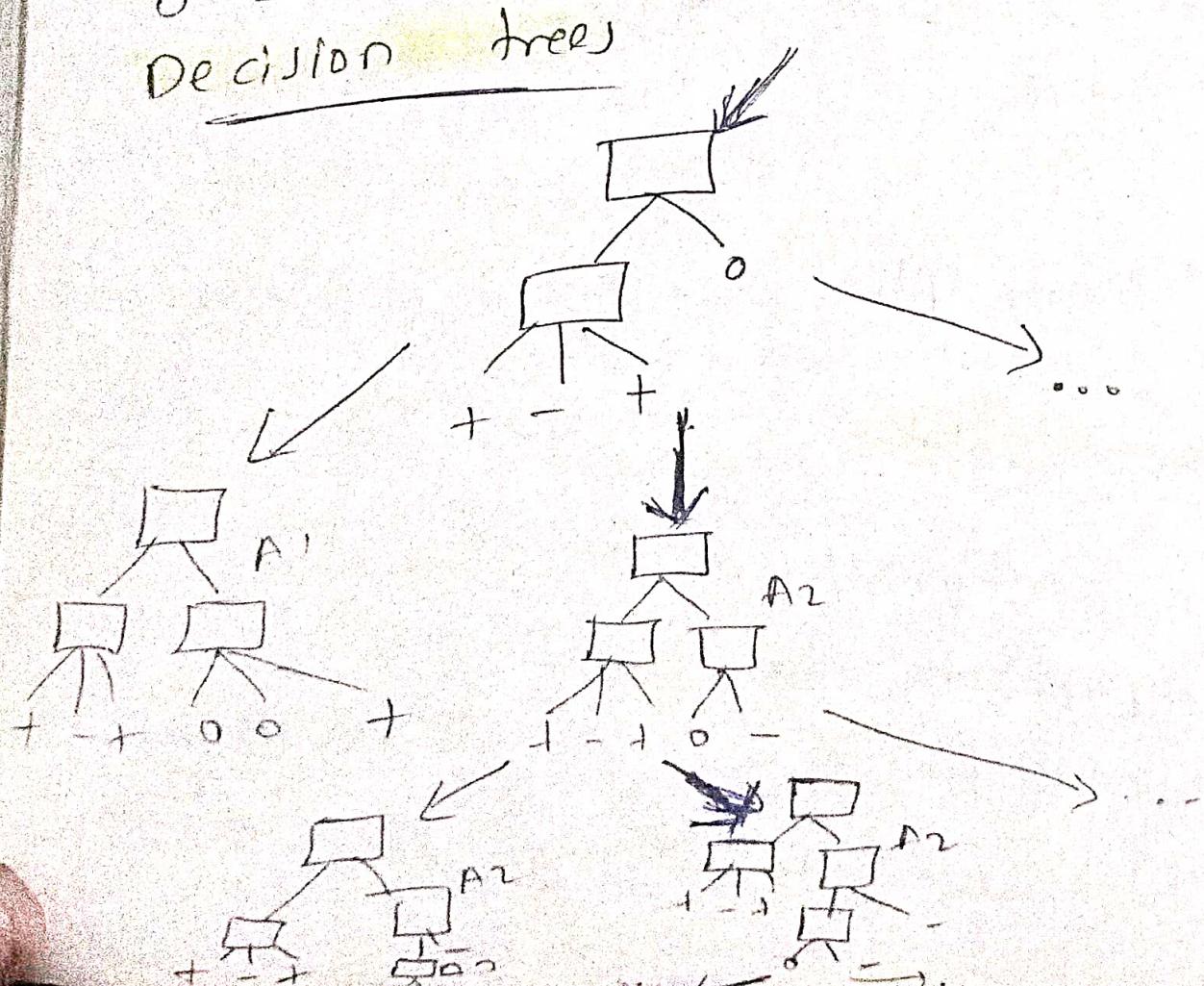
- Create a new boolean attribute A_c
looking for a threshold c ,

$$A_c = \begin{cases} \text{true} & \text{if } A_c < c \\ \text{false} & \text{otherwise} \end{cases}$$

- consider How to choose c ?
find the best cut

Temperature: 40 48 60 72 80 90
 play tennis: no no yes yes yes no

Hypothesis Space Search in
 Decision trees = 2



- Top Search a complete hypothesis space but does so incompletely since once it finds a good hypothesis it stops (cannot find others).
- Candidate-Elimination Search an incomplete hypothesis space (it can only represent some hypothesis) but does so completely.
- A preference bias is an inductive bias, where some hypothesis are preferred over others.
- A restriction bias is an inductive bias, where the set of hypothesis considered is restricted to a smaller set.

Inductive Bias in Decision Tree

Learning:-

- inductive bias is the set of assumptions that, together with the training data.
- It choose the first acceptable tree it encounters in its simple-to-complex, hill climbing search through the space of possible trees.
- The ID3 Search Strategy.
 - selects in favor of shorter trees over longer ones.
 - selects trees that place the attributes with highest information gain closest to the root.

Types of inductive bias

- preference bias or search bias
- restriction bias or language bias.

The induction bias of TDI 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 preference bias (or alternatively, a search bias).

The bias of the candidate elimination algorithm is in the form of a categorical restriction on the set of hypotheses considered. This form of bias is typically called restriction bias.
(language bias).

A preference bias is more desirable than a restriction bias.

Bias and Occam's Razor

Bias searches a complete hypothesis space incompletely, inductive bias is solely a consequence of the orderings of hypotheses by in search strategy.

Inductive bias often goes by the name of Occam's Razor

Occam's Razor :- prefer shorter hypotheses that fits the data

Issues in Decision Tree learning

- Determining how deeply to grow the decision tree, underfitting and overfitting.
- Handling continuous attributes.
- Choosing an appropriate attribute Selection measure.

- Handlings training data with missing attribute values.
- Handlings attributes with different costs.
- Improving computational efficiency.

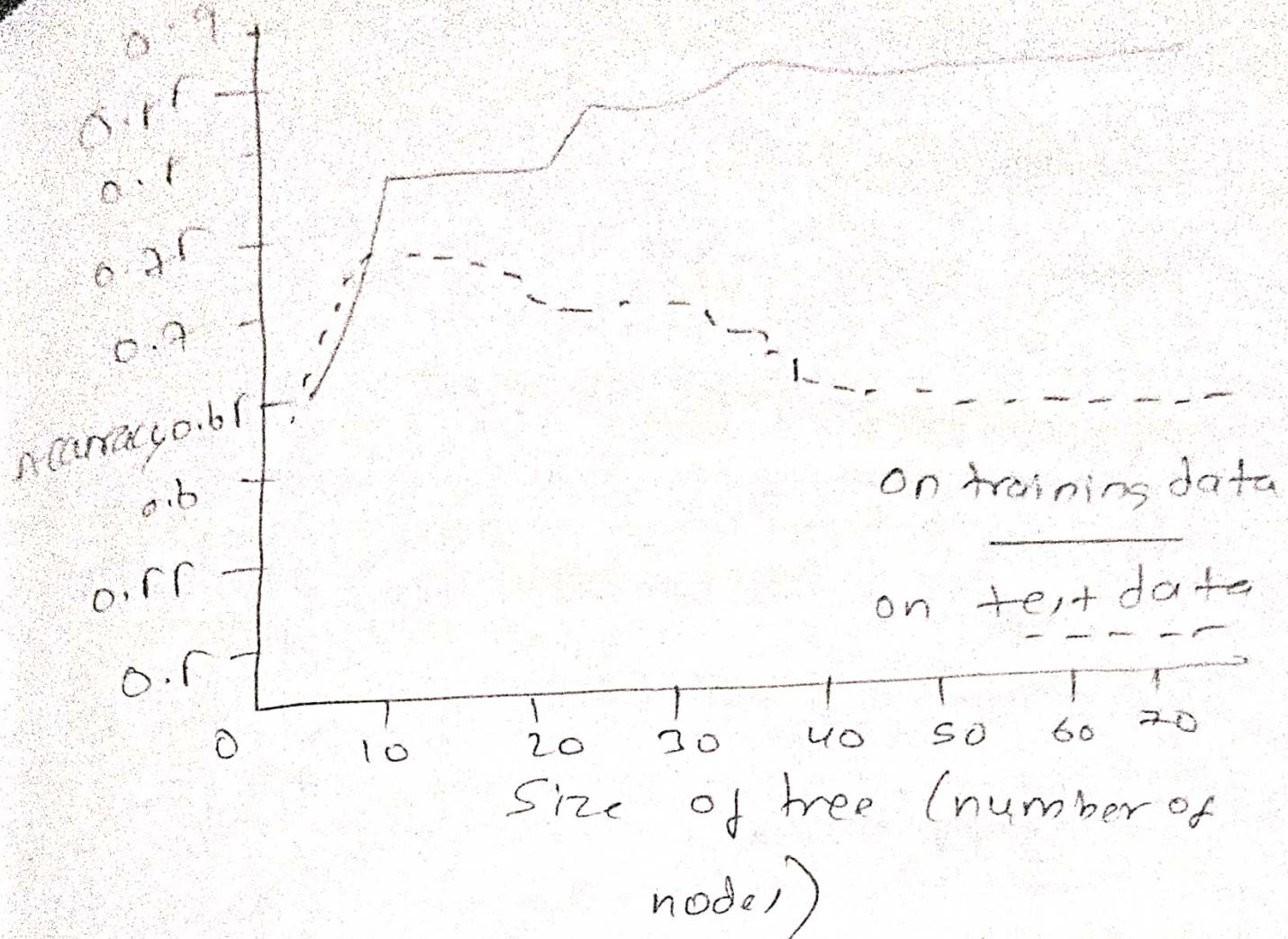
Overfitting :-

An algorithm can produce trees that overfit the training examples in the following two cases.

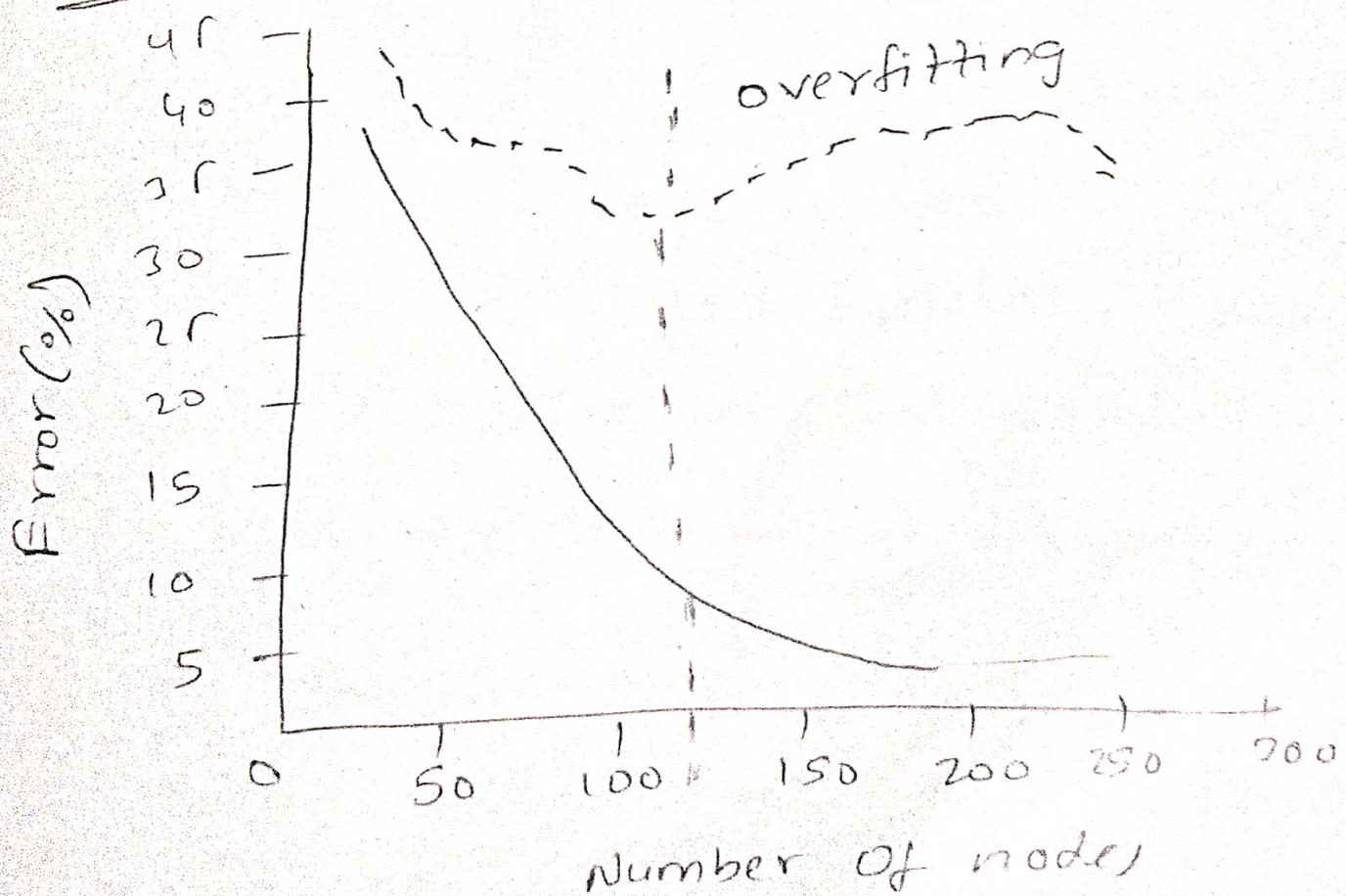
- There is noise in the data.
- The number of training examples is too small to produce a representative sample of the true target function.
- Learning a tree that classifies the training data perfectly, may not lead to the tree with the best generalization performance.

- There may be noise in the training data.
- May be based on insufficient data.
- A hypothesis, h is said to over-fit the training data if there is another hypothesis, h' , such that h has smaller error than h' on the training data but h' has larger error on the test data than h .



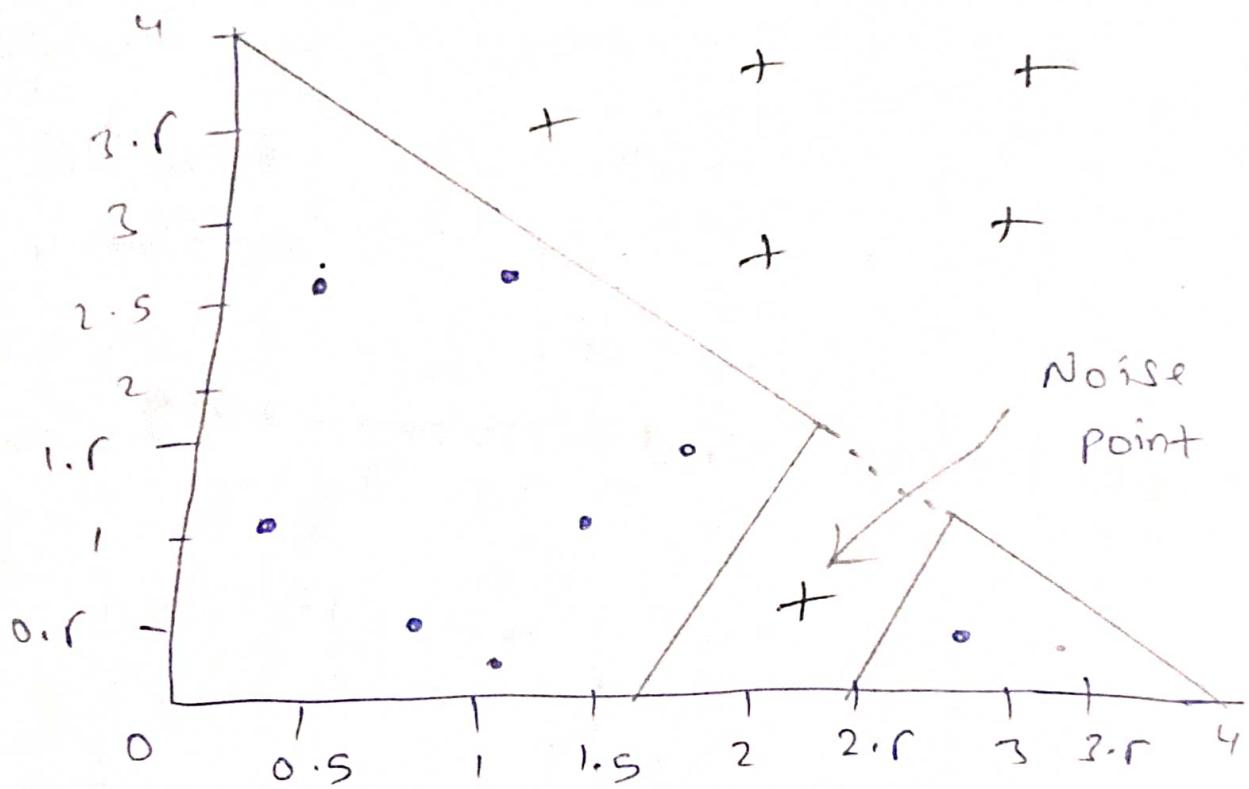


Underfitting and overfitting



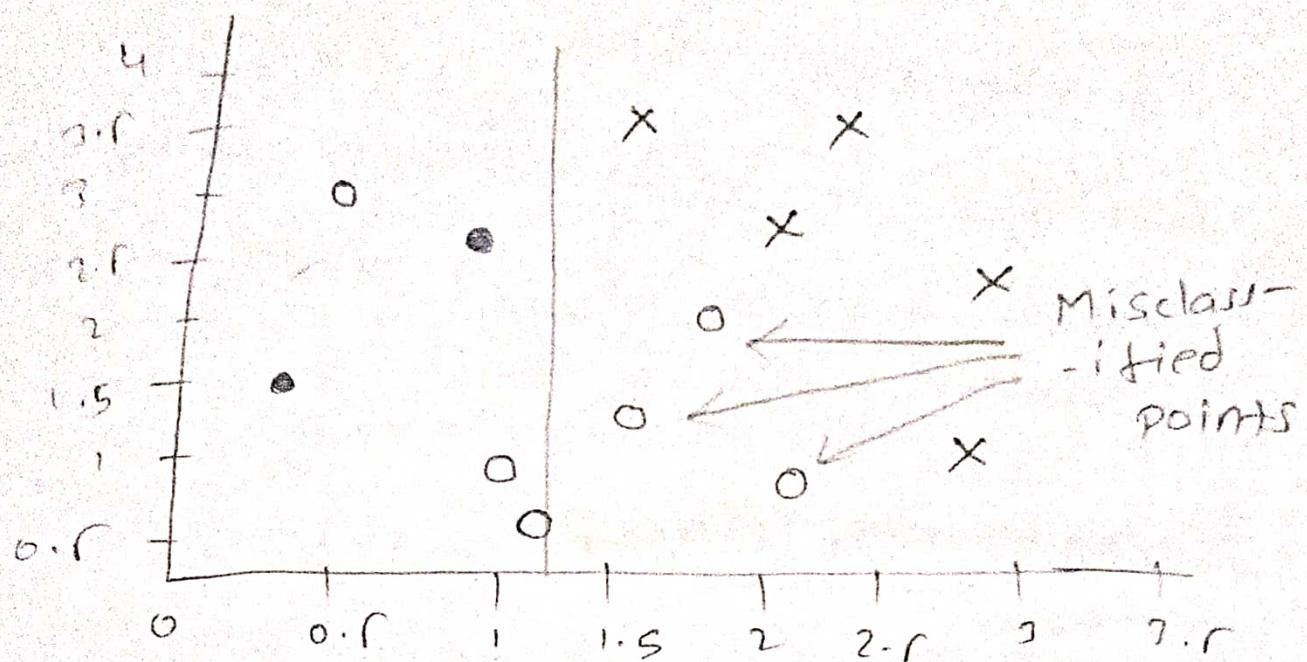
Underfitting :- When model is too simple, both training and test errors are large.

Overttting due to noise :-



Decision boundary is distorted
by noise point

overfitting due to insufficient examples



Lack of data points make it difficult to predict correctly the class labels of the region.

- * Overfitting results in decision trees that are more complex than necessary.
- * Training error no longer provides a good estimate of how well the tree will perform on previously unseen records.
- * Overfitting happens when a model is capturing idiosyncrasies of the data rather than generalities.

* often caused by too many parameters relative to the amount of training data.

* E.g - an order- n polynomial can intersect any $N+1$ data points.

Avoid overfitting:-

* There are several approaches to avoiding overfitting in decision tree learning. These can be grouped into two classes.

- * pruning :- Stop growing when data split not statistically significant
- * postpruning :- Grow full tree then remove nodes

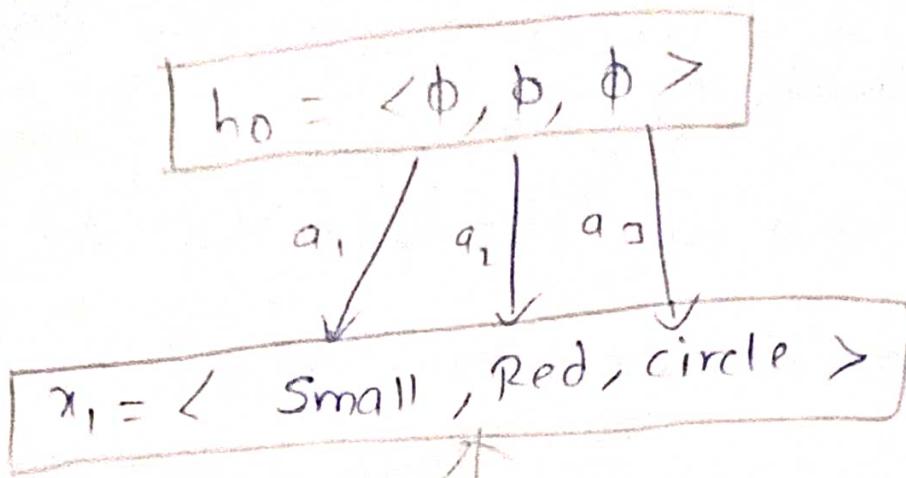
- methods for evaluating subtree to prune:
 - Minimum description length(MDL)
minimize: size(tree) + size
(misclassifications(tree))
 - measure performance over training data.
 - measure performance over separate validation data set

Questions:- (Unit-1)

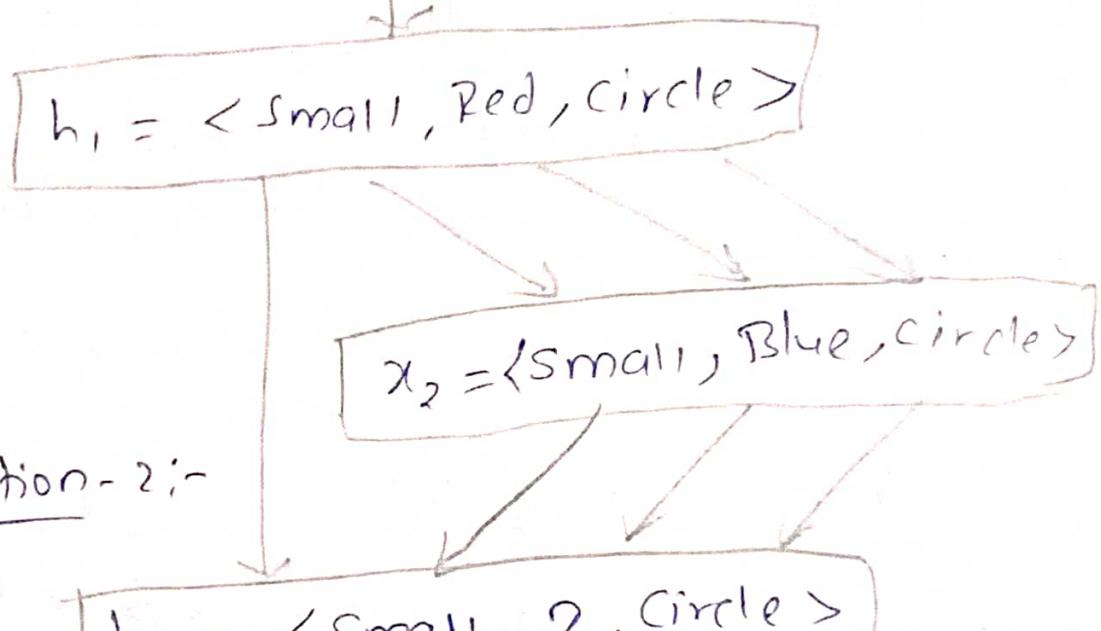
- Q) Find-S algorithm and find maximally specific hypothesis for the given training examples.

size	color	shape	class
Big	Red	circle	No
Small	Red	triangle	No
Small	Red	circle	Yes
Big	Blue	circle	No
Small	Blue	circle	Yes

Sol:-



iteration-1:-



iteration-2:-

Output:-

Note:- When there is No in logical and output values then no need to consider that Just check for only which has 'Yes'.

d) find-s Algorithm

example	citations	size	inlibrary	Price
1	some	small	no	affordable
2	many	big	no	expensive
3	some	big	always	expensive
4	many	medium	no	expensive
5	many	small	no	affordable
editions	buy			
many		no		
one		yes		
few		no		
many		yes		
many		yes		

a) How many concepts are possible for the instance space?

b) How many hypotheses can be expressed by the hypothesis language

c) Apply the find-s algorithm by hand on the given training set.

Consider the examples in the specified order and write down your hypothesis each time after observing an example.

$$9) \Rightarrow 2 * 3 * 2 * 2 * 3$$

$\Rightarrow 72$ concepts.

b)

by adding null and ? in each attribute then the total becomes

$$\Rightarrow 4 * 5 * 4 * 4 * 5$$

$\Rightarrow 1600$

Symmetrically Distinct Hypothesis

not consider null.
only ?

$$= (3 * 4 * 3 * 3 * 4) + 1 = 437$$

↑
all nulls
as 1

c) Step-1 :-

$$h_0 = (\phi, \phi, \phi, \phi, \phi)$$

Step-2 :-

$$x_1 = (\text{some, small, no, expensive, many})$$

- No

Negative Example Hence ignore

$$h_1 = (\phi, \phi, \phi, \phi, \phi)$$

$h_1 = (\phi, \phi, \phi, \phi, \phi)$

$x_2 = (\text{many}, \text{big}, \text{no}, \text{expensive}, \text{one})$ - Yes

$h_2 = (\text{many}, \text{big}, \text{no}, \text{expensive}, \text{one})$

$x_3 = (\text{some}, \text{big}, \text{always}, \text{expensive}, \text{few})$
- No

Negative example hence ignore

$h_3 = (\text{many}, \text{big}, \text{no}, \text{expensive}, \text{one})$

$x_4 = (\text{many}, \text{medium}, \text{no}, \text{expensive}, \text{many})$
- Yes

$h_4 = (\text{many}, ?, \text{no}, \text{expensive}, ?)$

$x_5 = (\text{many}, \text{small}, \text{no}, \text{affordable}, \text{many})$
- Yes

$\boxed{h_5 = (\text{many}, ?, \text{no}, ?, ?)}$

Final hypothesis

	Ex	Eyes	nose	head	FColor
1	Round	triangular	round	round	Purple
2	Square	Square	Square	square	Green
3	Square	Triangular	round	round	Yellow
4	Round	triangular	round	round	Green
5	Square	Square	round	round	Yellow

Hair Smile

yes	yes
yes	no
yes	yes
no	no
yes	yes

Do it yourself

(P) candidate algorithm

Ex	Sky	AirTemp	Humidity	Wind	Water
1	sunny	warm	normal	strong	warm
2	sunny	warm	high	strong	warm
3	Rainy	cold	high	strong	warm
4	sunny	warm	high	strong	cool

Forecast Enjoy Sport

same	yes
same	yes
change	no
change	yes

so:

$(\phi, \phi, \phi, \phi, \phi, \phi)$

$S_1 :$ (sunny, warm, normal, strong, warm, same)

$S_2 : S_3 :$ (sunny, warm, ?, strong, warm, same)

$\rightarrow S_4 :$ (sunny, warm, ?, strong, ?, ?)

$G_4 :$

(sunny, ?, ?, ?, ?, ?)

(?, warm, ?, ?, ?, ?)

(?, ?, ?, ?, ?, same)

(?, ?, ?, ?, cool, ?)

$G_3 :$ (sunny, ?, ?, ?, ?, ?)

(?, warm, ?, ?, ?)

(?, norm, ?, ?, ?)

$G_0 : G_1 : G_2 :$

(?, ?, ?, ?, ?, ?, ?)

S

(sunny, warm, ?, strong, ?)

(sunny, ?, ?, strong, ?)

(sunny, warm, ?, ?, ?)

(?, warm, ?, strong, ?)

G

(sunny, ?, ?, ?, ?)

(?, warm, ?, ?, ?)

Q) Candidate Elimination Algs

size	color	shape	label
Bis	Red	circle	No
Small	Red	triangle	No
Small	Red	circle	Yes
Bis	Blue	circle	Yes, No
Small	Blue	circle	Yes

Sol :- S0 : ($\emptyset, \emptyset, \emptyset$)

S1 : ($\emptyset, \emptyset, ?$)

S2 : ($\emptyset, \emptyset, ?$)

S3 : (Small, Red, circle)

S4 : (Small, Red, circle)

G1 : (Small, ?, circle)

G2 : (Small, ?, circle)

(Bis, ?, triangle) (? , Blue, triangle)

G3 : (Small, Blue, ?) (Small, ?, circle) (? , Blue, ?)

G4 : (Small, ?, ?) (? , Blue, ?) (? , ?, triangle)

G5 : (?, ?, ?)

55 : (small, ?, circle)

65 : (small, ?, circle)

5:6: (small, ?, circle)

It is perfectly learned
after

φ)

Example	citations	size	in library
1	some	small	No
2	many	big	No
3	many	medium	No
4	many	small	No

price	Editions	buy
Affordable	one	No
Affordable	many	Yes
Expensive	few	Yes
Expensive	many	Yes
Affordable		

Q)

Ex	Eyes	Nose	Head	Color	Hair
1	Round	triangle	Round	purple	Yes
2	Square	Square	Square	Green	Yes
3	Square	triangle	Round	Yellow	Yes
4	Round	triangle	Round	Green	No
5	Square	Square	Round	Yellow	Yes

Smile

Yes

No

Yes

No

Yes.

Q) Decision tree for Boolean functions.

a) $A \wedge \neg B$

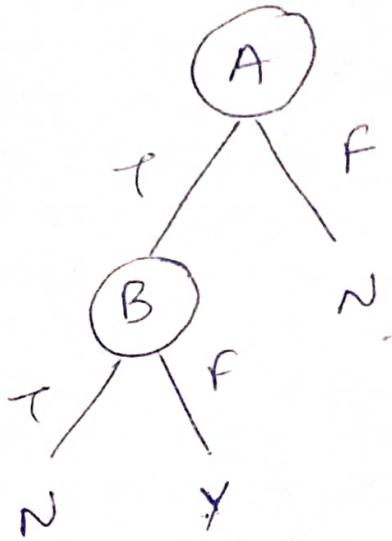
b) $A \vee [B \wedge C]$

c) $A \oplus B$

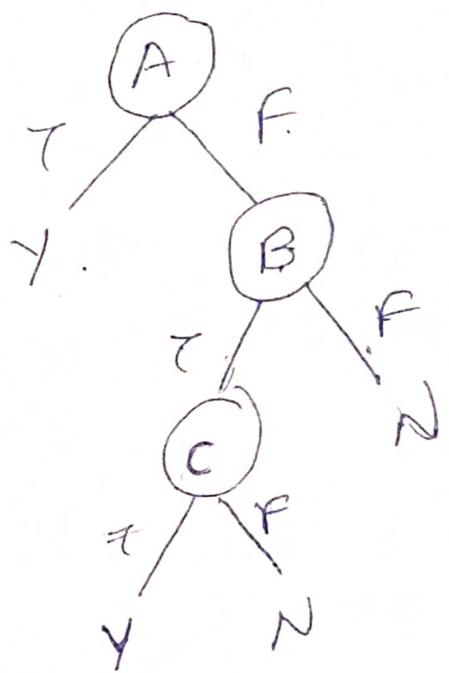
d) $[A \wedge B] \vee [C \wedge D]$

a) $A \wedge \neg B$

$\text{N} \rightarrow \text{No}$
 $\text{Y} \rightarrow \text{Yes}$

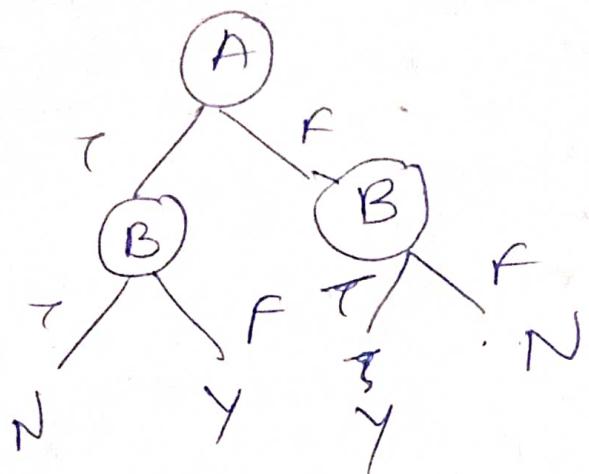


b) $\underline{A} \vee [B \wedge C]$

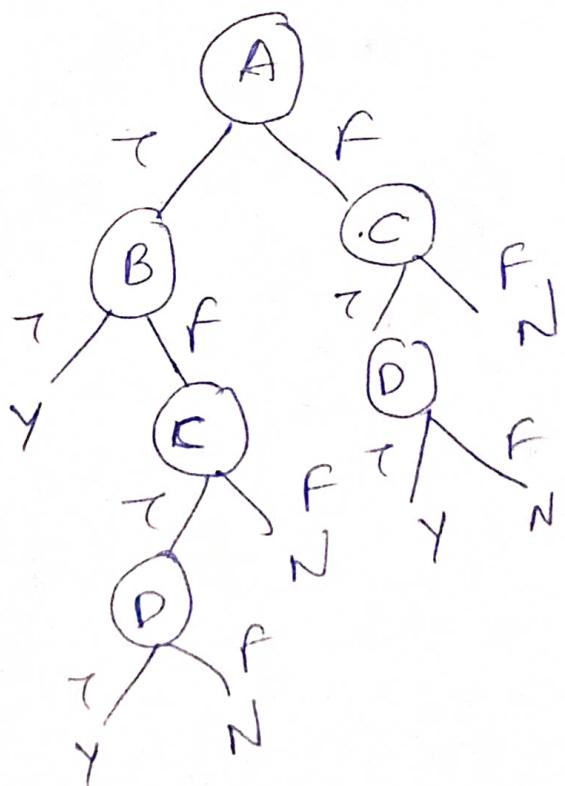


c) $A \oplus B$

$$(A \wedge \neg B) \vee (\neg A \wedge B)$$



d) $[A \wedge B] \vee [C \wedge D]$



Q) calculate Entropy and Information

Gain

Day	outlook	temp	humidity	wind	Play tennis
D1	sunny	Hot	High	weak	No
D2	sunny	Hot	High	strong	No
D3	overcast	Hot	High	weak	Yes
D4	Rain	Mild	Normal	weak	Yes
D5	Rain	cool	Normal	strong	No
D6	Rain	cool	Normal	strong	Yes
D7	overcast	cool	High	weak	No
D8	sunny	mild	Normal	weak	Yes
D9	sunny	cool	Normal	weak	Yes
D10	Rain	mild	Normal	strong	Yes
D11	sunny	mild	High	strong	Yes
D12	overcast	mild	Normal	weak	Yes
D13	overcast	Hot	Normal	strong	No
D14	Rain	mild	High	strong	No

Entropy :-

$$\text{Entropy}(S) = -P_0 \log_2 P_0$$

$$-P_1 \log_2 P_1$$

A

* if we have only positive examples and only negative examples then it is $\rightarrow 0$?

In:-
 $\text{Entropy}([4+, 0-]) = 0$

* if there is equal no of positive and negative then "1"
 $\text{Entropy}([7+, 7-]) = 1$

Information gain

$$\text{gain}(S, A) \equiv \text{Entropy}(S)$$

$$-\sum_{\text{v.e values (A)}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

Now for the table we find
for wind

values (wind) = weak, strong

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

$$S_{\text{weak}} \leftarrow [6+, 2-]$$

$$S_{\text{strong}} \leftarrow [3+, 3-]$$

$$\text{Gain}(S, \text{wind}) = \text{Entropy}(S) - \sum \frac{|S_i|}{|S|} \text{Entropy}_{(S_i)}$$

UE (weak, strong)

$$= \text{Entropy}(S) - \left(\frac{8}{14}\right) \text{Entropy}(S_{\text{weak}})$$

$$- \left(\frac{6}{14}\right) \text{Entropy}(S_{\text{strong}})$$

$$= 0.940 - \left(\frac{8}{14}\right) 0.811$$

$$- \left(\frac{6}{14}\right) 1.00$$

$$= 0.048$$

$$\boxed{\text{Entropy}([9+, 5-])}$$

$$= \frac{9}{14} \log_2 \left(\frac{9}{14}\right)$$

$$- \frac{5}{14} \log_2 \left(\frac{5}{14}\right)$$

$$= 0.94$$

Q) Given

$$P_1 = 0.1, P_2 = 0.2, P_3 = 0.3, P_4 = 0.4$$

Entropy?

Sol:- Entropy = $-\sum_{i=1}^n P_i \log_2(P_i)$

$$= -P_1 * \log_2(P_1) - P_2 * \log_2(P_2) - P_3 * \log_2(P_3) \\ - P_4 * \log_2(P_4)$$

$$= -0.1 * \log_2(0.1) - 0.2 * \log_2(0.2)$$

$$- 0.3 * \log_2(0.3) - 0.4 * \log_2(0.4)$$

$$= -0.1 * (-3.322) - 0.2 * (-2.322)$$

$$- 0.3 * (-1.736) - 0.4 * (-1.322)$$

$$= 0.3322 + 0.4644 + 0.5208 + 0.5284$$

$$= 1.8462$$

$$\log(0.1) = \frac{\log(1)}{10}$$

Q) Find Gini Index for splitting continuous attribute

Annual income	Label	Split point
60	No	
70	No	
75	No	< 80
85	Yes	≥ 80
90	Yes	
95	Yes	< 97.5
100	No	≥ 97.5
120	No	
125	No	
220	No	

} if not in order arrange in order

Annual income	Label	split points	Yes	No	Gini
60	No				
70	No				
75	No	< 80	0	3	
85	Yes	≥ 80	3	4	0.3427
90	Yes				
95	Yes	< 97.5	3	3	
100	No	≥ 97.5	0	4	0.30
120	No				
125	No				
220	No				

• split point = 80

$$\cdot \text{Gini}(< 80) = 1 - \sum_{i=1}^c (p_i)^2$$

$$= 1 - \left(\frac{0}{3}\right)^2 - \left(\frac{3}{3}\right)^2 = 0.0$$

$$\cdot \text{Gini}(\geq 80) = 1 - \sum_{i=1}^c (p_i)^2$$

$$= 1 - \left(\frac{2}{2}\right)^2 - \left(\frac{4}{2}\right)^2 = 0.4897$$

• Gini(s_0)

$$= w_1 * \text{Gini}(< 80) + w_2 * \text{Gini}(\geq 80)$$

$$= \frac{3}{10} * 0.0 + \frac{7}{10} * 0.4897$$

$$= 0.3427$$

• split point = 92.5

$$\cdot \text{Gini}(< 92.5) = 1 - \sum_{i=1}^c (p_i)^2$$

$$= 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{7}{6}\right)^2 = 0.6$$

$$\text{Gini}(\geq 97.5) = 1 - \left(\frac{6}{10}\right)^2 - \left(\frac{4}{10}\right)^2 \\ = 0.0$$

$$\text{Gini}(97.5) = w_1 * \text{Gini}(\geq 97.5) \\ + w_2 * \text{Gini}(\geq 97.0)$$

$$= \left(\frac{6}{10}\right) * 0.5 + \frac{4}{10} * 0.0 \\ = 0.30$$

So for this continuous

splittings Attribute 97.5
is best point

① Entropy

- Information Gain

- Gini Index

- Splittings Attribute

Instance	a ₁	a ₂	a ₃	target class
1	T	T	1.0	+
2	T	T	6.0	+
3	T	F	5.0	-
4	F	F	4.0	+
5	F	T	7.0	-
6	F	T	1.0	-
7	F	F	8.0	-
8	T	F	7.0	+
9	F	T	5.0	-

$$\textcircled{1} \quad \text{Entropy} = - \sum_{i=1}^n p_i (\log_2(p_i))$$

$$\text{Entropy}(I) = -\frac{4}{9} \log_2\left(\frac{4}{9}\right) - \frac{5}{9} \log_2\left(\frac{5}{9}\right)$$

$$= 0.9911$$

① information gain of a,

$$\text{Entropy}(S) = \sum_{i=1}^n p_i \log_2(p_i)$$

$$\begin{aligned}\text{Entropy}(S_T) &= -\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) \\ &= 0.311 + 0.1 = 0.411\end{aligned}$$

$$\begin{aligned}\text{Entropy}(S_F) &= -\frac{1}{5} \log_2\left(\frac{1}{5}\right) - \frac{4}{5} \log_2\left(\frac{4}{5}\right) \\ &= 0.4644 + 0.2576 \\ &= 0.722\end{aligned}$$

$\gamma \rightarrow \gamma$
 $s \rightarrow \emptyset$

$$\text{Gain}(a_1) = \text{Entropy}(S) - \sum_{v \in T(F)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\begin{aligned}\text{Gain}(a_1) &= \text{Entropy}(S) - \frac{4}{9} \text{Entropy}(S_T) \\ &\quad - \frac{5}{9} \text{Entropy}(S_F)\end{aligned}$$

$$\text{Gain}(a_1) = 0.2295$$

for information gain for a_2

$$\text{Entropy}(S_T) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$$
$$= 0.9709$$

$$\text{Entropy}(S_F) = -\frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right)$$
$$= 0.5 + 0.5 = 1.0$$

$$\text{Gain}(a_2) =$$

$$\text{Entropy}(S) - \frac{5}{9} \text{Entropy}(S_T)$$
$$- \frac{4}{9} \text{Entropy}(S_F)$$

$$= 0.0072$$

③ Gini Index - of attributes a_1

$$\text{Gini} = 1 - \sum_{i=1}^n (p_i)^2$$

$$\text{Gini}(T) = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 0.271$$

$$\text{Gini}(F) = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 = 0.32$$

$$\text{gini index}(a_1) = \sum_{v \in \{r, F\}} \frac{|S_v|}{|S|} \text{gini}(S_v)$$

$$\therefore \left(\frac{4}{9}\right) * \text{gini}(r) + \left(\frac{5}{9}\right) * \text{gini}(F) \\ = 0.3444$$

for similar for gini index
of a_2 also

④ best splitting attribute
between a_1 and a_2

$$\text{gain}(a_1) = 0.2291$$

$$\text{gain}(a_2) = 0.0072$$

should be high

Here a_1 is the best

split attribute

Now

$$\text{gini index}(a_1) = 0.3444$$

$$\text{gini index}(a_2) = 0.4849$$

should be low

Here a_1 is the best split attribute