

Classification

Classification is the process of building a model that predicts the class label for unseen example/tuple/sample.

It involves two steps

1) Learning (Training)

2) Testing (prediction)

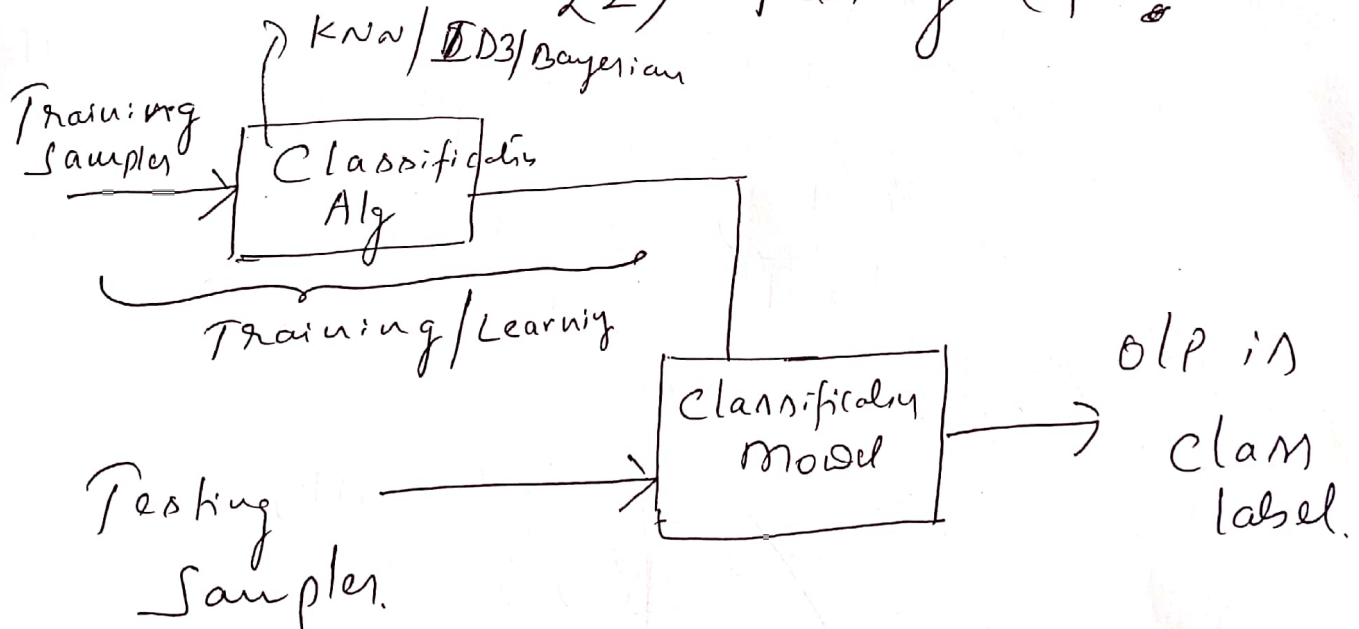


Fig:- A General Classification System.

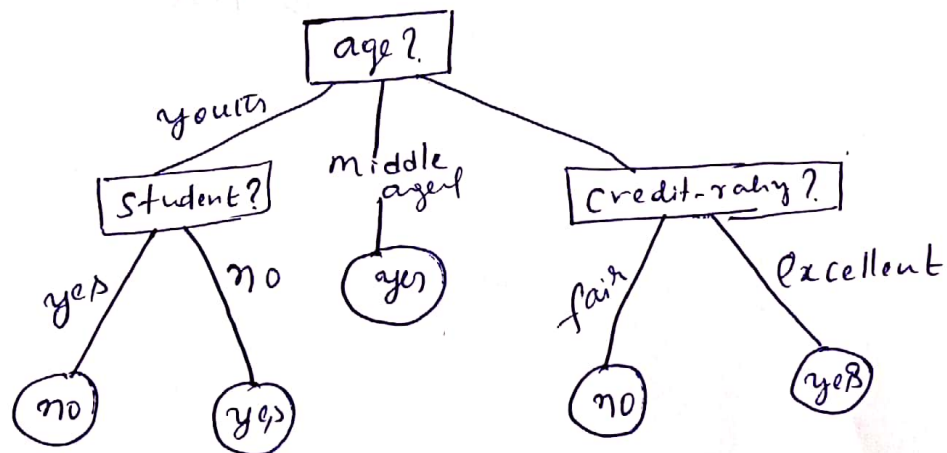
Classification Methods

- (1) Decision Tree Induction
- (2) Bayes Algorithm
- (3) Rule-based ~~Ab~~ Classification Algorithm
- (4) K-Nearest Neighbour Algorithm
- (5) Neural Network.

Decision Tree Induction

→ Decision tree induction is the learning of Decision Trees from class-labelled training examples (tuples).

→ A Decision tree is a flowchart-like tree structure, where each internal node (nonleaf node) denotes a test on an attribute, & Each leaf node holds a class label. The topmost node in tree is the root node.



Example ∴ Decision tree for the concept buys-computer

→ The Learning (training) and Classification (testing)

Steps of Decision tree including are

Simple & fast.

→ It is an example for Nonparametric Classifier.

→ Decision tree classifiers have good accuracy.

→ Decision Tree algorithms have been used for classification in many applications areas such as Medicine, manufacturing & products, financial analysis, & Molecular biology.

- In 1980, J Ross Quinlan Developed a Decision tree algorithm known as ID3.
- Later, he presented C4.5 (a successor of ID3), which became benchmark to compare with newer ~~Supp~~ Supervised algorithms.
- In 1984, a group of statisticians published the book Classification and Regression Trees (CART), which described the generation of binary Decision Trees.
- ID3, C4.5 and CART adopt a Greedy approach in which Decision trees are constructed in top-down recursive manner.

for Example

Consider D - Data partition

Let

Attribute list

Say $\{A_1, A_2, A_3, A_4\}$

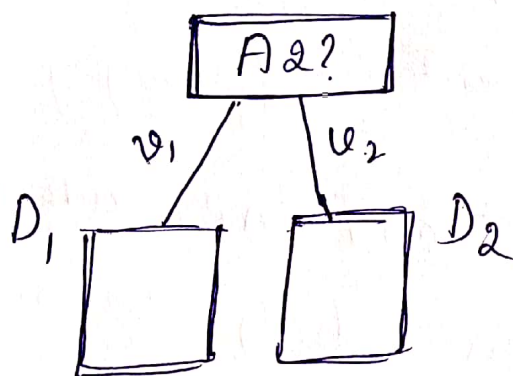
	A_1, A_2, A_3, A_4	Class
T_1		yes
T_2		yes
\vdots	\vdots	\vdots
T_n		no

Let

No. of
Classes - 2
yes, no

Then Select the best ^{splitting} attribute such that
the resultant partitions are
pure

Let A_2 is selected as best splitting attribute



Then

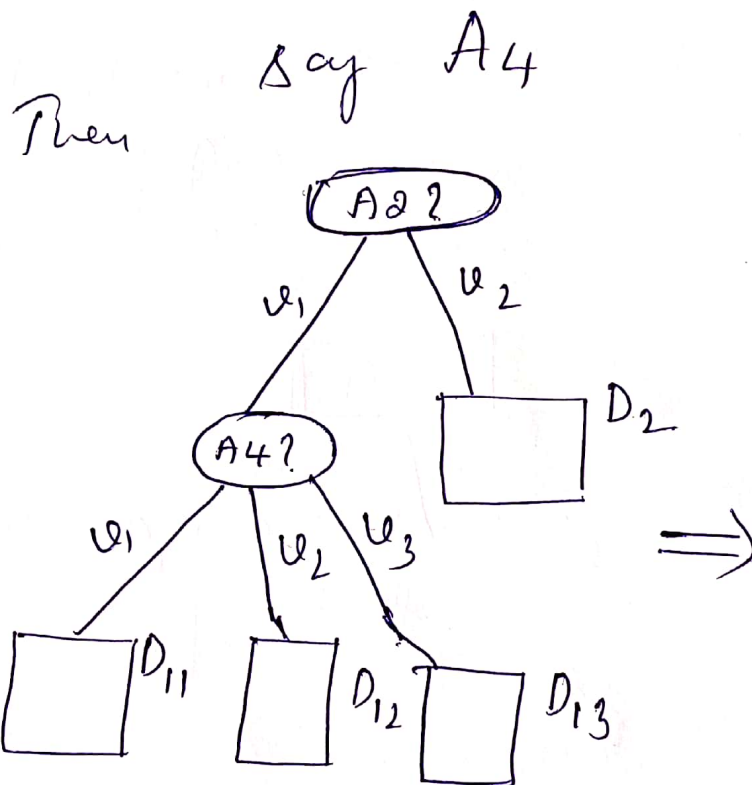
Algorithm will check first D_1 if
pure
or
impure?

Assume ^{partition} resultant D_1 is impure.

Then Again Split D_1

Therefore ~~Select~~ Select best split attribute
from $\{A_1, A_2, A_3, A_4\} - A_2$

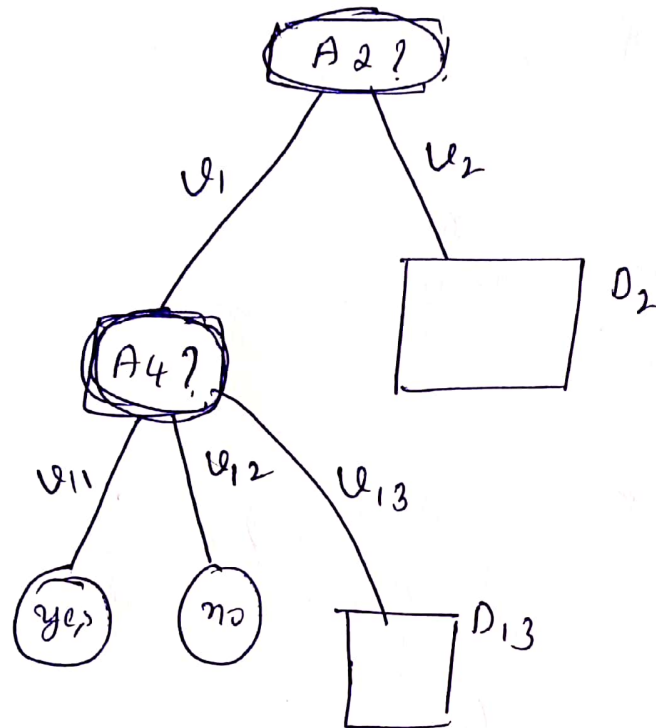
Select from $\rightarrow \{A_1, A_3, A_4\}$
best split attribute



Assume D_{11}, D_{12}, D_{13} are pure

Say D_{11} contains only yes
 $D_{12} \rightarrow$ no
 $D_{13} \rightarrow$ no

Then ,

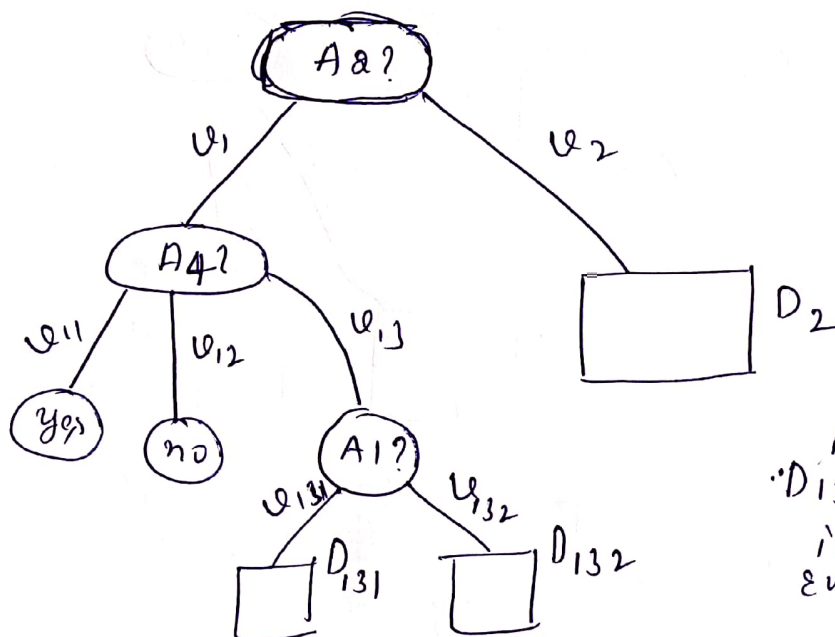


Then Select the best Split attribute.

among remaining attributes

$\{A_1, A_3\}$ to split D_{13}

Say best in A_1



Assume D_{132} is empty

Assume D_{131} is pure (all tuples are no (seller))

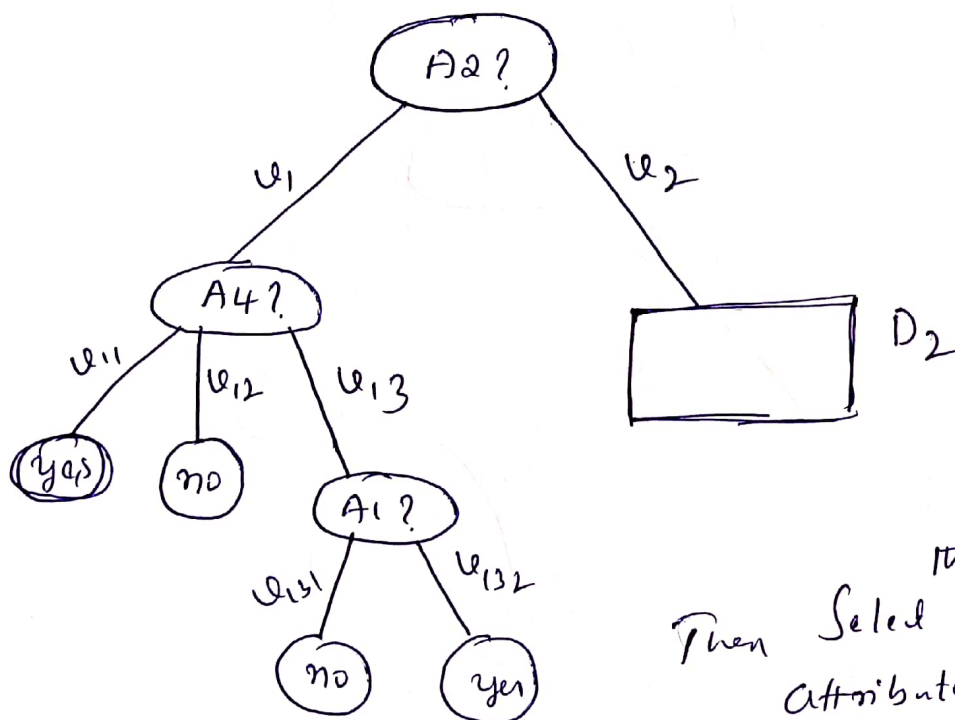
Then

label D_{131} as "no" class label

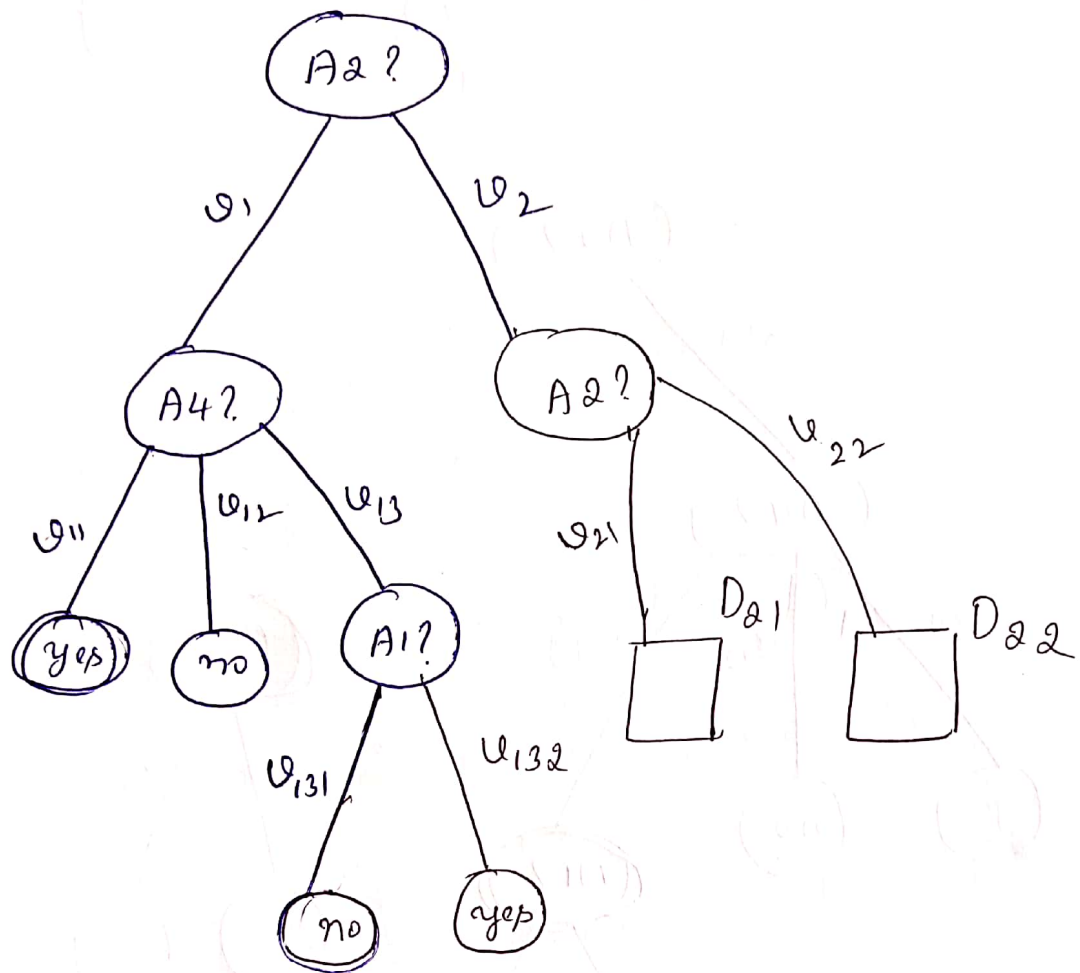
Since D_{132} is Empty (no tuples)

Then we need to consider
majority class in the
above partition D_{13}

based on that, we have
to label ~~set~~ ^{in D_{13}}
Assume majority class is "yes"
Then label D_{132} as "yes"



Then select the best splitting
attribute
to split partition D_2
from list $\{A3\}$



Assume if partition D_{21} is pure
(all tuples are
say "no" labelled)

Then label D_{21} as "no"

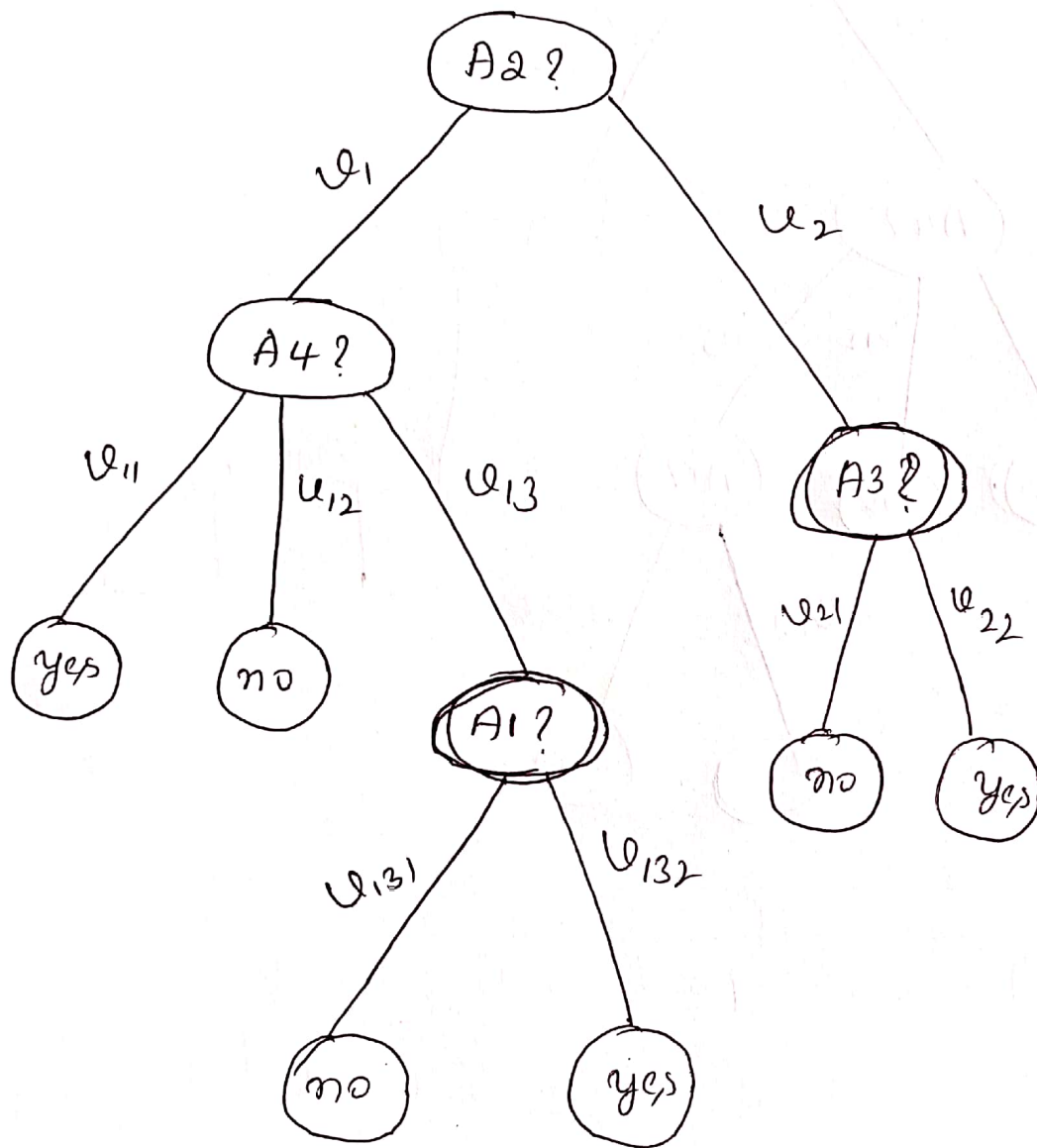
& if D_{22} is impure

Again we have to split D_{22}
but No attributes left

In this situation, we have to apply
majority voting in D_{22}

Assume say majority vote in D_{22} is "yes"

Turn label D_{g2} as "yes"



Algorithm: Generate-Decision-Tree

// Generate Decision Tree from training tuples of Data partition(D)

Input:

⊗ Data partition (D)
// Set of Training tuples & associated class labels.

⊗ attribute-list
// Set of attributes.

⊗ Attribute-Selection-method.

Output: A Decision Tree

Method

1. Create a Node N
2. if tuples in D are all of the same class C then
3. return N as a leaf node labelled with class C
4. if attribute-list is empty then
5. return N as a leaf node with majority class in D
- 6.

6. apply Attribute-Selection-method(D , attribute-list)
to find the "best" Splitting-Criterion;
7. Label node N with Splitting-Criterion;
8. if Splitting-attribute is discrete-valued and
multiways splits allowed then
9. attribute-list \leftarrow attribute-list
- Splitting-attribute;
10. for each outcome j of Splitting-Criterion
11. let D_j be The Set of Data tuples in D
satisfying outcome j ;
12. if D_j is Empty Then
13. attach a leaf labelled with
majority class in D
to node N ;
14. else attach The node returned by
Generate-decision-tree(D_j , attribute-list)
to node N ;
- ~~End for~~
- End for
15. return N ;