

# Classification Techniques

---

## ❖ Base Classifiers

❖ Decision Tree based Methods

❖ Rule-based Methods

❖ Nearest-neighbor

❖ Neural Networks, Deep Neural Nets

❖ Naïve Bayes and Bayesian Belief Networks

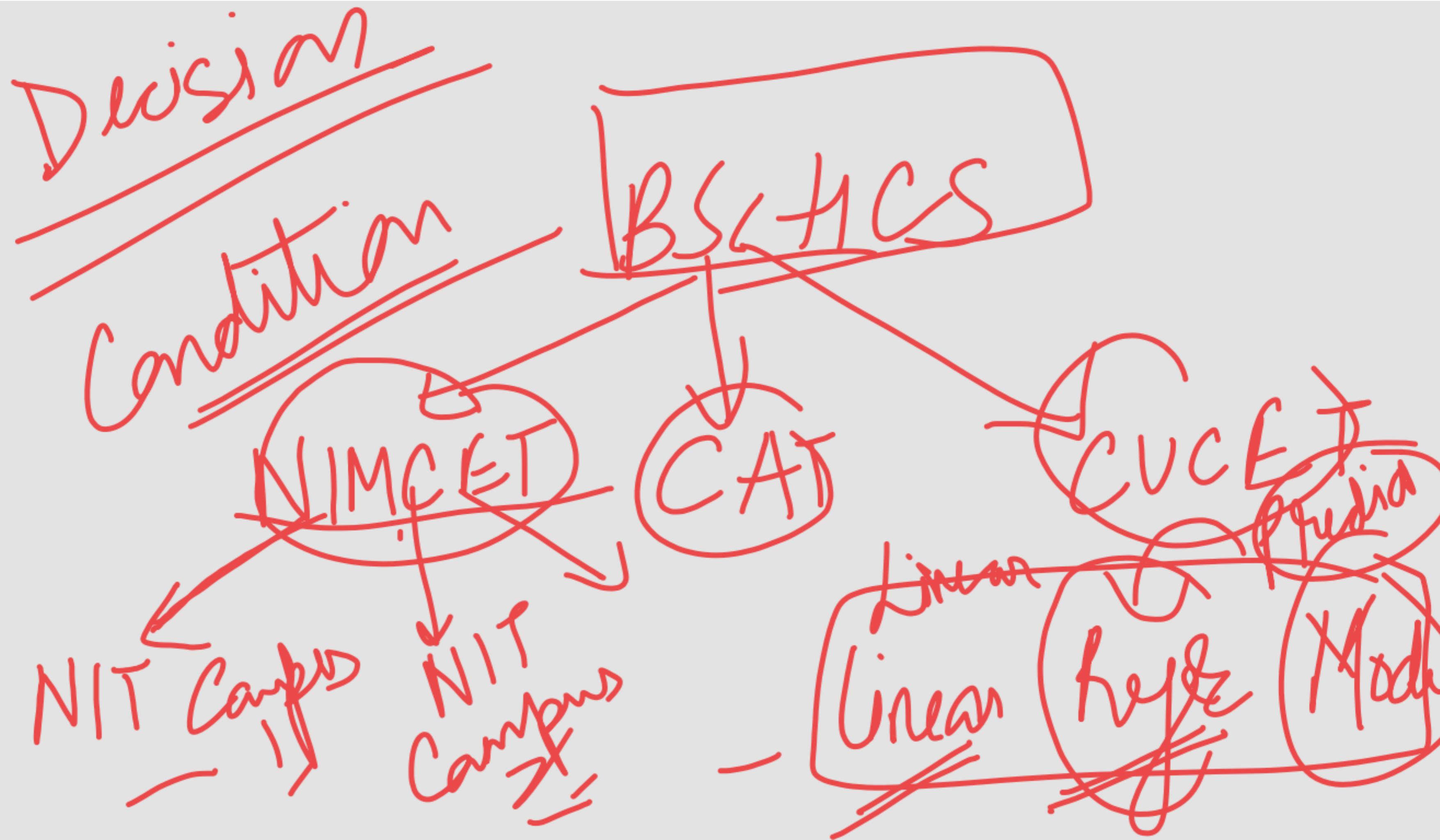
❖ Support Vector Machines

*Logistic Regression*

*Clustering*

## ❖ Ensemble Classifiers

❖ Boosting, Bagging, Random Forests

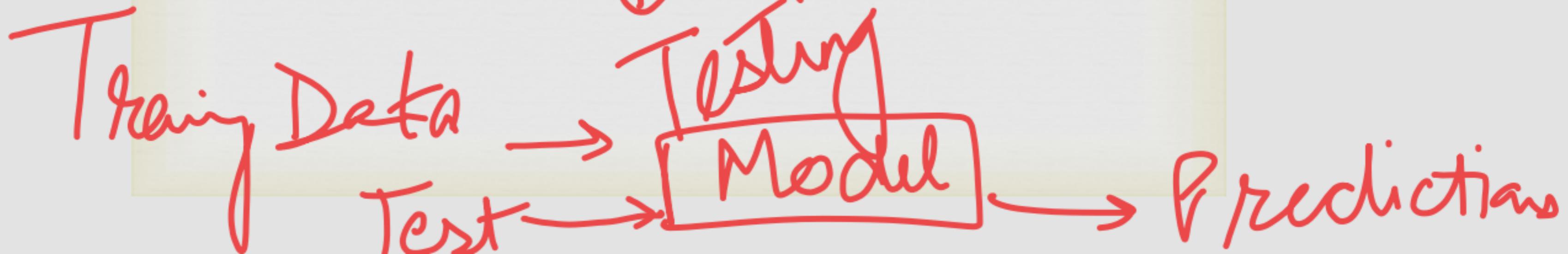




# Classification: Two Phases

B

1. Create a specific model by evaluating the training data. This step has an input the training data (including defined classification for each tuple) and as output a definition of model developed.
2. Apply the model developed in step 1 by classifying tuples from the target database.





❖ AIM:

- ❖ To build up a model/learning algorithm
- ❖ To identify a model that best fits the relationship between the attribute set and class label of the input data.
- ❖ The model generated by a learning algorithm should both fit the input data well and correctly predict the class labels of records it has never seen before. Therefore, a key objective of the learning algorithm is to build models with good generalization capability; i.e., models that accurately predict the class labels of previously unknown records.

Training data

Test data

Series of questions & their possible answers can be organized in form of decision tree

## Decision Tree-based Algorithms



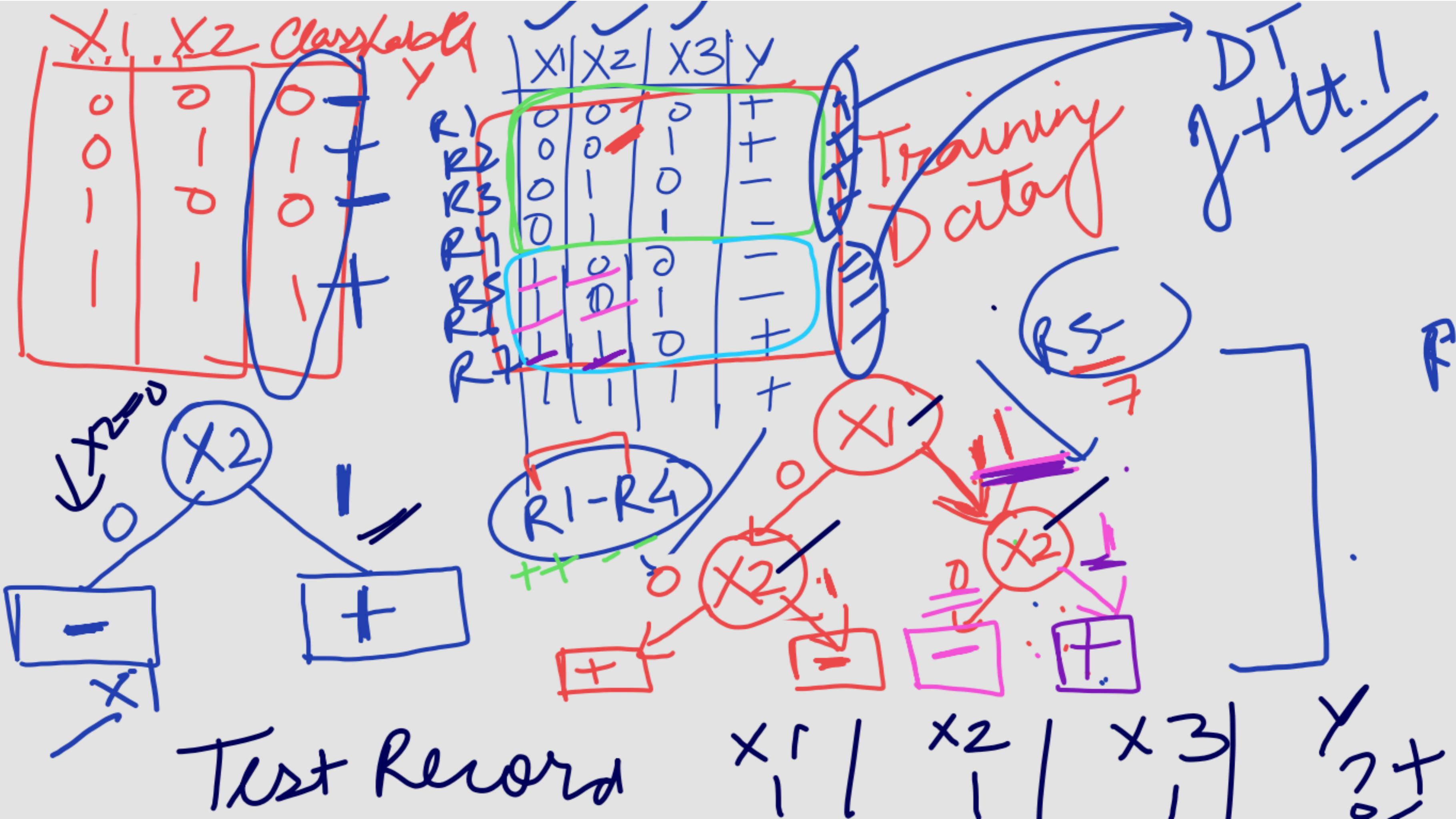
- Most useful approach for classification problems.
- Tree is constructed to model the classification process.
- Divide the search space into rectangular region.
- Tuple is classified based on the region into which it falls.

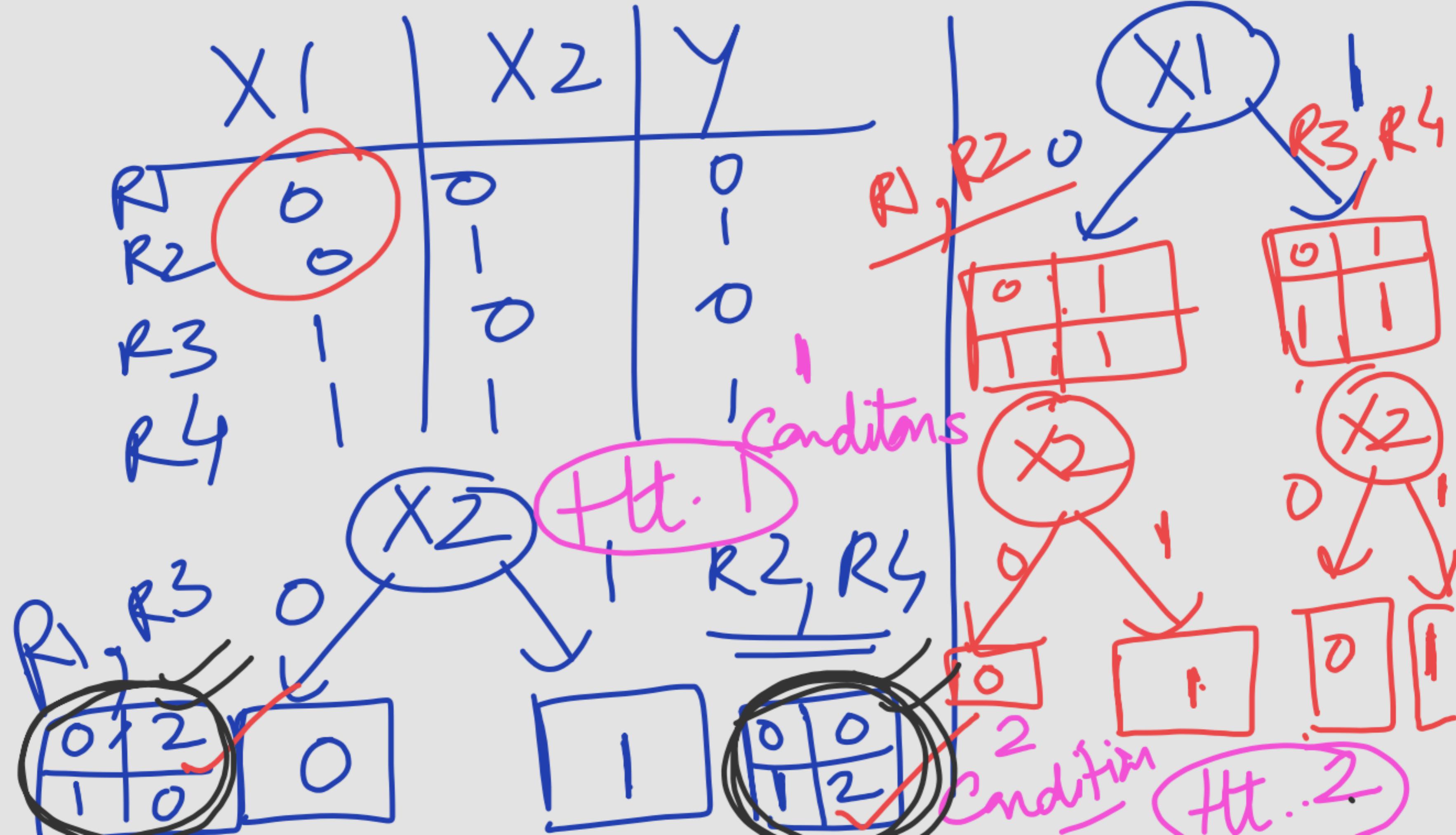


### Basic Steps:

1. Decision tree induction: Construct a DT using training data.
2. For each tuple  $t_i \in D$ , apply the DT to determine its class.

A1	A2	T
a	b	0
b	c	0
b	c	0
a		1





# Decision Tree Properties



Given a database  $D = \{t_1, t_2, \dots, t_n\}$  where  $t_i = \langle t_{i1}, \dots, t_{ih} \rangle$  and the database schema contains the following attributes  $\{A_1, A_2, \dots, A_h\}$ . Also given is a set of classes  $C = \{C_1, C_2, \dots, C_m\}$ . A decision tree (DT) of classification tree is a tree associated with D that has the following properties:

- Each internal node is labelled with an attribute,  $A_i$ .
- Each arc is labelled with a predicate that can be applied to the attribute associated with the parent.
- Each leaf node is labelled with a class,  $C_j$ .

wolf



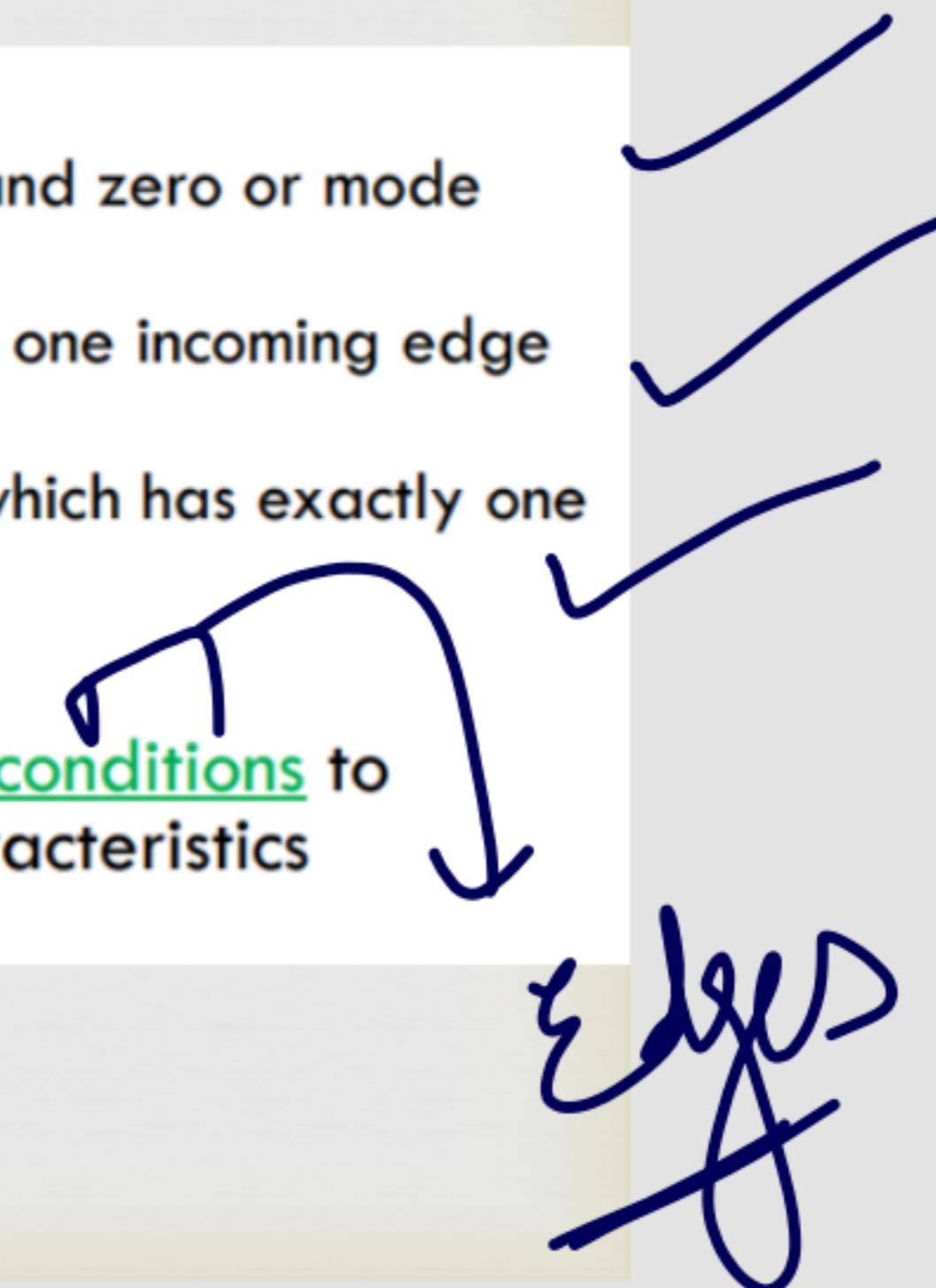
condition

- A decision tree has three types of nodes:

1. A root node that has no incoming edges and zero or mode outgoing edges
2. Internal nodes, each of which has exactly one incoming edge and two or more outgoing edges
3. Leaf nodes (or terminal nodes), each of which has exactly one incoming edge and no outgoing edges

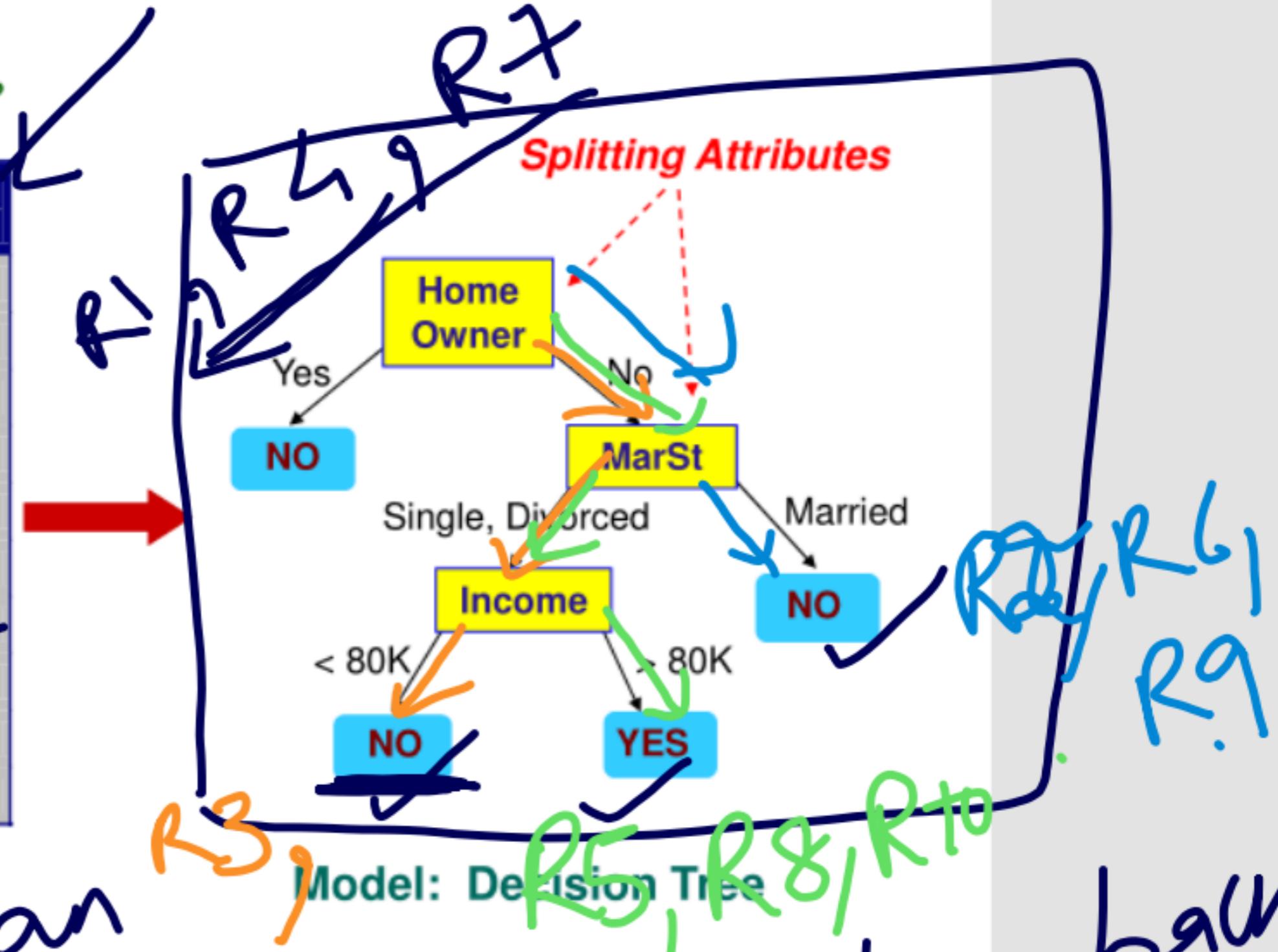
- Each leaf node is assigned a class label

- Non-terminal nodes contain attribute test conditions to separate records that have different characteristics

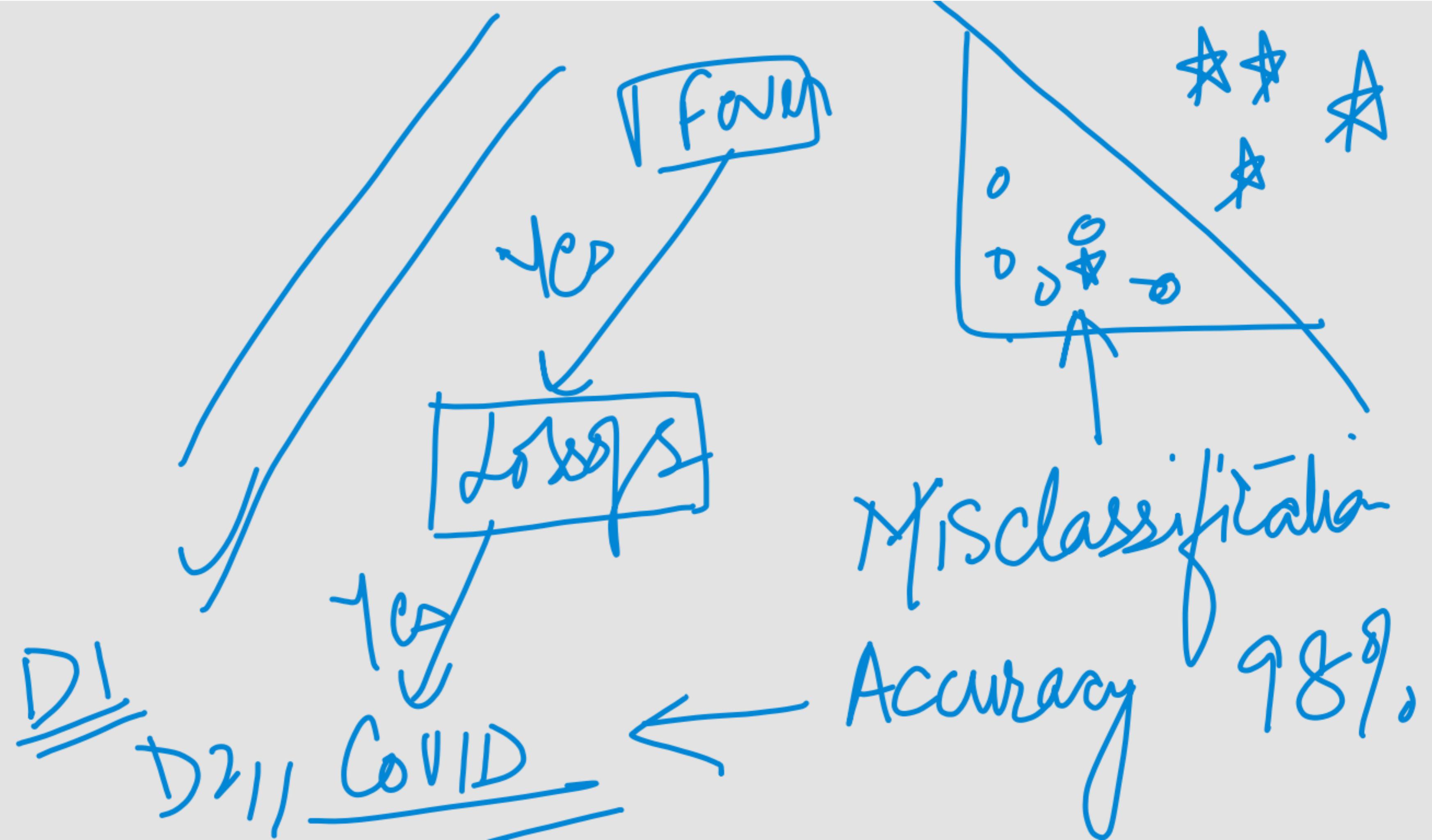


# Example of a Decision Tree

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



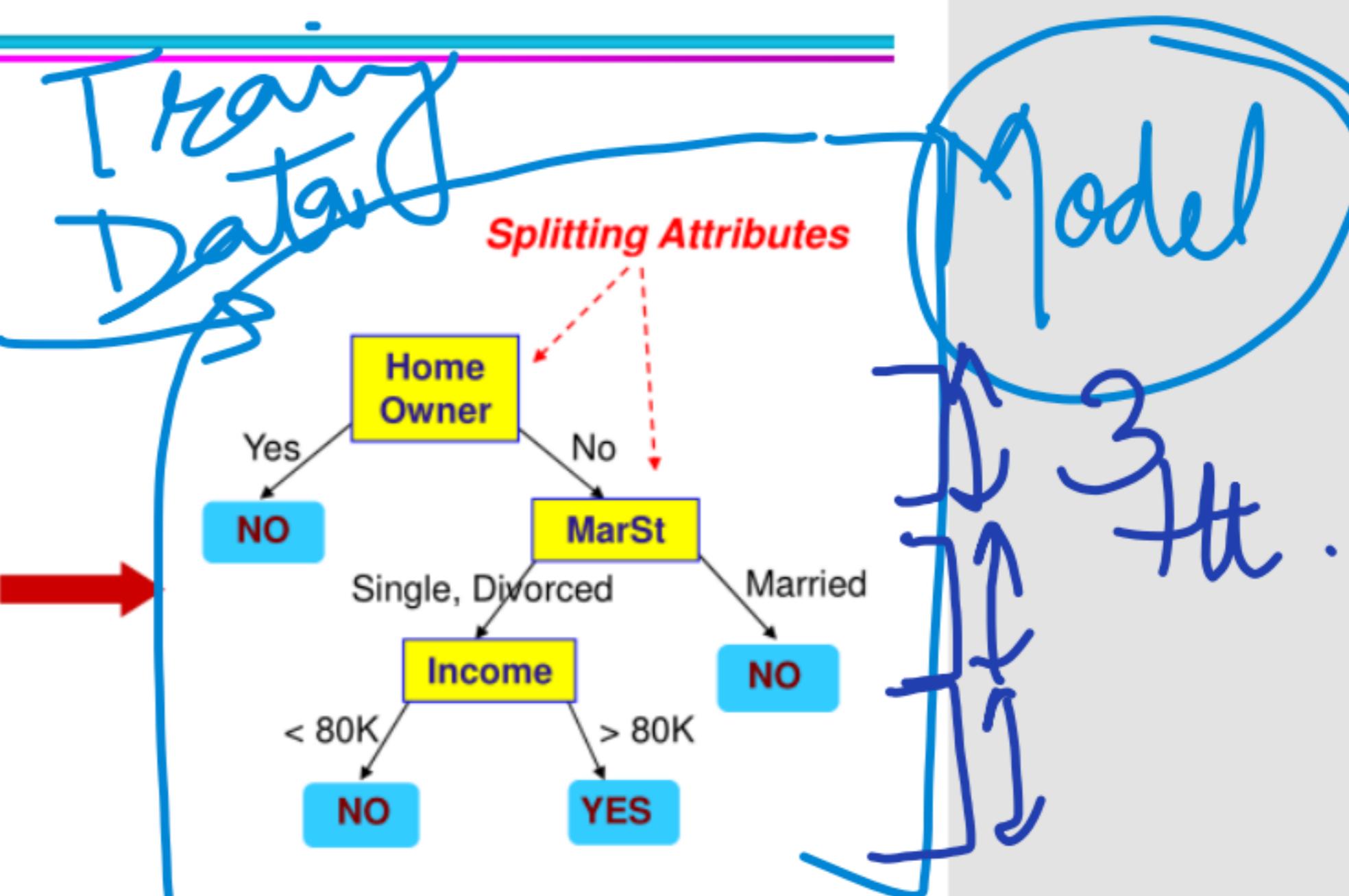
# No Yes, Defaulters: Will they Pay back the loan?



# Example of a Decision Tree

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

categorical  
categorical  
continuous  
class



Training Data

Model: Decision Tree

09/21/2020

Introduction to Data Mining, 2nd Edition

Purdue

Home Owner

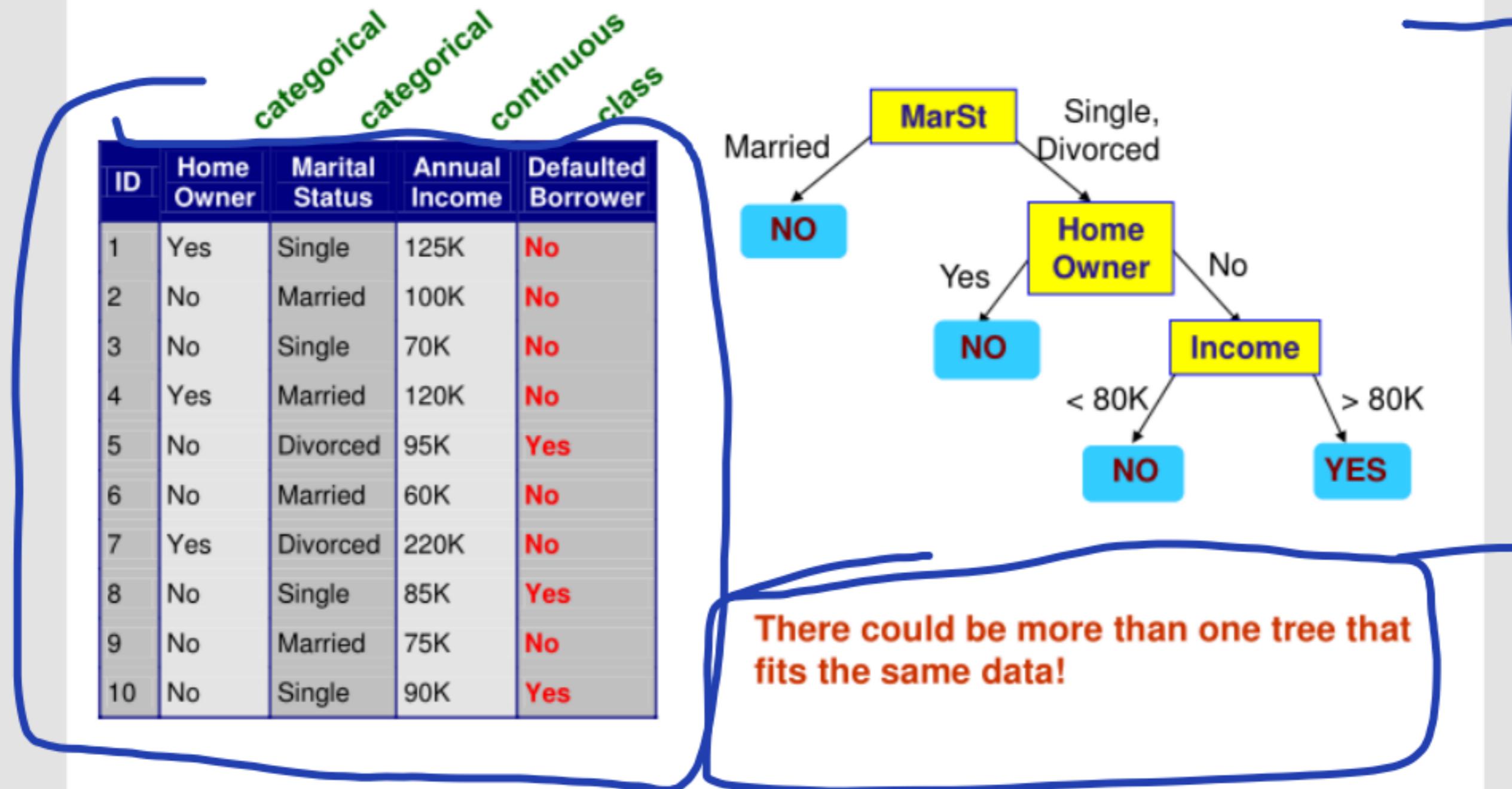
No

MS  
s  
A I  
60k

Defaulted Borrower

No 75.

# Another Example of Decision Tree



---

---

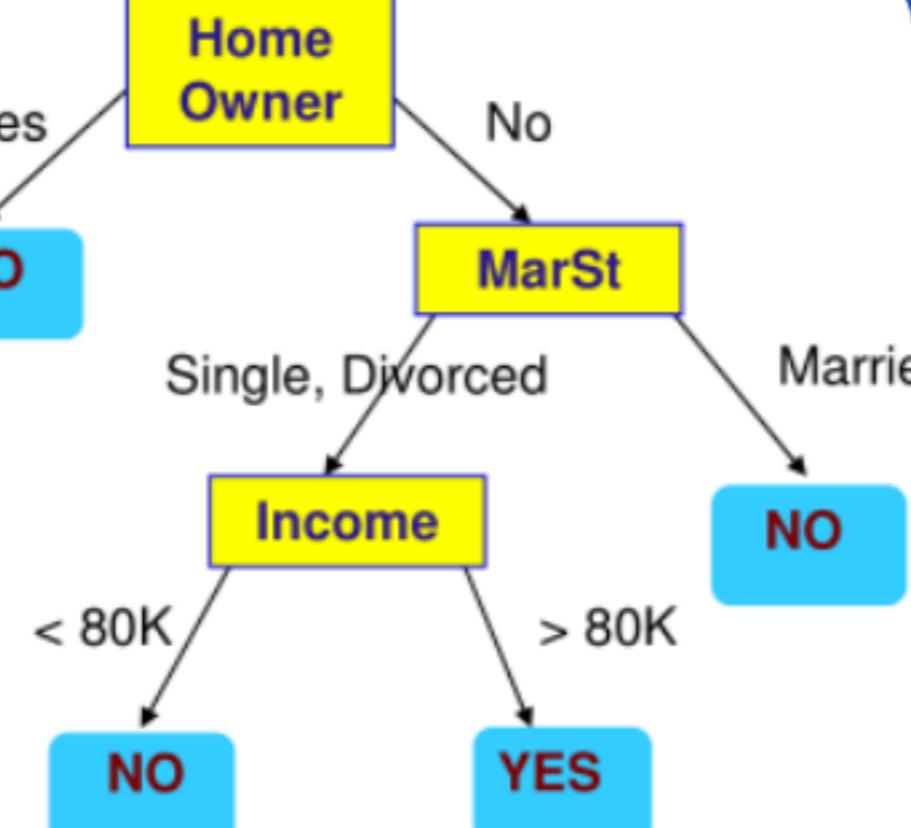
## **Classifying a test record once a decision tree is constructed**

Classifying a test record is straightforward once a decision tree has been constructed. Starting from the root node, we apply the test condition to the record and follow the appropriate branch based on the outcome of the test. This will lead us either to another internal node, for which a new test condition is applied, or to a leaf node. The class label associated with the leaf node is then assigned to the record. As an illustration, Figure 4.5 traces the path in

# Apply Model to Test Data

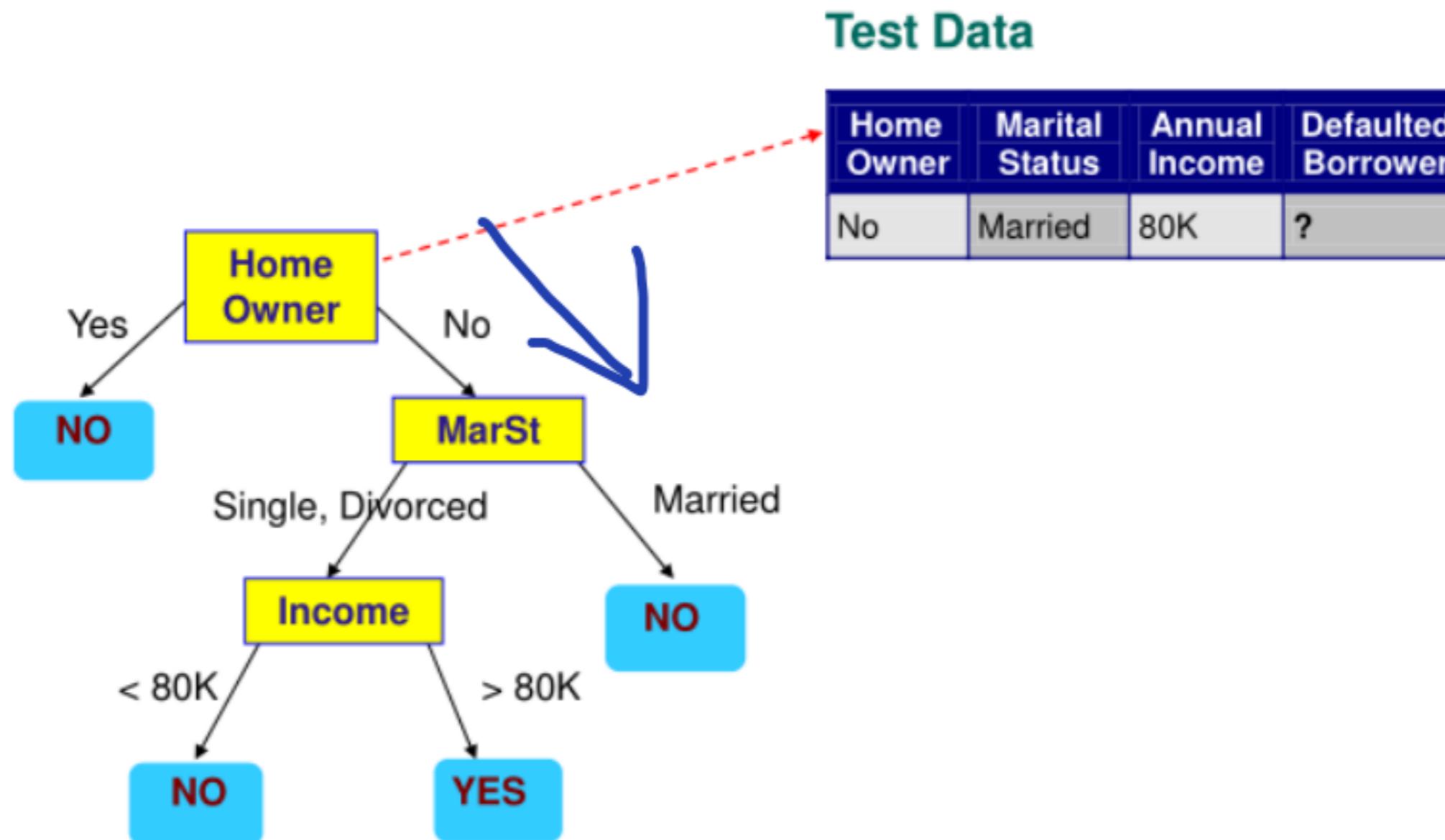
Test Data

Start from the root of tree.



Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?

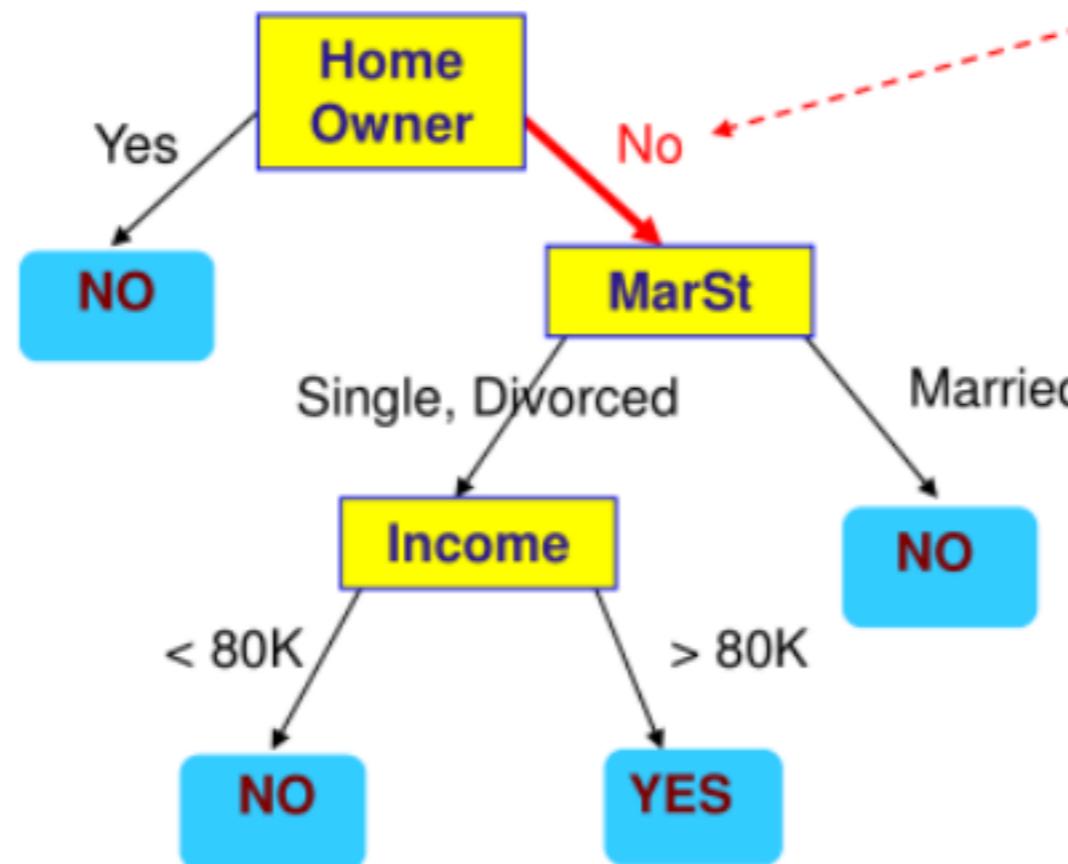
# Apply Model to Test Data



# Apply Model to Test Data

Test Data

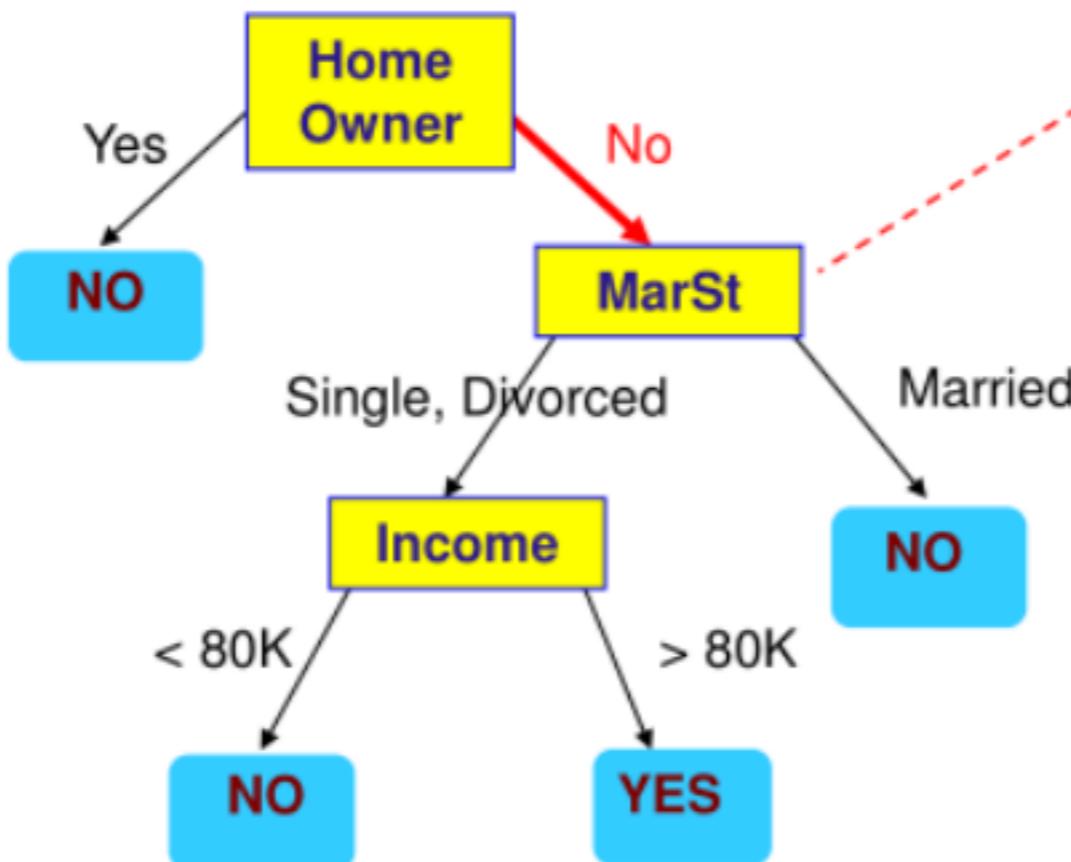
Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?



# Apply Model to Test Data

Test Data

Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?



# Apply Model to Test Data

Test Data

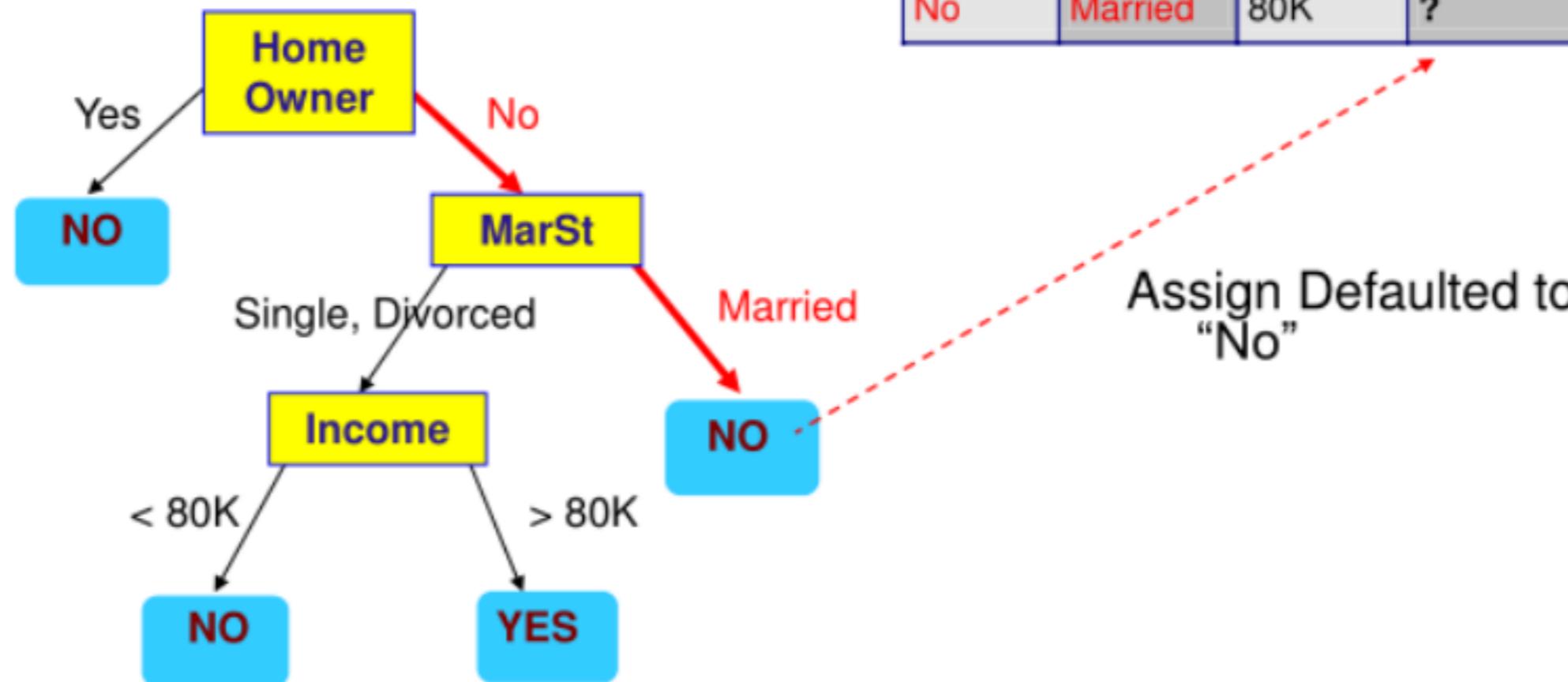
Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?



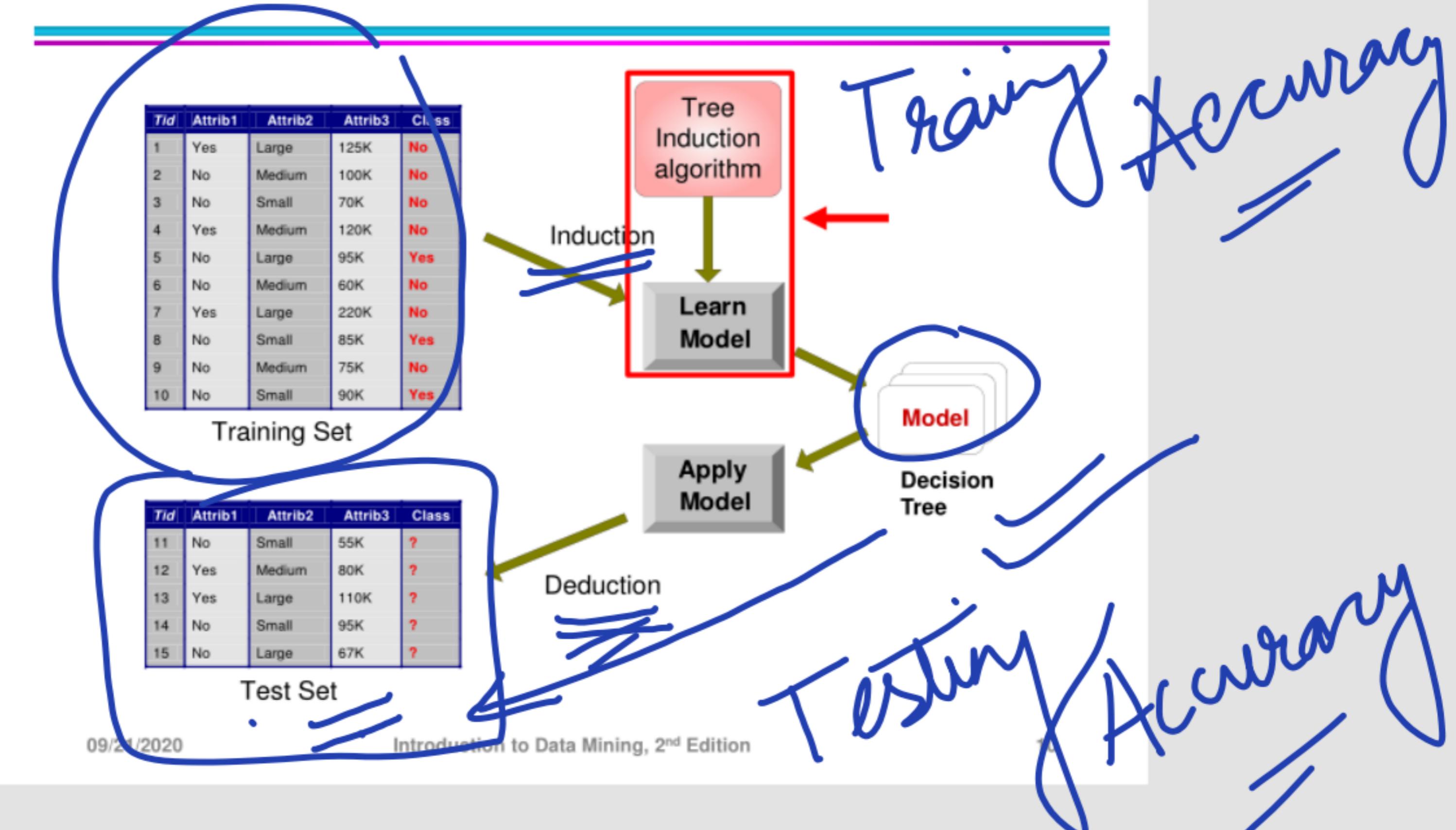
# Apply Model to Test Data

## Test Data

Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?



# Decision Tree Classification Task



# Decision Tree Induction

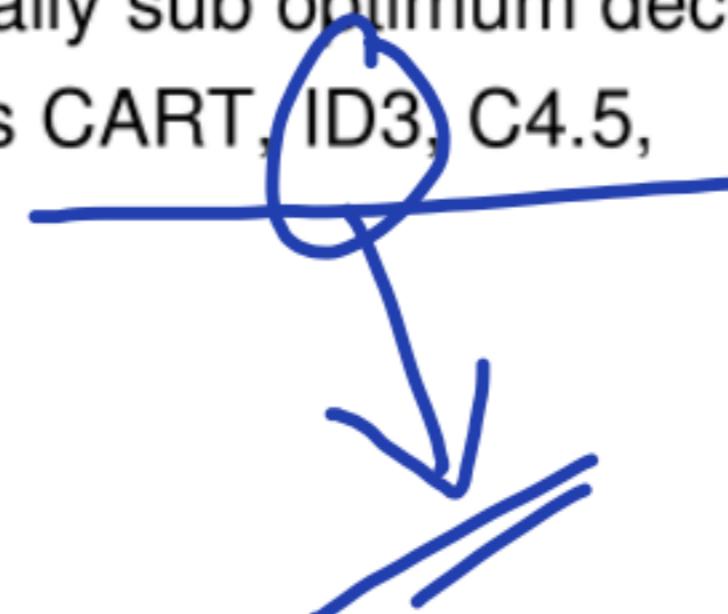
---

- Finding optimal tree is computationally expensive
  - Exponentially many decision tree.

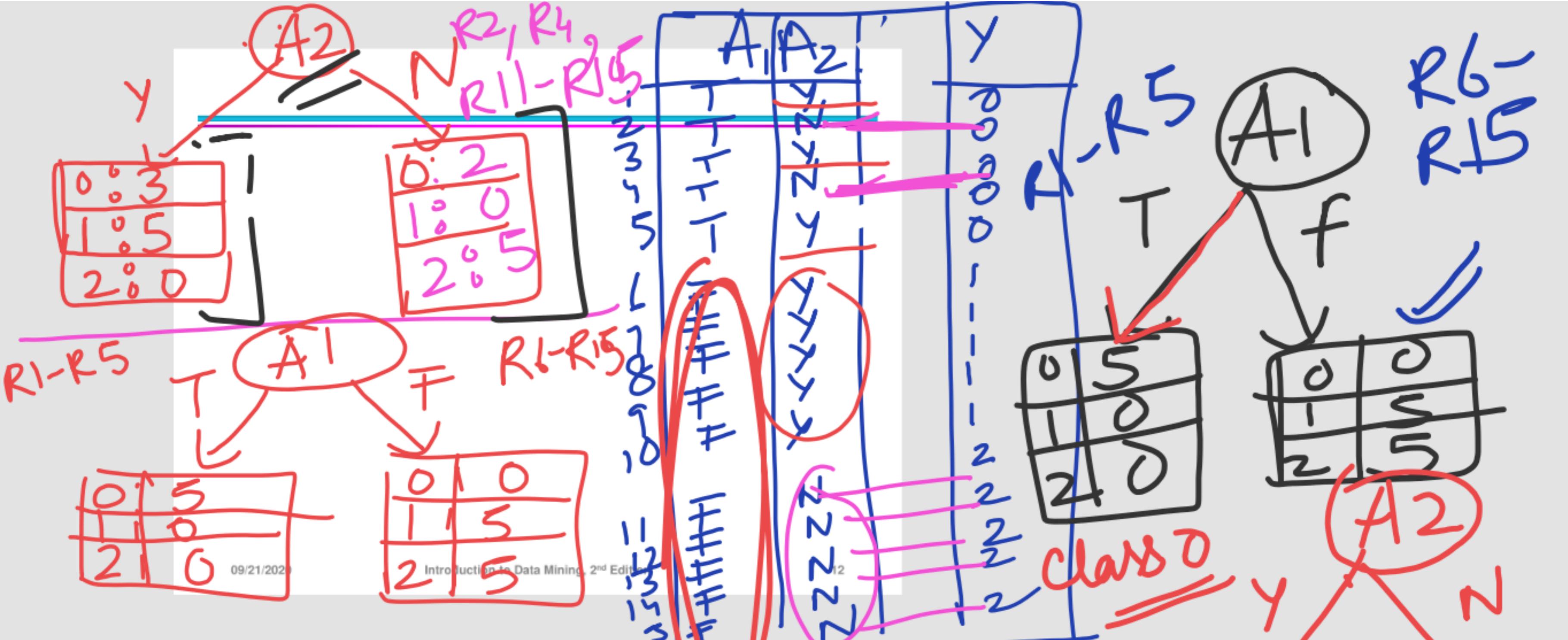
- Many Algorithms:

Hunt's Algorithm (one of the earliest)

- ◆ Employ greedy strategy that grows a decision tree by making series of locally sub optimum decisions.
- ◆ Basis of existing algos CART, ID3, C4.5, SLIQ, SPRINT







$$A_1 = \bar{F}$$

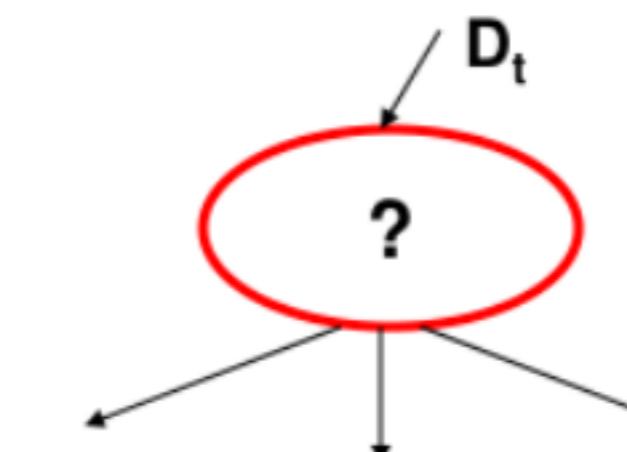
0	0
1	5
2	0

0	0
1	0
2	5

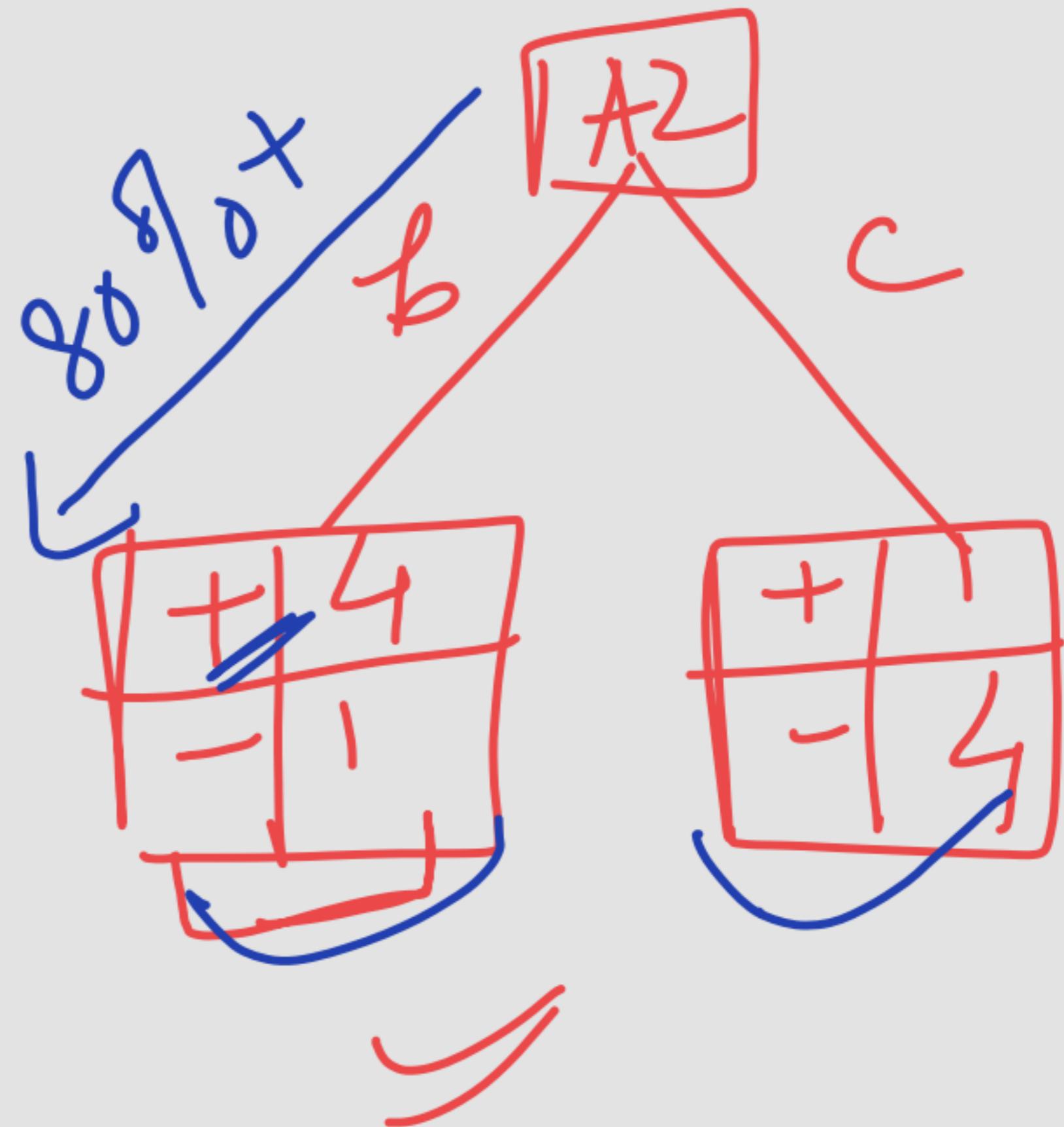
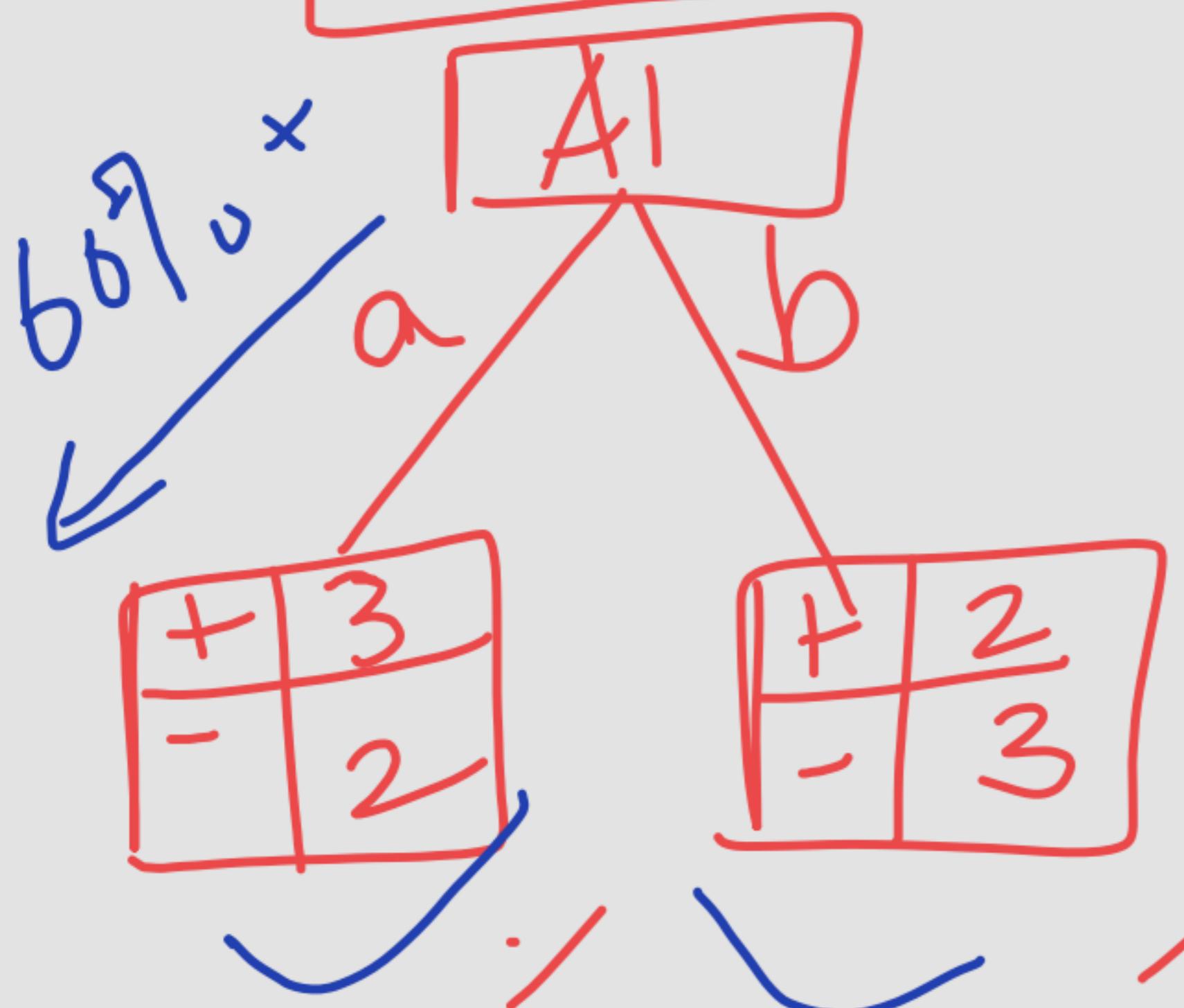
# General Structure of Hunt's Algorithm

- | Let  $D_t$  be the set of training records that reach a node  $t$
- | General Procedure:
  - If  $D_t$  contains records that belong to the same class  $y_t$ , then  $t$  is a leaf node labeled as  $y_t$
  - If  $D_t$  contains records that belong to more than one class, use an attribute test to split the data into smaller subsets. Recursively apply the procedure to each subset.

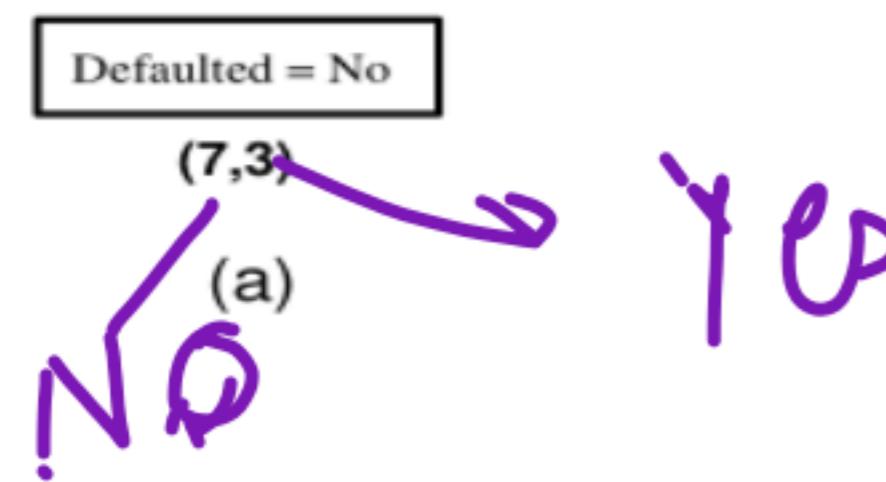
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Need not  
split it further



# Hunt's Algorithm

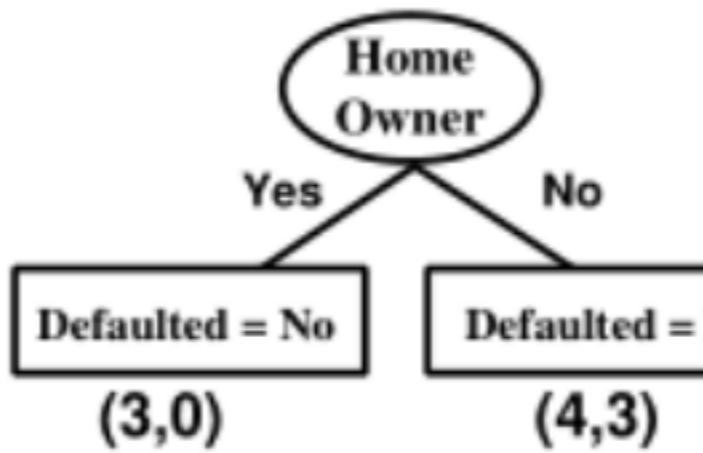


ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

# Hunt's Algorithm

Defaulted = No  
(7,3)

(a)



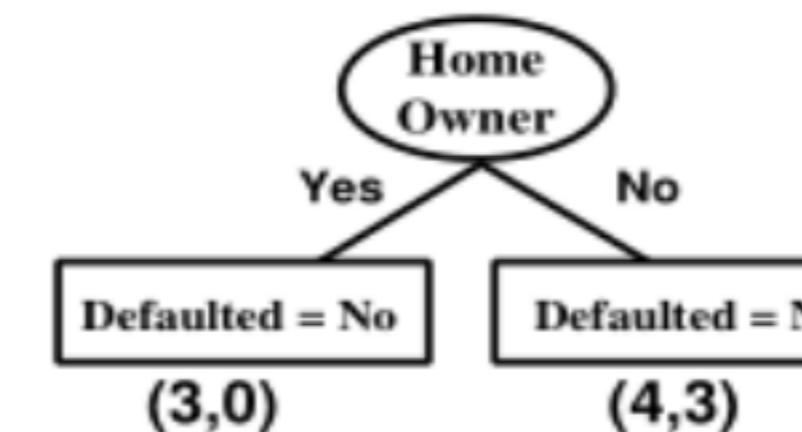
(b)

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

# Hunt's Algorithm

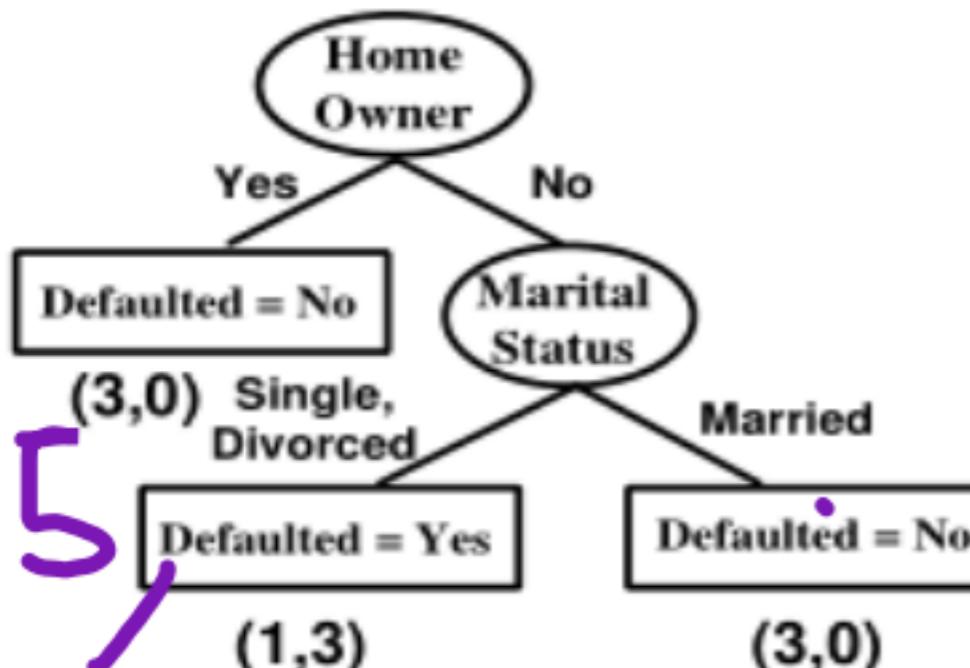
Defaulted = No  
(7,3)

(a)



(b)

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

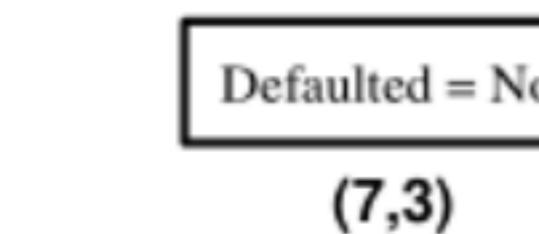


(c)

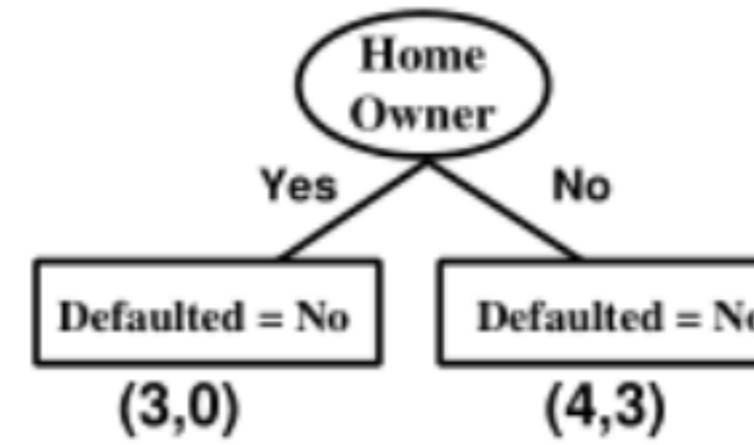
2,6,9

3,5  
8,10

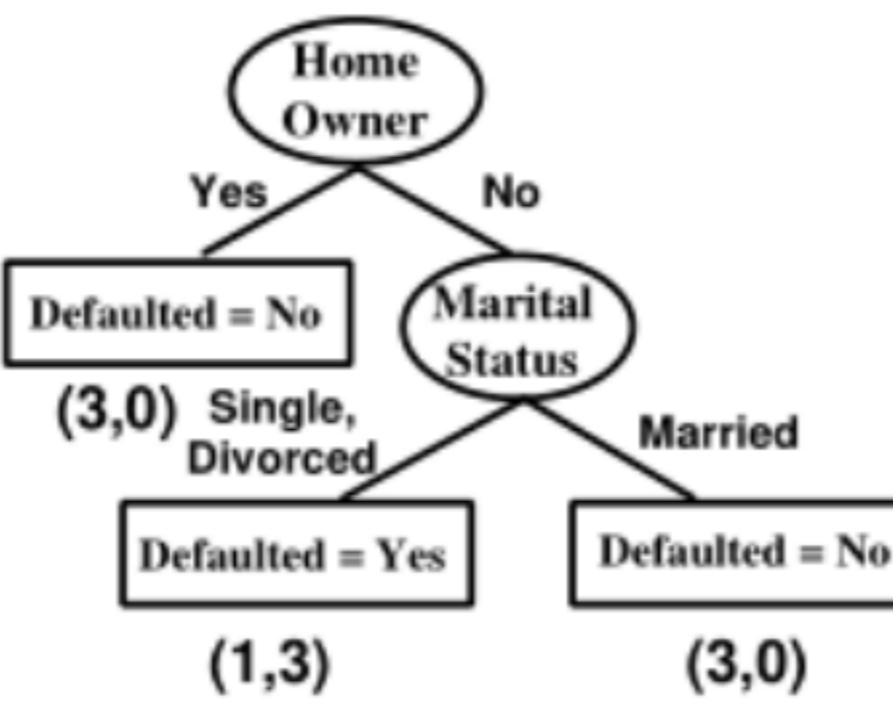
# Hunt's Algorithm



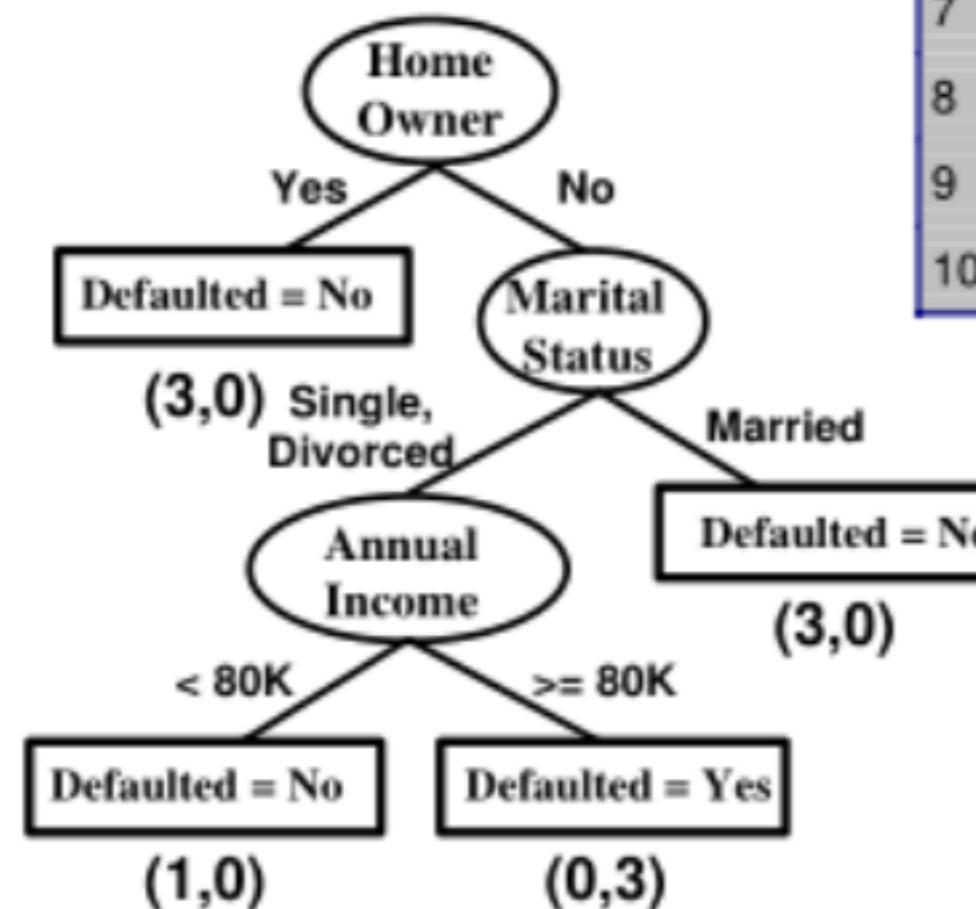
(a)



(b)



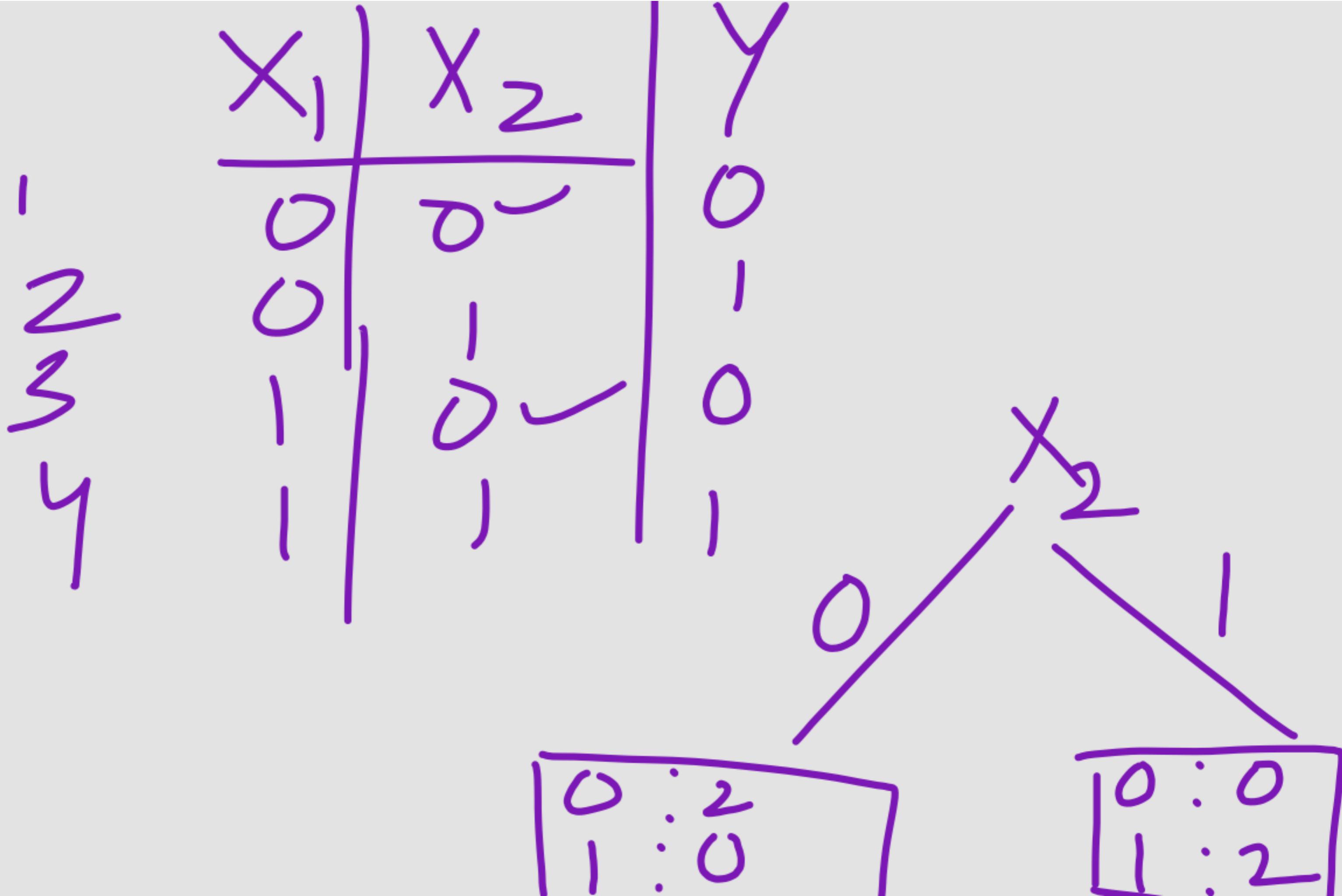
(c)



(d)

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

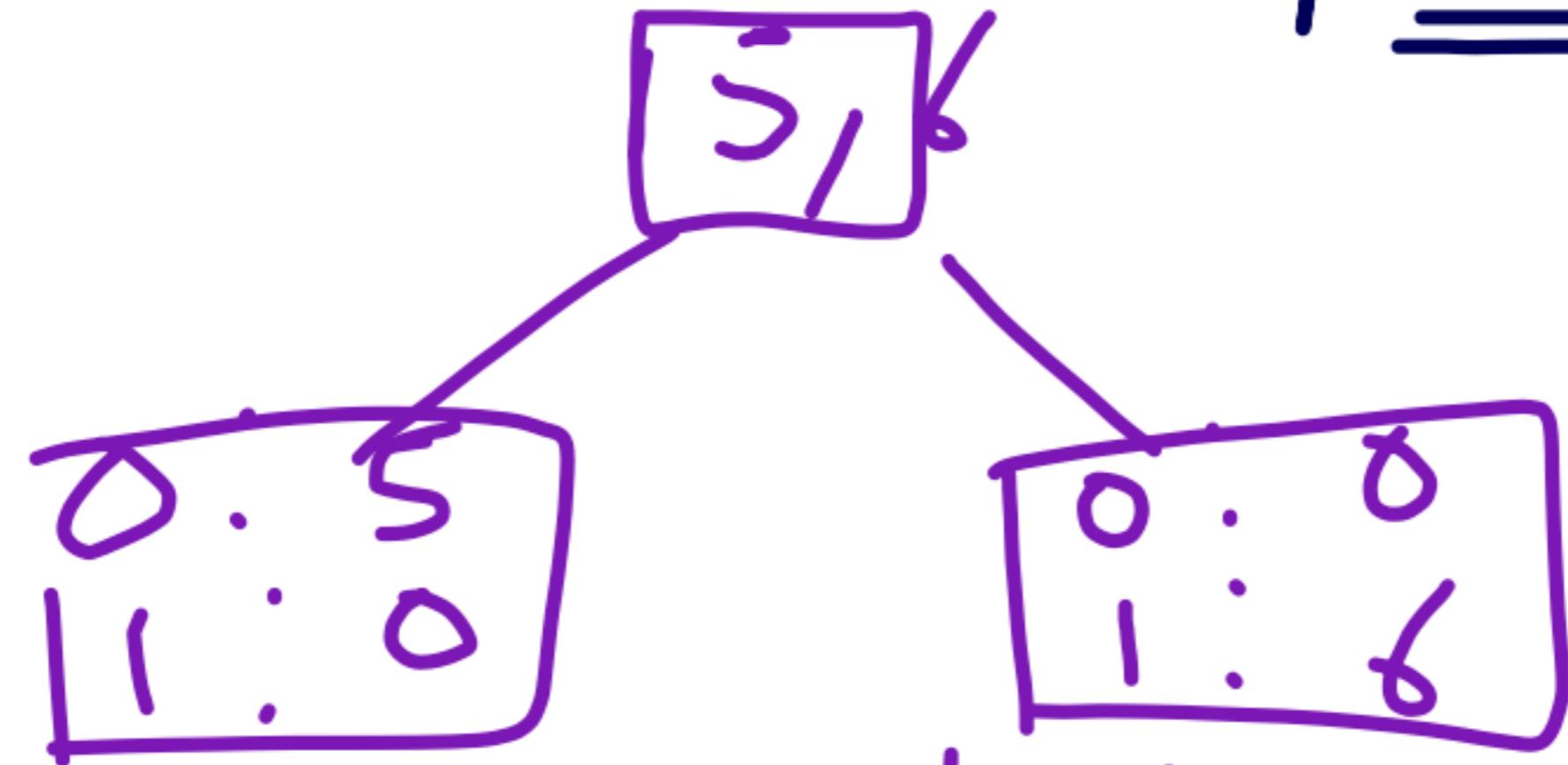
Node Impurity



- Tree begins with single node whose label reflects the majority class value
  - Tree needs to be refined because root node contains records from both classes
  - Divide records recursively into smaller subsets
-

List two cases when splitting is not possible / required

Case 1:



Both nodes  
are pure

Case 2:

$x_1$	$x_2$	$y$
1	1	0
1	1	0
1	1	0
1	1	0

# Design Issues of Decision Tree Induction

## I How should training records be split?

- Greedy strategy: split the records based on some attribute test (always choose immediate best option)
- Need to evaluate the “goodness” of each attribute test and select the best one.
  - Method for specifying test condition
    - ◆ depending on attribute types
  - Measure for evaluating the goodness of a test condition

## How should the splitting procedure stop?

- Stop splitting if all the records belong to the same class or have identical attribute values
- Early termination

20

→ Depth

# Methods for Expressing Test Conditions

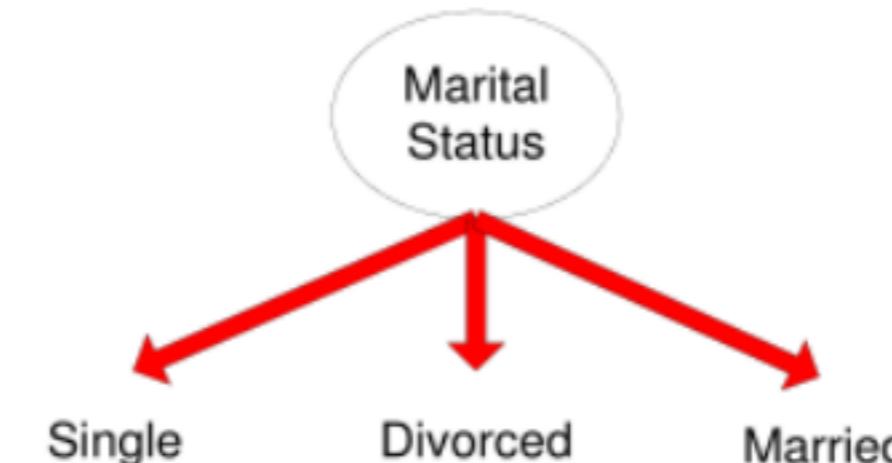
---

- | Depends on attribute types
  - Binary (**Test condition only generates two potential outcomes**)
  - Nominal
  - Ordinal
  - Continuous
  
- | Depends on number of ways to split
  - 2-way split
  - Multi-way split

# Test Condition for Nominal Attributes

- **Multi-way split:**

- Use as many partitions as distinct values.



- **Binary split:**

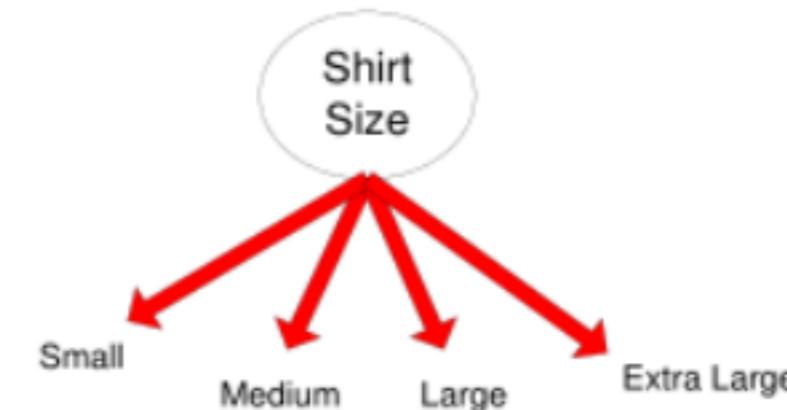
- Divides values into two subsets



# Test Condition for Ordinal Attributes

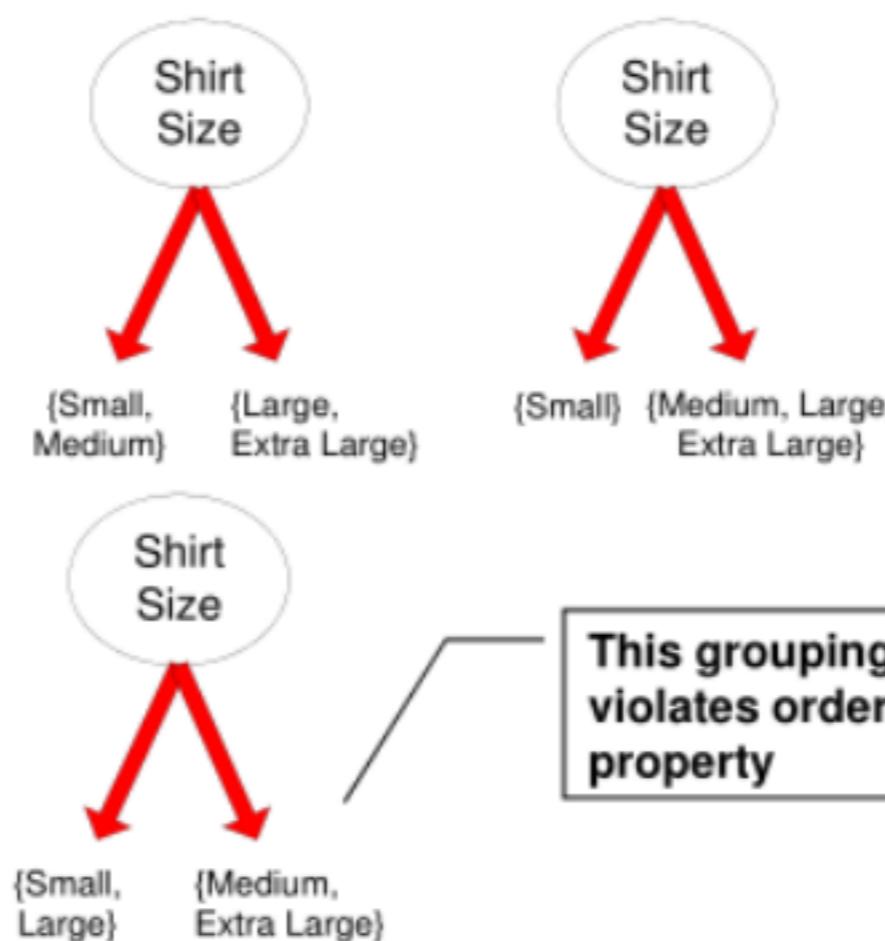
## | Multi-way split:

- Use as many partitions as distinct values

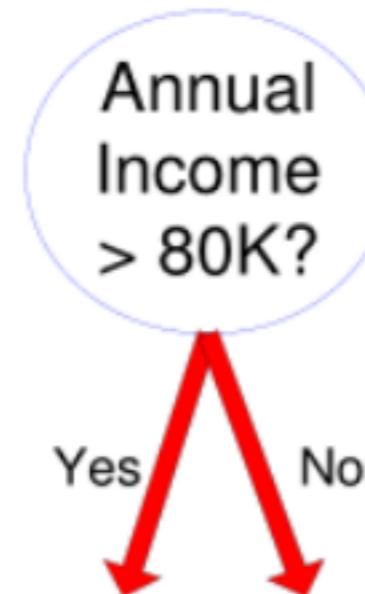


## | Binary split:

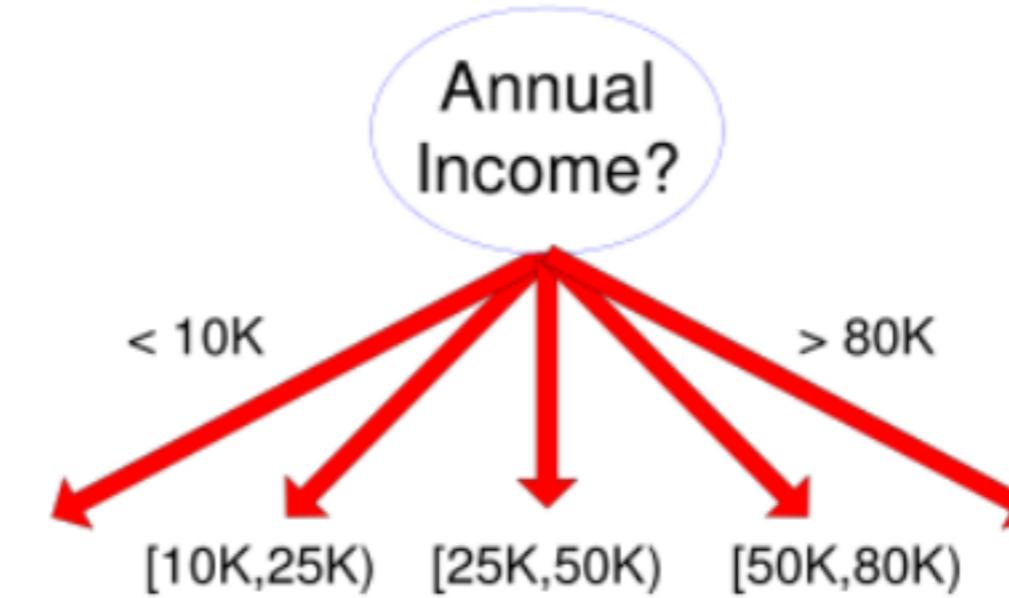
- Divides values into two subsets
- Preserve order property among attribute values



# Test Condition for Continuous Attributes



(i) Binary split



(ii) Multi-way split

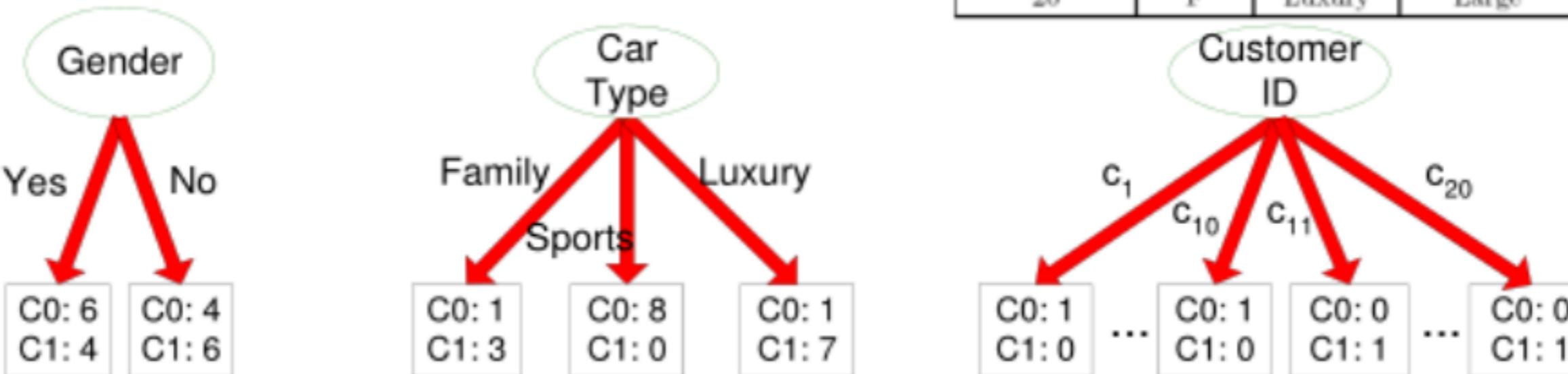
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

# How to determine the Best Split

Customer Id	Gender	Car Type	Shirt Size	Class
-------------	--------	----------	------------	-------

1	M	Family	Small	C0
2	M	Sports	Medium	C0
3	M	Sports	Medium	C0
4	M	Sports	Large	C0
5	M	Sports	Extra Large	C0
6	M	Sports	Extra Large	C0
7	F	Sports	Small	C0
8	F	Sports	Small	C0
9	F	Sports	Medium	C0
10	F	Luxury	Large	C0
11	M	Family	Large	C1
12	M	Family	Extra Large	C1
13	M	Family	Medium	C1
14	M	Luxury	Extra Large	C1
15	F	Luxury	Small	C1
16	F	Luxury	Small	C1
17	F	Luxury	Medium	C1
18	F	Luxury	Medium	C1
19	F	Luxury	Medium	C1
20	F	Luxury	Large	C1

**Before Splitting:** 10 records of class 0,  
10 records of class 1



**Which test condition is the best?**

"\ Test condition that results in a larger no. of outcomes may not be desirable bcz no. of records associated with each partition is too small to enable us to make any reliable predictions";

# How to determine the Best Split

---

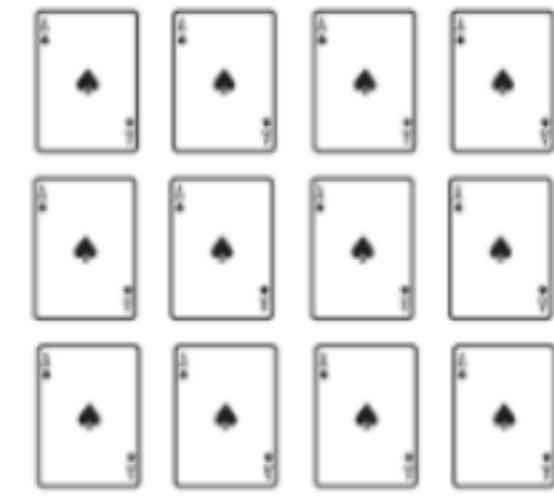
- | Greedy approach:
  - Nodes with **purer** class distribution are preferred
- | Need a measure of node impurity:

C0: 5
C1: 5

High degree of impurity

C0: 9
C1: 1

Low degree of impurity



Little impure.



Very pure. Not impure.



Completely impure.

## Measures of Node Impurity (Node t)

### | Gini Index

$$Gini\ Index = 1 - \sum_{i=0}^{c-1} p_i(t)^2$$

Where  $p_i(t)$  is the frequency of class  $i$  at node  $t$ , and  $c$  is the total number of classes

### | Entropy

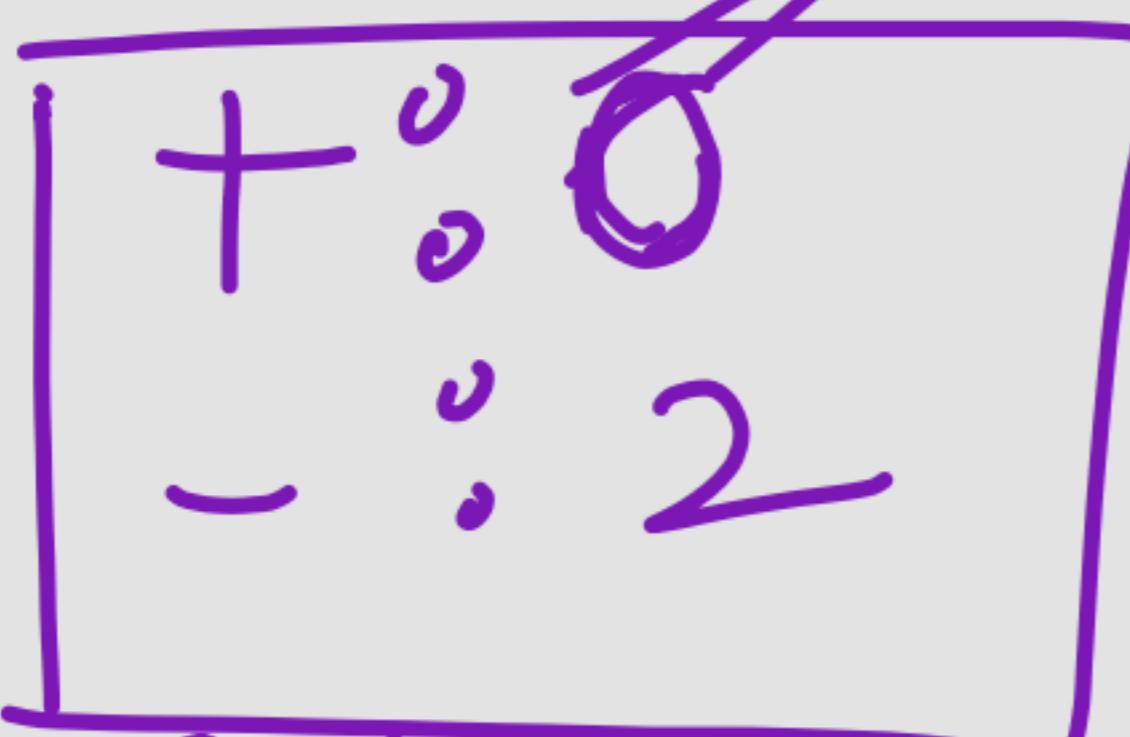
$$Entropy = - \sum_{i=0}^{c-1} p_i(t) \log_2 p_i(t)$$

### | Misclassification error

$$Classification\ error = 1 - \max[p_i(t)]$$

Measures The amount of uncertainty of random Variable

$$\sum p_i \log p_i$$



$$-\frac{1}{2} \log_2 \frac{1}{2} - \frac{2}{2} \log \frac{2}{2}$$

$$-\frac{5}{6} \log_2 \frac{5}{6} - \frac{1}{6} \log \frac{1}{6}$$

Determine the entropy in each of the following cases.

Node $N_1$	Count
Class=0	0
Class=1	6

Node $N_2$	Count
Class=0	1
Class=1	5

Node $N_3$	Count
Class=0	3
Class=1	3

$$\begin{aligned} & \rightarrow -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \\ & - \frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \\ \Rightarrow & \quad \frac{1}{2} + \frac{1}{2} = 1 \end{aligned}$$

Node $N_1$	Count
Class=0	0
Class=1	6

$$\text{Gini} = 1 - (0/6)^2 - (6/6)^2 = 0$$

$$\text{Entropy} = -(0/6) \log_2(0/6) - (6/6) \log_2(6/6) = 0$$

$$\text{Error} = 1 - \max[0/6, 6/6] = 0$$

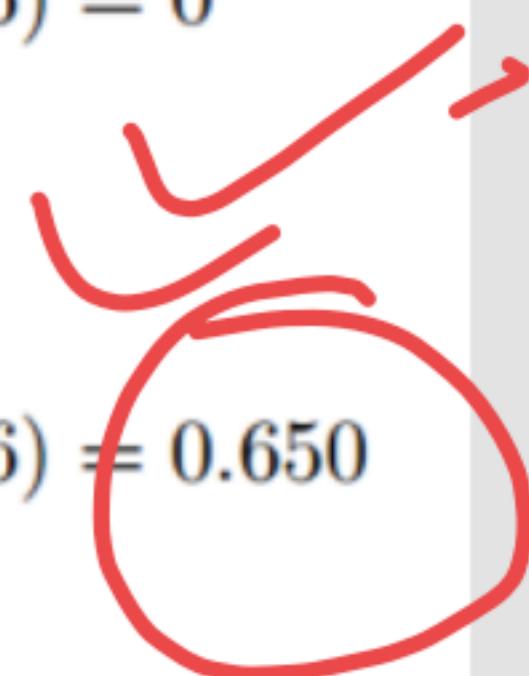
•

Node $N_2$	Count
Class=0	1
Class=1	5

$$\text{Gini} = 1 - (1/6)^2 - (5/6)^2 = 0.278$$

$$\text{Entropy} = -(1/6) \log_2(1/6) - (5/6) \log_2(5/6) = 0.650$$

$$\text{Error} = 1 - \max[1/6, 5/6] = 0.167$$

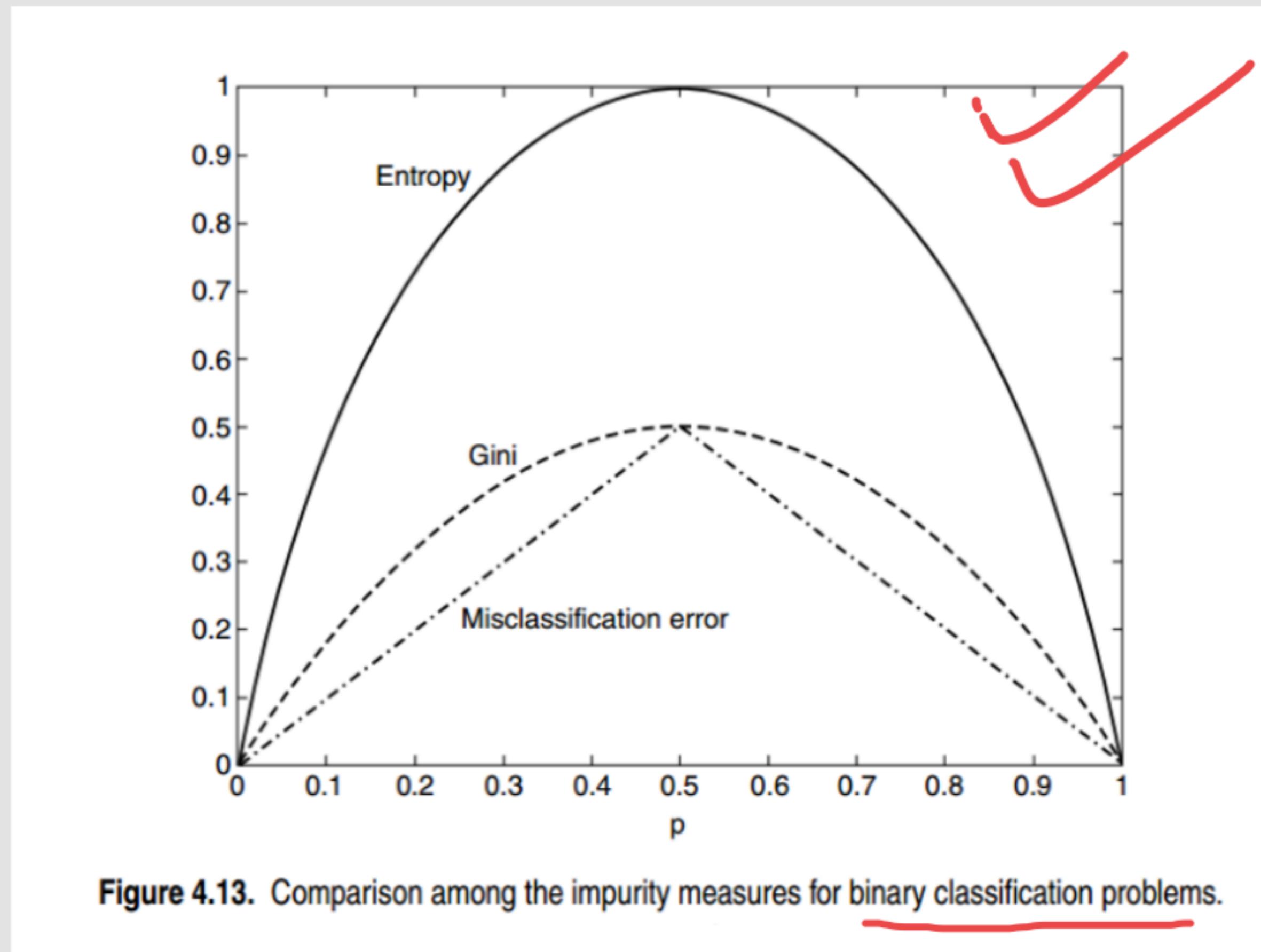


Node $N_3$	Count
Class=0	3
Class=1	3

$$\text{Gini} = 1 - (3/6)^2 - (3/6)^2 = 0.5$$

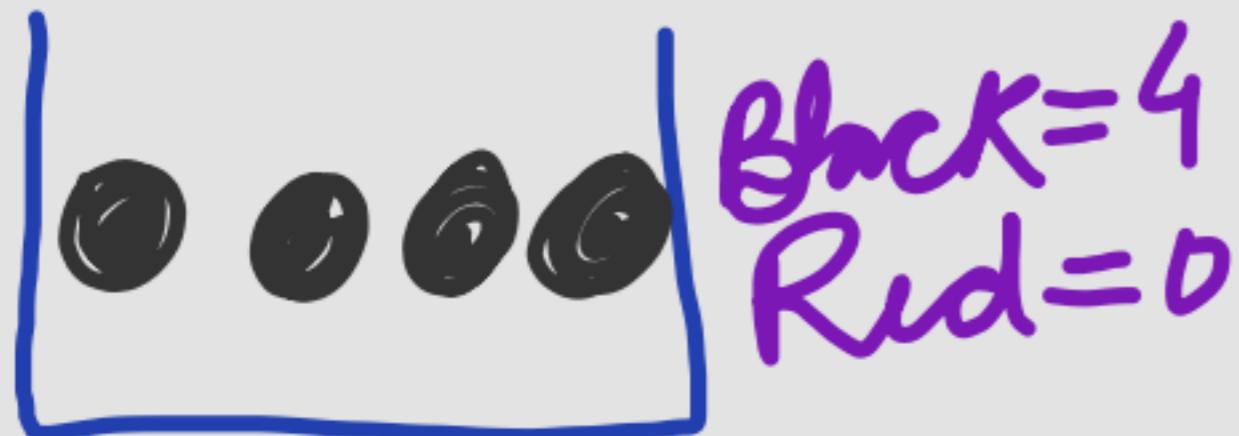
$$\text{Entropy} = -(3/6) \log_2(3/6) - (3/6) \log_2(3/6) = 1$$

$$\text{Error} = 1 - \max[3/6, 3/6] = 0.5$$



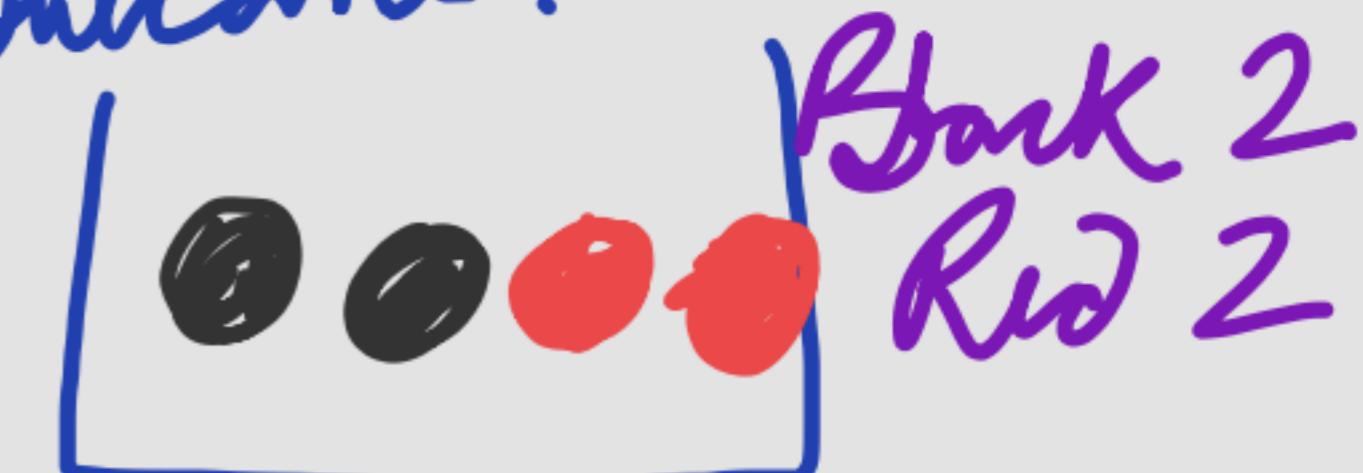
## A. Additional Slide

Entropy - Derived from information theory  
Shanon Entropy : No. of bits needed to encode the outcome.



$I = 0$  (Pure)

No. of bits needed to encode the outcome



$I = 1$  (Impure)

"If we are less certain about the outcome, we may need more info to describe it"

# Finding the Best Split

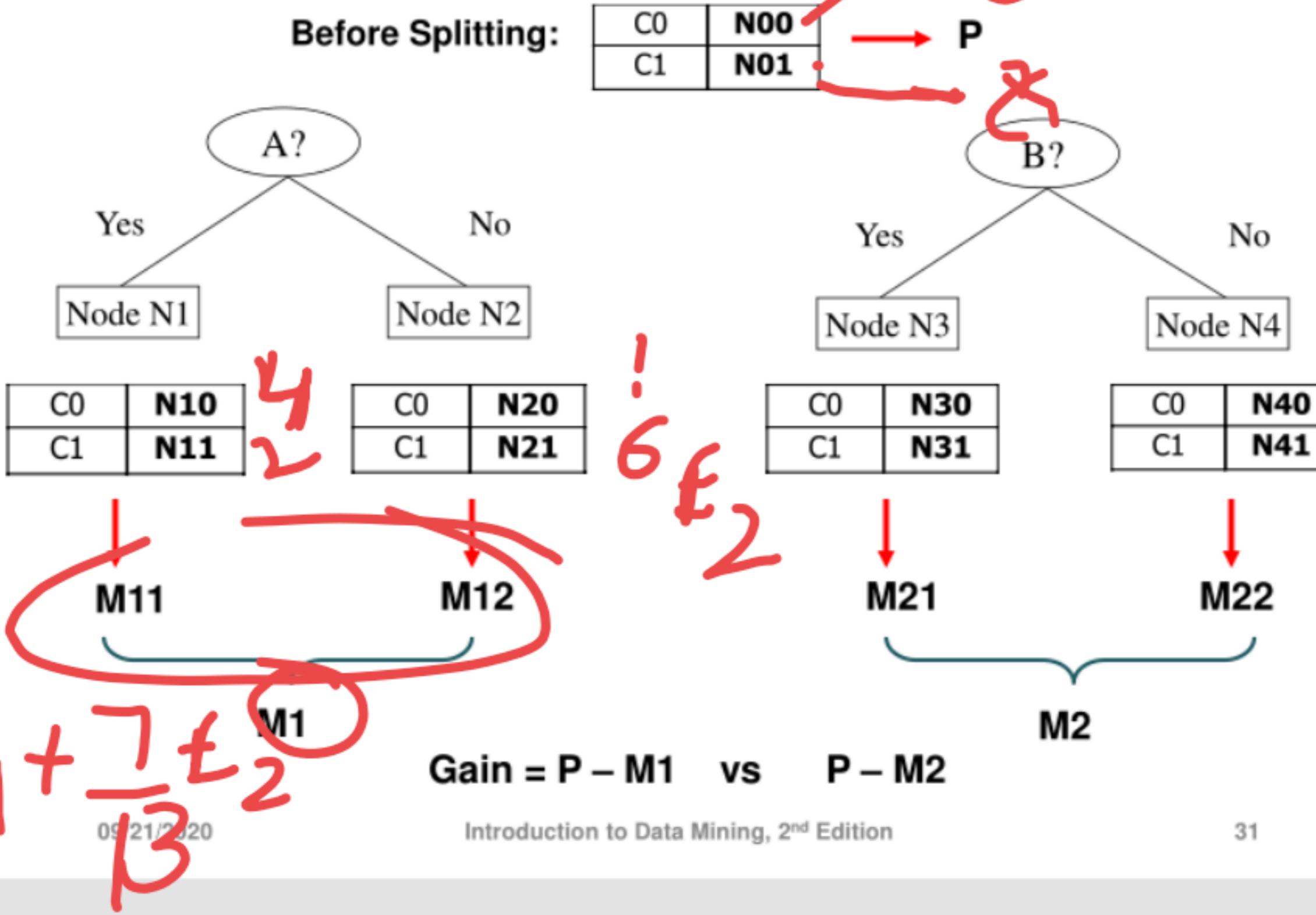
---

1. Compute impurity measure ( $P$ ) before splitting
2. Compute impurity measure ( $M$ ) after splitting
  - |A. Compute impurity measure of each child node
  - |B.  $M$  is the weighted impurity of child nodes
3. Choose the attribute test condition that produces the highest gain

$$\text{Gain} = P - M$$

or equivalently, lowest impurity measure after splitting ( $M$ )

# Finding the Best Split



## Gain

- Gain: "goodness of the split"
- Comparing:
  - degree of impurity of parent node (before splitting)
  - degree of impurity of the child nodes (after splitting), weighted
- larger gain => better split (better test condition)

$$\Delta(\text{gain}) = I(\text{parent}) - \sum_{j=1}^k \frac{N(v_j)}{N} I(v_j)$$

- $I(n)$  = impurity measure at node  $n$
- $I(v_i)$  = impurity measure at child node  $v_i$
- $k$  = number of attribute values
- $N(v_i)$  = total number of records at child node  $v_i$
- $N$  = total number of records at parent node

Footnote: [Information Gain: term used when entropy is used as the impurity measure]

# Splitting of Continuous Attributes

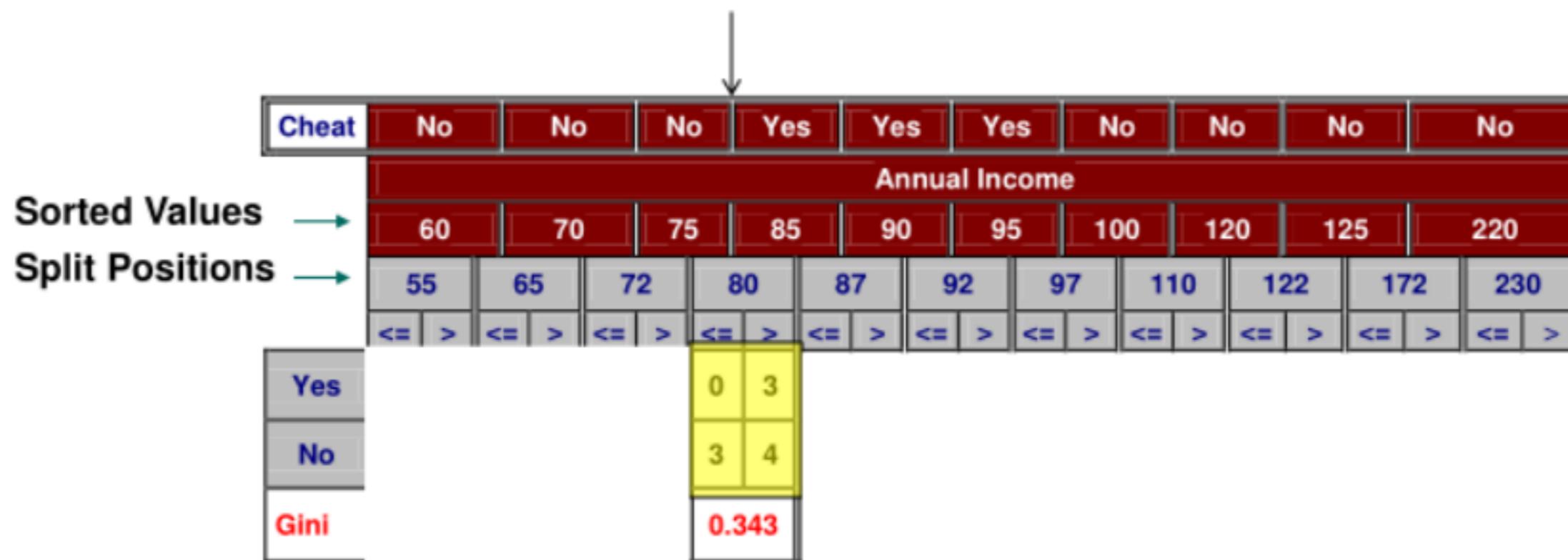
- By sorting records by the continuous attribute, this improves to  $O(n \log n)$ 
  - Candidate Split position are midpoints between two adjacent, different, class values
  - Only need to consider split positions: 80 and 97

Tid	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Class	No	No	No	Yes	Yes	Yes	No	No	No	No
	Annual Income									
Sorted Values →										
Split Positions →	60	70	75	85	90	95	100	120	125	220
	<=	>	<=	>	<=	>	<=	>	<=	>
Yes	0	3	0	3	0	3	1	2	2	0
No	0	7	1	6	2	5	3	4	3	4
Gini	0.420	0.400	0.375	0.343	0.417	0.400	0.300	0.343	0.375	0.400

## Continuous Attributes: Computing Gini Index...

- | For efficient computation: for each attribute,
  - Sort the attribute on values
  - Linearly scan these values, each time updating the count matrix and computing gini index
  - Choose the split position that has the least gini index



---

**Algorithm 4.1** A skeleton decision tree induction algorithm.

---

**TreeGrowth** ( $E, F$ )

- 1: **if**  $\text{stopping\_cond}(E, F) = \text{true}$  **then**
- 2:    $\text{leaf} = \text{createNode}()$ .
- 3:    $\text{leaf.label} = \text{Classify}(E)$ .
- 4:   **return**  $\text{leaf}$ .
- 5: **else**
- 6:    $\text{root} = \text{createNode}()$ .
- 7:    $\text{root.test\_cond} = \text{find\_best\_split}(E, F)$ .
- 8:   let  $V = \{v | v \text{ is a possible outcome of } \text{root.test\_cond}\}$ .
- 9:   **for** each  $v \in V$  **do**
- 10:      $E_v = \{e \mid \text{root.test\_cond}(e) = v \text{ and } e \in E\}$ .
- 11:      $\text{child} = \text{TreeGrowth}(E_v, F)$ .
- 12:     add  $\text{child}$  as descendent of  $\text{root}$  and label the edge  $(\text{root} \rightarrow \text{child})$  as  $v$ .
- 13:   **end for**
- 14: **end if**
- 15: **return**  $\text{root}$ .

---