

## Classification: Decision Tree

CSE 4334 / 5334 Data Mining  
Spring 2019

**Won Hwa Kim**

(Slides courtesy of Pang-Ning Tan, Michael Steinbach and Vipin Kumar, and Jiawei Han, Micheline Kamber and Jian Pei)



## Classification: Definition



Given a collection of records (*training set*)

- Each record contains a set of *attributes*, one of the attributes is the *class*.

Find a *model* for class attribute as a function of the values of other attributes.

Goal: previously unseen records should be assigned a class as accurately as possible.

- A *test set* is used to determine the accuracy of the model. Usually, the given data set is divided into training and test sets, with training set used to build the model and test set used to validate it.

2

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Illustrating Classification Task

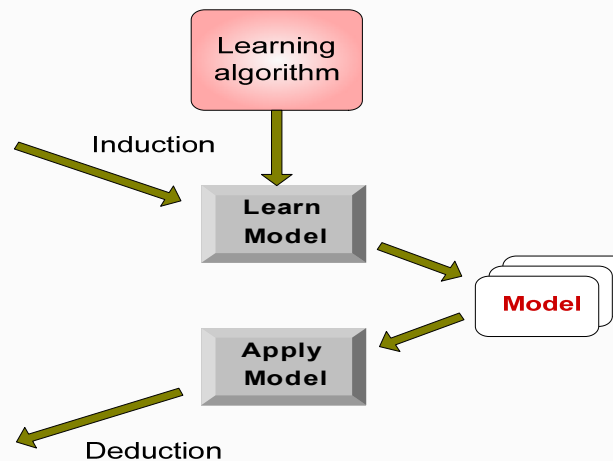


Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



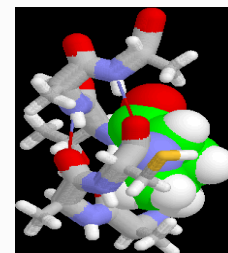
3

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Examples of Classification Task



- Predicting tumor cells as benign or malignant
- Classifying credit card transactions as legitimate or fraudulent
- Classifying secondary structures of protein as alpha-helix, beta-sheet, or random coil
- Categorizing news stories as finance, weather, entertainment, sports, etc



4

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.



# Classification vs. Prediction

## Classification

- Predicts categorical class labels
- Most suited for nominal attributes
- Less effective for ordinal attributes

## Prediction

- models continuous-valued functions or ordinal attributes, i.e., predicts unknown or missing values
- E.g., Linear regression

5

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.



# Supervised vs. Unsupervised Learning

## Supervised learning (e.g., classification)

- Supervision: The training data (observations, measurements, etc.) are accompanied by labels indicating the class of the observations
- New data is classified based on the training set

## Unsupervised learning (e.g., clustering)

- The class labels of training data is unknown
- Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

6

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Classification Techniques



- Decision Tree based Methods
- Rule-based Methods
- Memory based reasoning
- Neural Networks
- Naïve Bayes and Bayesian Belief Networks
- Support Vector Machines

7

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

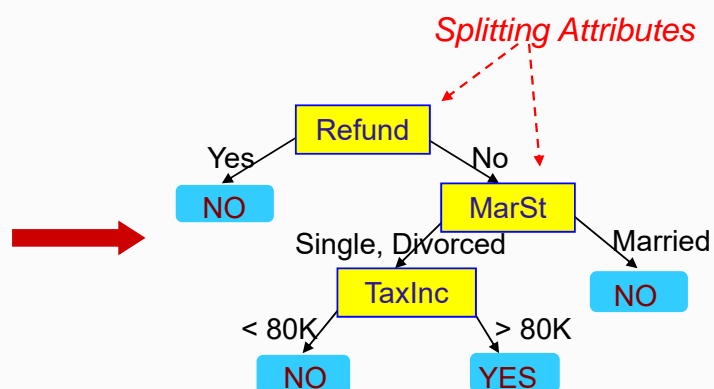
## Example of a Decision Tree



*categorical*  
*categorical*  
*continuous*  
*class*

Tid	Refund	Marital Status	Taxable Income	Cheat
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

Training Data



Model: Decision Tree

8

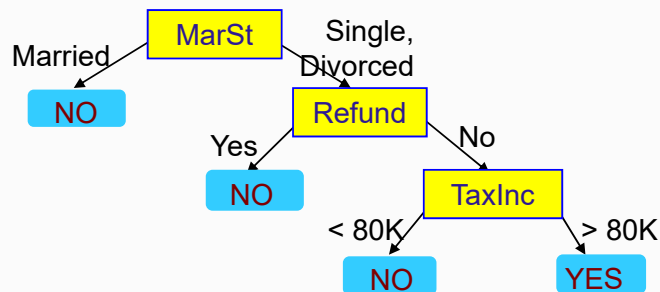
Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Another Example of Decision Tree



*categorical*  
*categorical*  
*continuous*  
*class*

Tid	Refund	Marital Status	Taxable Income	Cheat
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



There could be more than one tree that fits the same data!

9

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Decision Tree Classification Task

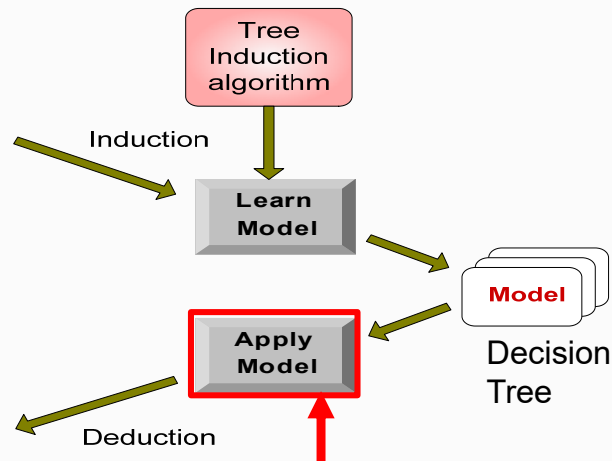


Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Training Set

Refund	Marital Status	Taxable Income	Cheat
No	Single	75K	?
Yes	Married	50K	?
No	Married	150K	?
Yes	Divorced	90K	?
No	Single	40K	?
No	Married	80K	?

Test Set



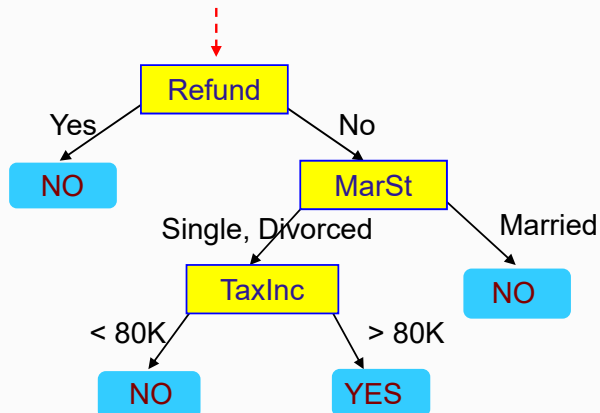
10

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Apply Model to Test Data



Start from the root of tree.



Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

11

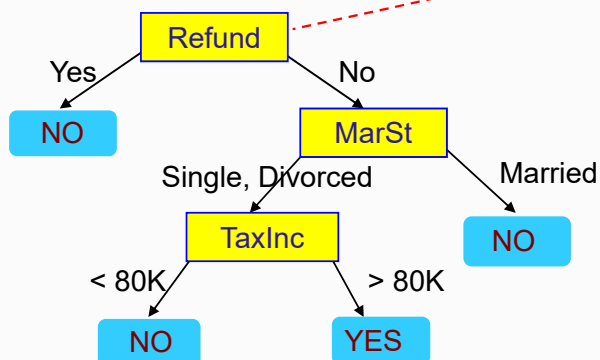
Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Apply Model to Test Data



Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



12

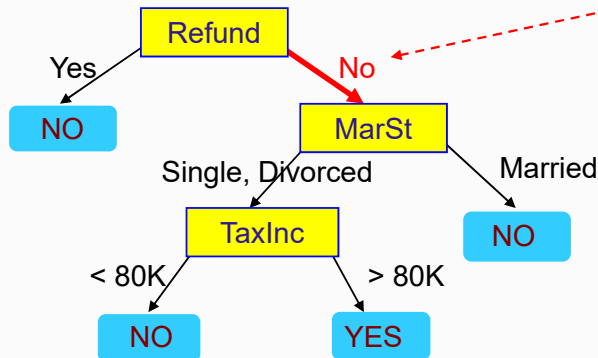
Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Apply Model to Test Data



Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



13

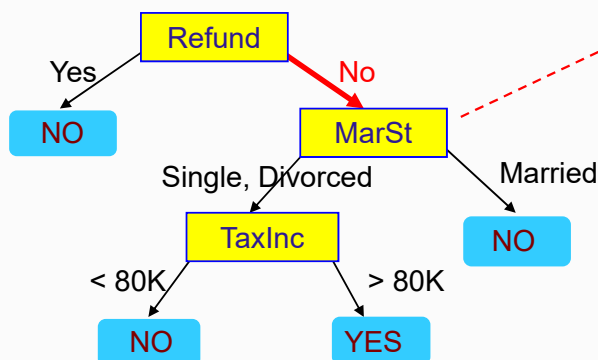
Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Apply Model to Test Data



Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



14

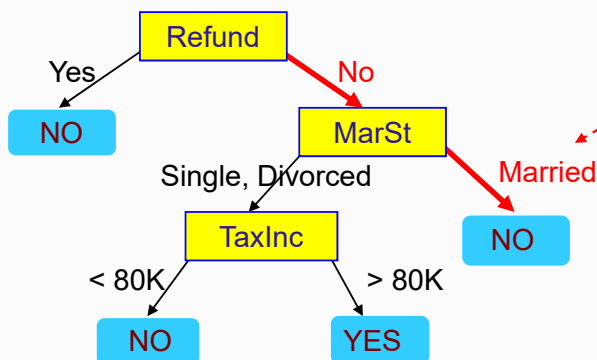
Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Apply Model to Test Data



Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



15

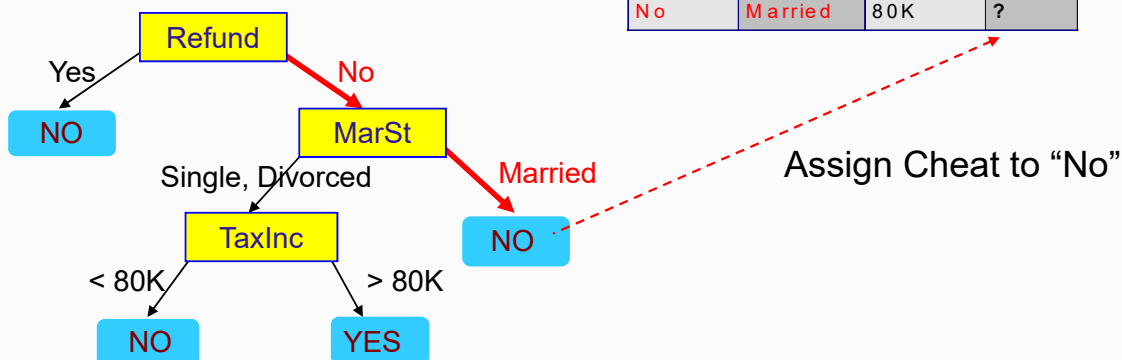
Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Apply Model to Test Data



Test Data

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



16

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.



# Decision Tree Classification Task

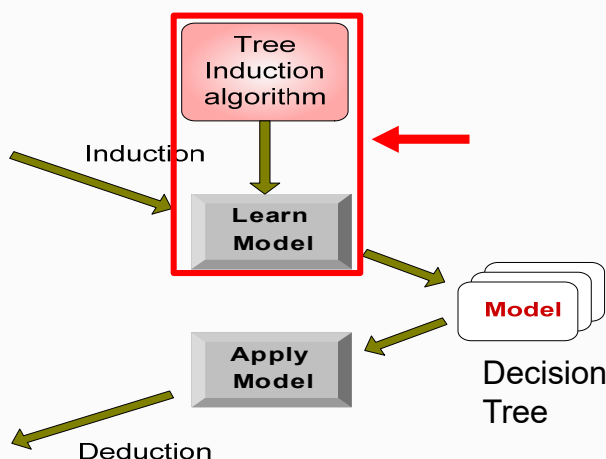


Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Training Set

Refund	Marital Status	Taxable Income	Cheat
No	Single	75K	?
Yes	Married	50K	?
No	Married	150K	?
Yes	Divorced	90K	?
No	Single	40K	?
No	Married	80K	?

Test Set



17

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Decision Tree Induction



### Large search space

- Exponential size, with respect to the set of attributes
- Finding the optimal decision tree is computationally infeasible

### Efficient algorithm for accurate suboptimal decision tree

- Greedy strategy
- Grow the tree by making locally optimally decisions in selecting the attributes

18

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Decision Tree Induction



## Many Algorithms:

- Hunt's Algorithm (one of the earliest)
- CART
- ID3, C4.5
- SLIQ, SPRINT

19

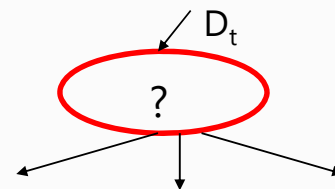
Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# 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 the same class  $y_t$ , then  $t$  is a leaf node labeled as  $y_t$ .
  - If  $D_t$  is an empty set, then  $t$  is a leaf node labeled by the majority class among the records of  $D_t$ 's parent node.
  - If  $D_t$  contains records that have identical values on all attributes but the class attribute, then  $t$  is a leaf node labeled by the majority class among  $D_t$ 's records.
  - If none of the above conditions is satisfied, use **an attribute test** to split the data into smaller subsets. Recursively apply the procedure to each subset.

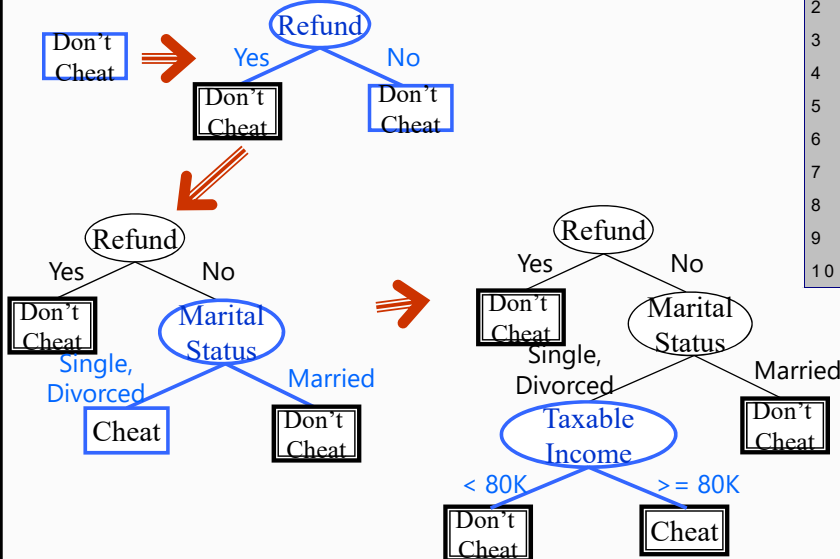
Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



20

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Hunt's Algorithm



Tid	Refund	Marital Status	Taxable Income	Cheat
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

21

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Tree Induction

### Greedy strategy.

- Split the records based on an attribute test that optimizes certain criterion.

### Issues

- Determine how to split the records
  - How to specify the attribute test condition?
  - How to determine the best split?
- Determine when to stop splitting

22

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Tree Induction



## Greedy strategy.

- Split the records based on an attribute test that optimizes certain criterion.

## Issues

- Determine how to split the records
  - How to specify the attribute test condition?
  - How to determine the best split?
- Determine when to stop splitting

23

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# How to Specify Test Condition?



## Depends on attribute types

- Categorical vs. Numeric
  - Categorical attributes: Nominal, Ordinal
  - Numeric attributes: Interval, Ratio
- Discrete vs. Continuous

## Depends on number of ways to split

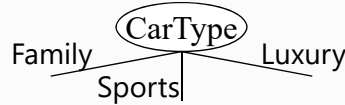
- 2-way split
- Multi-way split

24

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Splitting Based on Nominal Attributes

**Multi-way split:** Use as many partitions as distinct values.



**Binary split:** Divides values into two subsets.  
Need to find optimal partitioning.

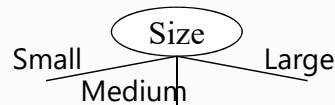


25

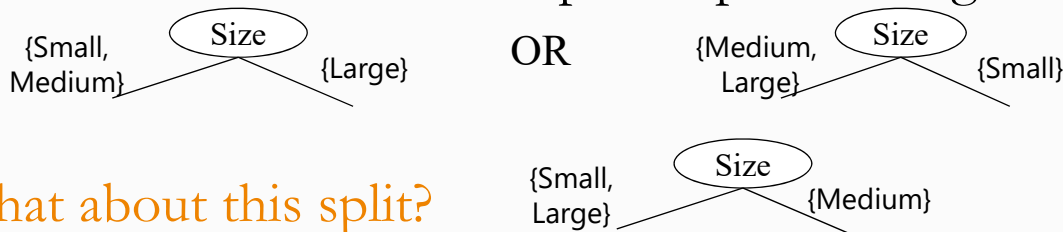
Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Splitting Based on Ordinal Attributes

**Multi-way split:** Use as many partitions as distinct values.



**Binary split:** Divides values into two subsets.  
Need to find optimal partitioning.



What about this split?

26

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Splitting Based on Continuous Attribute

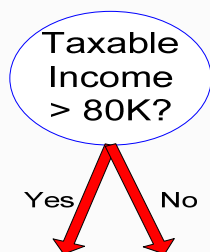
## Different ways of handling

- **Discretization** to form an ordinal categorical attribute
  - Static – discretize once at the beginning
  - Dynamic – ranges can be found by equal interval bucketing, equal frequency bucketing (percentiles), or clustering.
- **Binary Decision:**  $(A < v)$  or  $(A \geq v)$ 
  - consider all possible splits and finds the best cut
  - can be more compute intensive

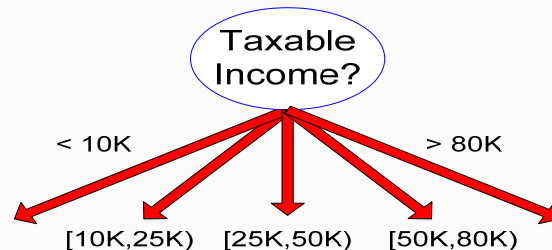
27

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Splitting Based on Continuous Attribute



(i) Binary split



(ii) Multi-way split

28

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Tree Induction



## Greedy strategy.

- o Split the records based on an attribute test that optimizes certain criterion.

## Issues

- o Determine how to split the records
  - How to specify the attribute test condition?
  - **How to determine the best split?**
- o Determine when to stop splitting

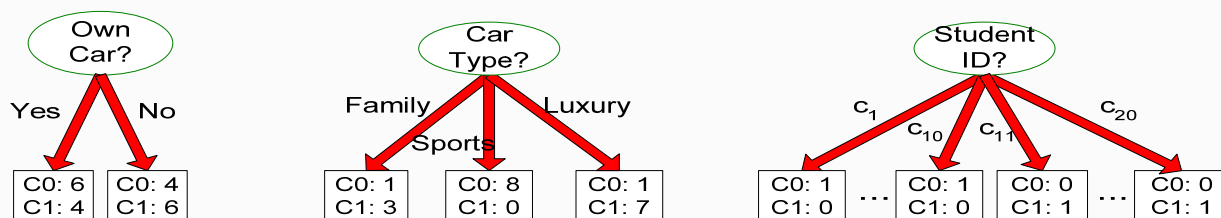
29

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# How to determine the Best Split



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



Which test condition is the best?

30

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# How to determine the Best Split



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

C0: 5
C1: 5

Non-homogeneous,  
High degree of impurity

C0: 9
C1: 1

Homogeneous,  
Low degree of impurity

31

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Measures of Node Impurity



Gini Index

Entropy

Misclassification error

32

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.



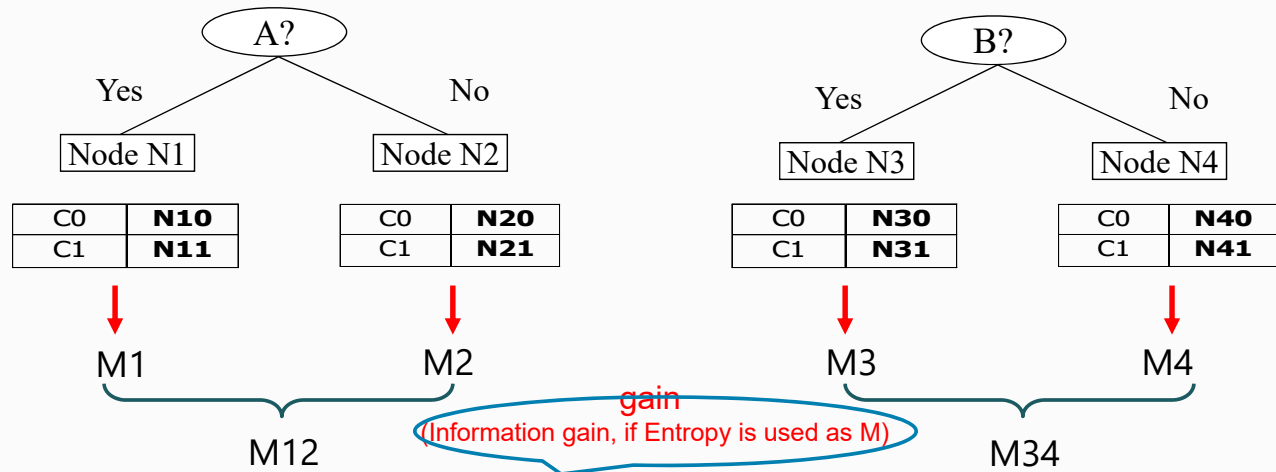
# How to Find the Best Split



Before Splitting:

C0	<b>N00</b>
C1	<b>N01</b>

→ M0



Gain = M0 – M12 vs M0 – M34

33

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Measure of Impurity: GINI



Gini Index for a given node t :

$$GINI(t) = 1 - \sum_j [p(j | t)]^2$$

(NOTE:  $p(j | t)$  is the relative frequency of class j at node t).

- Maximum  $(1 - 1/n_c)$  when records are equally distributed among all classes, implying least interesting information
- Minimum (0.0) when all records belong to one class, implying most interesting information

C1	<b>0</b>
C2	<b>6</b>
<b>Gini= 0.000</b>	

C1	<b>1</b>
C2	<b>5</b>
<b>Gini= 0.278</b>	

C1	<b>2</b>
C2	<b>4</b>
<b>Gini= 0.444</b>	

C1	<b>3</b>
C2	<b>3</b>
<b>Gini= 0.500</b>	

34

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Examples for computing GINI



$$GINI(t) = 1 - \sum_j [p(j|t)]^2$$

C 1	<b>0</b>
C 2	<b>6</b>

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$Gini = 1 - P(C1)^2 - P(C2)^2 = 1 - 0 - 1 = 0$$

C 1	<b>1</b>
C 2	<b>5</b>

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$Gini = 1 - (1/6)^2 - (5/6)^2 = 0.278$$

C 1	<b>2</b>
C 2	<b>4</b>

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$Gini = 1 - (2/6)^2 - (4/6)^2 = 0.444$$

35

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Splitting Based on GINI



- Used in CART, SLIQ, SPRINT.
- When a node p is split into k partitions (children), the quality of split is computed as,

$$GINI_{split} = \sum_{i=1}^k \frac{n_i}{n} GINI(i)$$

where,  $n_i$  = number of records at child i,  
 $n$  = number of records at node p.

36

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Binary Attributes: Computing GINI Index



- Splits into two partitions
- Effect of Weighing partitions:
  - Larger and Purer Partitions are sought for.

Decision Tree Split: B?

```

    B?
   /  \
 Yes    No
Node N1 Node N2
  
```

Parent Table:

	Parent
C 1	6
C 2	6
<b>Gini = 0.500</b>	

Node N1 Table:

	N1	N2
C1	5	1
C2	2	4
<b>Gini=0.371</b>		

Node N2 Table:

	N1	N2
C1	5	1
C2	2	4
<b>Gini=0.371</b>		

Gini(N1) calculation:

$$\begin{aligned} \text{Gini}(N1) &= 1 - (5/7)^2 - (2/7)^2 \\ &= 0.408 \end{aligned}$$

Gini(N2) calculation:

$$\begin{aligned} \text{Gini}(N2) &= 1 - (1/5)^2 - (4/5)^2 \\ &= 0.32 \end{aligned}$$

Gini(Children) calculation:

$$\begin{aligned} \text{Gini(Children)} &= 7/12 * 0.408 + \\ &\quad 5/12 * 0.32 \\ &= 0.371 \end{aligned}$$

37

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Categorical Attributes: Computing Gini Index



- For each distinct value, gather counts for each class in the dataset
- Use the count matrix to make decisions

Multi-way split

	CarType		
	Family	Sports	Luxury
C 1	1	2	1
C 2	4	1	1
<b>Gini</b>	<b>0.393</b>		

Two-way split (find best partition of values)

	CarType	
	{Sports, Luxury}	{Family}
C 1	3	1
C 2	2	4
<b>Gini</b>	<b>0.400</b>	

	CarType	
	{Sports}	{Family, Luxury}
C 1	2	2
C 2	1	5
<b>Gini</b>	<b>0.419</b>	

38

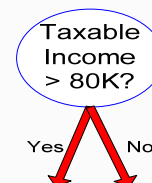
Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Continuous Attributes: Computing Gini Index



- Use Binary Decisions based on one value
- Several Choices for the splitting value
  - Number of possible splitting values  
= Number of distinct values
- Each splitting value has a count matrix associated with it
  - Class counts in each of the partitions,  $A < v$  and  $A \geq v$
- Simple method to choose best  $v$ 
  - For each  $v$ , scan the database to gather count matrix and compute its Gini index
  - Computationally Inefficient! Repetition of work.

Tid	Refund	Marital Status	Taxable Income	Cheat
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



39

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## 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

		Cheat	No		No		No		Yes		Yes		Yes		No		No		No		No			
			Taxable Income																					
			60		70		75		85		90		95		100		120		125		220			
Sorted Values Split Positions	→		55		65		72		80		87		92		97		110		122		172		230	
			<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>
		Yes	0	3	0	3	0	3	0	3	1	2	2	1	3	0	3	0	3	0	3	0	3	0
		No	0	7	1	6	2	5	3	4	3	4	3	4	3	4	4	3	5	2	6	1	7	0
		Gini	0.420		0.400		0.375		0.343		0.417		0.400		<u>0.300</u>		0.343		0.375		0.400		0.420	

40

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Alternative Splitting Criteria based on INFO



### Entropy at a given node t:

$$\text{Entropy}(t) = - \sum_j p(j | t) \log p(j | t)$$

(NOTE:  $p(j | t)$  is the relative frequency of class  $j$  at node  $t$ ).

- Measures homogeneity of a node.
  - Maximum ( $\log n_c$ ) when records are equally distributed among all classes implying least information
  - Minimum (0.0) when all records belong to one class, implying most information
- Entropy based computations are similar to the GINI index computations

41

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Examples for computing Entropy



$$\text{Entropy}(t) = - \sum_j p(j | t) \log_2 p(j | t)$$

C 1	<b>0</b>
C 2	<b>6</b>

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$\text{Entropy} = - 0 \log_2 0 - 1 \log_2 1 = - 0 - 0 = 0$$

C 1	<b>1</b>
C 2	<b>5</b>

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$\text{Entropy} = - (1/6) \log_2 (1/6) - (5/6) \log_2 (5/6) = 0.65$$

C 1	<b>2</b>
C 2	<b>4</b>

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$\text{Entropy} = - (2/6) \log_2 (2/6) - (4/6) \log_2 (4/6) = 0.92$$

42

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.



## Why is that $0 \log 0 = 0$ ?

$$\lim_{x \rightarrow 0} x \log_2(x) = \lim_{x \rightarrow 0} \frac{\ln(x)}{\ln(2)} \cdot \frac{x^{-1}}{x^{-1}} = \lim_{x \rightarrow 0} \frac{\ln(x)}{-x^{-2}} = \lim_{x \rightarrow 0} \frac{-x}{\ln(2)} = 0$$

L'Hospital's Rule (Wikipedia)

If

$$\begin{aligned} \lim_{x \rightarrow c} f(x) = \lim_{x \rightarrow c} g(x) = 0 \text{ or } \pm \infty, \text{ and} \\ \lim_{x \rightarrow c} \frac{f'(x)}{g'(x)} \text{ exists, and} \\ g'(x) \neq 0 \text{ for all } x \text{ in } I \text{ with } x \neq c, \end{aligned}$$

then

$$\lim_{x \rightarrow c} \frac{f(x)}{g(x)} = \lim_{x \rightarrow c} \frac{f'(x)}{g'(x)}.$$

43

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Splitting Based on INFO...



### Information Gain:

$$GAIN_{split} = Entropy(p) - \left( \sum_{i=1}^k \frac{n_i}{n} Entropy(i) \right)$$

Parent Node, p is split into k partitions;

$n_i$  is number of records in partition i

- Measures Reduction in Entropy achieved because of the split. Choose the split that achieves most reduction (maximizes GAIN)
- Used in ID3 and C4.5
- Disadvantage: Tends to prefer splits that result in large number of partitions, each being small but pure.

44

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Splitting Based on INFO...



### Gain Ratio:

$$GainRatio_{split} = \frac{GAIN_{Split}}{SplitINFO}$$

$$SplitINFO = - \sum_{i=1}^k \frac{n_i}{n} \log \frac{n_i}{n}$$

Parent Node, p is split into k partitions

$n_i$  is the number of records in partition i

- Adjusts Information Gain by the entropy of the partitioning (SplitINFO). Higher entropy partitioning (large number of small partitions) is penalized!
- Used in C4.5
- Designed to overcome the disadvantage of Information Gain

45

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Splitting Criteria based on Classification Error



### Classification error at a node t :

$$Error(t) = 1 - \max_i P(i | t)$$

Measures misclassification error made by a node.

- Maximum ( $1 - 1/n_c$ ) when records are equally distributed among all classes, implying least interesting information
- Minimum (0.0) when all records belong to one class, implying most interesting information

46

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Examples for Computing Error



$$\text{Error}(t) = 1 - \max_i P(i | t)$$

C 1	<b>0</b>
C 2	<b>6</b>

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$\text{Error} = 1 - \max(0, 1) = 1 - 1 = 0$$

C 1	<b>1</b>
C 2	<b>5</b>

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$\text{Error} = 1 - \max(1/6, 5/6) = 1 - 5/6 = 1/6$$

C 1	<b>2</b>
C 2	<b>4</b>

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$\text{Error} = 1 - \max(2/6, 4/6) = 1 - 4/6 = 1/3$$

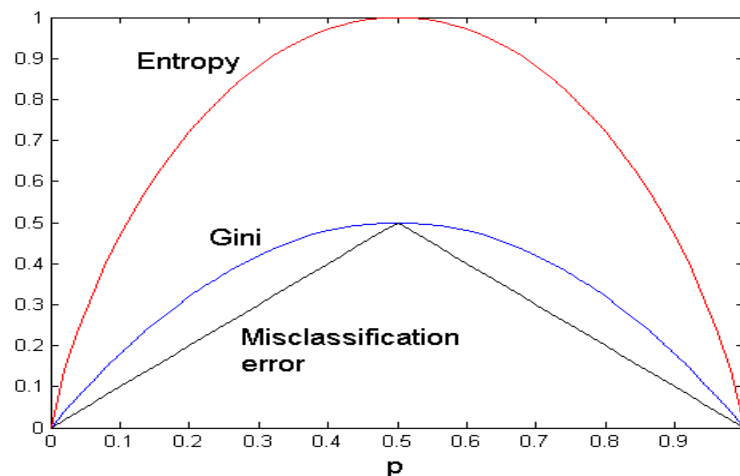
47

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Comparison among Splitting Criteria



For a 2-class problem:

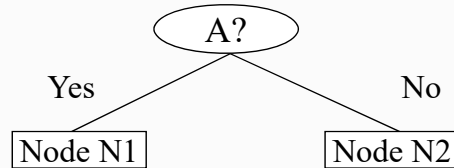


48

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.



# Misclassification Error vs Gini



	Parent
C1	7
C2	3
<b>Gini = 0.42</b>	

$$\begin{aligned} \text{Gini}(N1) &= 1 - (3/3)^2 - (0/3)^2 \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{Gini}(N2) &= 1 - (4/7)^2 - (3/7)^2 \\ &= 0.489 \end{aligned}$$

	N1	N2
C1	3	4
C2	0	3
<b>Gini=0.342</b>		

$$\begin{aligned} \text{Gini(Children)} &= 3/10 * 0 \\ &+ 7/10 * 0.489 \\ &= 0.342 \end{aligned}$$

**Gini improves !!**

49

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

# Tree Induction



## Greedy strategy.

- Split the records based on an attribute test that optimizes certain criterion.

## Issues

- Determine how to split the records
  - How to specify the attribute test condition?
  - How to determine the best split?
- Determine when to stop splitting

50

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Stopping Criteria for Tree Induction

- Stop expanding a node when all the records belong to the same class
- Stop expanding a node when all the records have similar attribute values
  - What to do? majority voting
- Early termination, e.g., when the information gain is below a threshold.

51

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Decision Tree Based Classification

### Advantages:

- Inexpensive to construct
- Extremely fast at classifying unknown records
- Easy to interpret for small-sized trees
- Accuracy is comparable to other classification techniques for many simple data sets

52

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Example: C4.5



Simple depth-first construction.

Uses Information Gain

Sorts Continuous Attributes at each node.

Needs entire data to fit in memory.

Unsuitable for Large Datasets.

- o Needs out-of-core sorting.

You can download the software from:

<http://www.cse.unsw.edu.au/~quinlan/c4.5r8.tar.gz>