

## Gini Index (CART, IBM IntelligentMiner)

- If a data set  $D$  contains examples from  $n$  classes, gini index,  $gini(D)$  is defined as

$$gini(D) = 1 - \sum_{j=1}^n p_j^2$$

where  $p_j$  is the relative frequency of class  $j$  in  $D$

- If a data set  $D$  is split on  $A$  into two subsets  $D_1$  and  $D_2$ , the gini index  $gini_A(D)$  is defined as

$$gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$$

- Reduction in Impurity:

$$\Delta gini(A) = gini(D) - gini_A(D)$$

- The attribute provides the smallest  $gini_{split}(D)$  (or the largest reduction in impurity) is chosen to split the node
- The Gini index considers a binary split for each attribute

18

### Binary Splitting using GINI Index

Age has 3 values, so

- Possible splits are:

- $\{<=30, 31-40\}, \{>40\}$
- $\{<=30, >40\}, \{31-40\}$
- $\{31-40, >40\}, \{<=30\}$

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

$$Gini_{\{<=30, 31-40\}, \{>40\}} = \frac{9}{14} \{1 - (\frac{6}{9})^2 - (\frac{3}{9})^2\} + \frac{5}{14} \{1 - (\frac{3}{5})^2 - (\frac{2}{5})^2\}$$

$$= 0.48$$

$$Gini_{\{<=30, >40\}, \{31-40\}} = \frac{10}{14} \{1 - (\frac{5}{10})^2 - (\frac{5}{10})^2\} + \frac{4}{14} \{1 - (\frac{4}{4})^2 - 0\}$$

$$= 0.357$$

$$Gini_{\{31-40, >40\}, \{<=30\}} = \frac{9}{14} \{1 - (\frac{7}{9})^2 - (\frac{2}{9})^2\} + \frac{5}{14} \{1 - (\frac{2}{5})^2 - (\frac{3}{5})^2\}$$

$$= 0.39$$

19

- Thus, amongst the three possible splits of age, the gini of  $\{<=30, >40\}, \{31-40\}$  is the lowest.
- So, now there are two possible outcomes of age.
- Let middle\_aged = 31-40 ✓

youth =  $<=30$

senior =  $>40$

- Now the two possible outcomes of age are

1.  $\{y,s\}$
2.  $\{m_a\}$

age	income	student	credit_rating	buys_computer
{y,s}	high	no	fair	no
{y,s}	high	no	excellent	no
{m}	high	no	fair	yes
{y,s}	medium	no	fair	yes
{y,s}	low	yes	fair	yes
{y,s}	low	yes	excellent	no
{m}	low	yes	excellent	yes
{y,s}	medium	no	fair	no
{y,s}	low	yes	fair	yes
{y,s}	medium	yes	fair	yes
{y,s}	medium	yes	excellent	yes
{m}	medium	no	excellent	yes
{m}	high	yes	fair	yes
{y,s}	medium	no	excellent	no

20

- Binary splitting for income as income has three possible values:

1. Low
2. Medium
3. High

age	income	student	credit_rating	buys_computer
{y,s}	high	no	fair	no
{y,s}	high	no	excellent	no
{m}	high	no	fair	yes
{y,s}	medium	no	fair	yes
{y,s}	low	yes	fair	yes
{y,s}	low	yes	excellent	no
{m}	low	yes	excellent	yes
{y,s}	medium	no	fair	no
{y,s}	low	yes	fair	yes
{y,s}	medium	yes	fair	yes
{y,s}	medium	yes	excellent	yes
{m}	medium	no	excellent	yes
{m}	high	yes	fair	yes
{y,s}	medium	no	excellent	no

- $Gini_{\{low,high\},\{medium\}} = 0.458$

- $Gini_{\{low,medium\},\{high\}} = 0.443$  ✓

- $Gini_{\{medium,high\},\{low\}} = 0.45$

21

## Computation of Gini Index

- Ex. D has 9 tuples in buys\_computer = "yes" and 5 in "no"

$$gini(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.459$$

- Gini<sub>student</sub> =  $\frac{7}{14}\{1 - (\frac{6}{7})^2 - (\frac{1}{7})^2\} + \frac{7}{14}\{1 - (\frac{3}{7})^2 - (\frac{4}{7})^2\} = 0.367$
- Gini<sub>CR</sub> =  $\frac{8}{14}\{1 - (\frac{6}{8})^2 - (\frac{2}{8})^2\} + \frac{6}{14}\{1 - (\frac{3}{6})^2 - (\frac{3}{6})^2\} = 0.429$
- Gini<sub>Income</sub> = 0.443
- Gini<sub>age</sub> = 0.358

age	income	student	credit_rating	buys_computer
{y,s}	high	no	fair	no
{y,s}	high	no	excellent	no
{m}	high	no	fair	yes
{y,s}	medium	no	fair	yes
{y,s}	low	yes	fair	yes
{y,s}	low	yes	excellent	no
{m}	low	yes	excellent	yes
{y,s}	medium	no	fair	no
{y,s}	low	yes	fair	yes
{y,s}	medium	yes	fair	yes
{y,s}	medium	yes	excellent	yes
{m}	medium	no	excellent	yes
{m}	high	yes	fair	yes
{y,s}	medium	no	excellent	no

## Computing Information-Gain for Continuous-Valued Attributes

- Let attribute A be a continuous-valued attribute
- Must determine the *best split point* for A
  - Sort the value A in increasing order
  - Typically, the midpoint between each pair of adjacent values is considered as a possible *split point*
    - $(a_i + a_{i+1})/2$  is the midpoint between the values of  $a_i$  and  $a_{i+1}$
  - The point with the *minimum expected information requirement* for A is selected as the split-point for A
- Split:
  - D1 is the set of tuples in D satisfying  $A \leq \text{split-point}$ , and D2 is the set of tuples in D satisfying  $A > \text{split-point}$

A	Class
100	No
75	No
90	Yes
95	Yes
60	No
120	No
220	No
70	No
85	Yes
125	No

7 --> no  
3 --> yes

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

A	Class
>97	No
<=97	No
<=97	Yes
<=97	Yes
<=97	No
>97	No
>97	No
<=97	No
<=97	Yes
>97	No

24

## Gain Ratio for Attribute Selection (C4.5)

- Information measure is biased towards attributes with a large number of values
- C4.5 (a successor of ID3) uses gain ratio to overcome the problem (normalization to information gain)

$$SplitInfo_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right)$$

$$GainRatio(A) = Gain(A) / SplitInfo(A)$$

- Ex.  $SplitInfo_{income}(D) = - \frac{4}{14} \times \log_2 \left( \frac{4}{14} \right) - \frac{6}{14} \times \log_2 \left( \frac{6}{14} \right) - \frac{4}{14} \times \log_2 \left( \frac{4}{14} \right) = 1.557$ 
  - gain\_ratio(income) = 0.029/1.557 = 0.019
- The attribute with the maximum gain ratio is selected as the splitting attribute

25



## Comparing Attribute Selection Measures

- The three measures, in general, return good results but
  - **Information gain:**
    - biased towards multivalued attributes
  - **Gain ratio:**
    - tends to prefer unbalanced splits in which one partition is much smaller than the others
  - **Gini index:**
    - biased to multivalued attributes
    - has difficulty when # of classes is large
    - tends to favor tests that result in equal-sized partitions and purity in both partitions

26

## Practice with GINI index

Attribute 1	Attribute 2	Attribute 3	Class
A	70	T	C1
A	90	T	C2
A	85	F	C2
A	95	F	C2
A	70	F	C1
B	90	T	C1
B	78	F	C1
B	65	T	C1
B	75	F	C1
C	80	T	C2
C	70	T	C2
C	80	F	C1
C	80	F	C1
C	96	F	C1

27

## Practice with Gain Ratio

Height	Hair	Eyes	Class
Short	Blond	Blue	+
Tall	Blond	Brown	-
Tall	Red	Blue	+
Short	Dark	Blue	-
Tall	Dark	Blue	-
Tall	Blond	Blue	+
Tall	Dark	Brown	-
Short	Blond	Brown	-