

Module : Decision Tree

Instructor: Dr. Darshan
Ingle



HI Y'ALL !!!

STARTING SOON



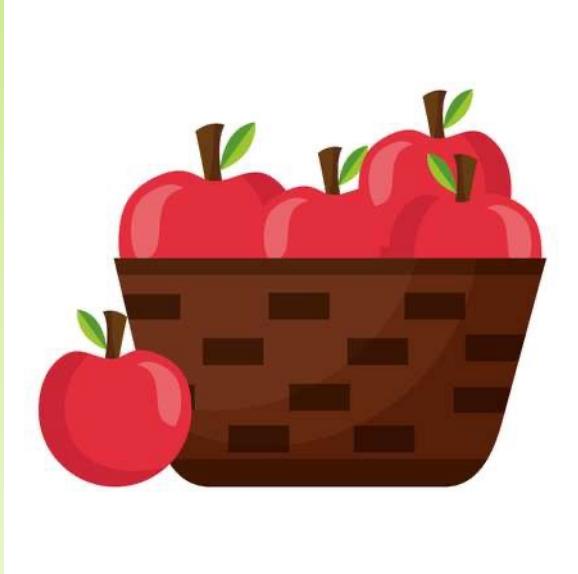
Entropy & Information Gain

- Entropy is the measure of **randomness** or **impurity in** the **data** set.
- Entropy **uses** the **concept of homogeneity**.

Things to Remember:

- **If** samples are completely **homogeneous, then** the **entropy** of that attribute **will be zero**.
- **If samples are equally divided, then entropy will be one.**
- So out of the heterogeneous options **we** need to **select** the **ones having maximum homogeneity**.

Purity vs Impurity in Data



Impurity = 0

Homogeneous Data

Impurity ≠ 0



Heterogeneous Data

$$n=14 < \begin{matrix} \text{Yes}=9 \\ \text{No}=5 \end{matrix}$$

$$P(\text{Yes}) = 9/14$$

$$P(\text{No}) = 5/14$$

Steps to estimate Entropy & Information Gain

- Calculate the expected information needed to classify a tuple in data set (D) is given by –

$m=2$ (\because there are 2 unique classes)

$$\text{Entropy}(D) = - \sum_{i=1}^m p_i \log_2(p_i) = - \sum_{i=1}^2 p_i \cdot \log_2 p_i$$

- We will check how many tuples are **yes** and **no** in target variable in the below data set.

	x_1	x_2	x_3	x_4	$y = \{ \text{no}, \text{yes} \}$
A	age	income	student	credit_rating	buys_computer
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle-aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle-aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle-aged	medium	no	excellent	yes
13	middle-aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

$$= - [p_1 \cdot \log_2 p_1 + p_2 \cdot \log_2 p_2]$$

$$= - [P_{\text{Yes}} \cdot \log_2 P_{\text{Yes}} + P_{\text{No}} \cdot \log_2 P_{\text{No}}]$$

$$\text{Entropy}(D) = - p(\text{yes}) \times \log_2(p(\text{yes})) - p(\text{no}) \times \log_2(p(\text{no}))$$

$$\text{Entropy}(D) = - 9/14 \times \log_2(9/14) - 5/14 \times \log_2(5/14)$$

$$\text{Entropy}(D) = 0.94$$

$$n=14 \leftarrow \text{Yes}=9 \\ \text{No}=4$$

Steps to estimate Entropy & Information Gain Cont.

Age:	Yes		No	Total
	4	2	3	5
M	4	0	4	
S	3	2	5	
			$\sum = 14$	

- How much more information would we still need (after the partitioning) to arrive at an exact classification?

This amount is measured by –

$$\text{Entropy}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times \text{Entropy}(D_j)$$

$$\text{Age} = \{4, M, S\}$$

- We will check how many tuples are **yes** and **no** in target variable along with that particular **predictor variable**.

	A	B	C	D	E
1	age	income	student	credit_rating	buys_computer
2	youth	high	no	fair	no
3	youth	high	no	excellent	no
4	middle-aged	high	no	fair	yes
5	senior	medium	no	fair	yes
6	senior	low	yes	fair	yes
7	senior	low	yes	excellent	no
8	middle-aged	low	yes	excellent	yes
9	youth	medium	no	fair	no
10	youth	low	yes	fair	yes
11	senior	medium	yes	fair	yes
12	youth	medium	yes	excellent	yes
13	middle-aged	medium	no	excellent	yes
14	middle-aged	high	yes	fair	yes
15	senior	medium	no	excellent	no

$$E = - \sum_{i=1}^v P_i \cdot \log_2 P_i$$

$$E_4 = - [P_{4_{\text{Yes}}} \cdot \log_2 P_{4_{\text{Yes}}} + P_{4_{\text{No}}} \cdot \log_2 P_{4_{\text{No}}}]$$

$$\text{Entropy}_{\text{Age}}(D) = \frac{|D_{\text{youth}}|}{|D|} \times \text{Entropy}(D_{\text{youth}}) + \frac{|D_{\text{middle-aged}}|}{|D|} \times \text{Entropy}(D_{\text{middle-aged}}) + \frac{|D_{\text{senior}}|}{|D|} \times \text{Entropy}(D_{\text{senior}})$$

$$\text{Entropy}_{\text{Age}}(D) = 5/14 \times [-2/5 \log_2(2/5) - 3/5 \log_2(3/5)] + 4/14 \times [-4/4 \log_2(4/4)] + 5/14 \times [-3/5 \log_2(3/5) - 2/5 \log_2(2/5)]$$

$$\text{Entropy}_{\text{Age}}(D) = 0.629$$

$$E_{\text{Senior}} = - \sum_{i=1}^2 P_i \cdot \log_2 P_i = - [P_{S_{\text{Yes}}} \cdot \log_2 P_{S_{\text{Yes}}} + P_{S_{\text{No}}} \cdot \log_2 P_{S_{\text{No}}}]$$

$$= - \left[\frac{3}{5} \cdot \log_2 \frac{3}{5} + \frac{2}{5} \cdot \log_2 \frac{2}{5} \right]$$

Steps to estimate Entropy & Information Gain Cont.

- Final step is to calculate Information Gain –



$$\text{Information Gain}(A) = \text{Entropy}(D) - \text{Entropy}_A(D)$$

$$\text{Information Gain}(\text{age}) = 0.940 - 0.629 = 0.248$$

Example on Buys_Computer Data set

	x_1	x_2	x_3	x_4	y
	A	B	C	D	E
1	age	income	student	credit_rating	buys_computer
2	youth	high	no	fair	no
3	youth	high	no	excellent	no
4	middle-aged	high	no	fair	yes
5	senior	medium	no	fair	yes
6	senior	low	yes	fair	yes
7	senior	low	yes	excellent	no
8	middle-aged	low	yes	excellent	yes
9	youth	medium	no	fair	no
10	youth	low	yes	fair	yes
11	senior	medium	yes	fair	yes
12	youth	medium	yes	excellent	yes
13	middle-aged	medium	no	excellent	yes
14	middle-aged	high	yes	fair	yes
15	senior	medium	no	excellent	no

It is a classification problem



- Predictors are: age, income, student, credit_rating. Target variable is buys_computer.

Calculation Work

E_y

①

$$\text{Entropy}(D) = -\frac{9}{14} \times \log_2(9/14) - \frac{5}{14} \times \log_2(5/14) = 0.94$$

E_{x1}

②

$$\text{Entropy}_{\text{age}}(D) = \frac{|D_{\text{youth}}|}{|D|} \times \text{Entropy}(D_{\text{youth}}) + \frac{|D_{\text{middle}}|}{|D|} \times \text{Entropy}(D_{\text{middle}}) + \frac{|D_{\text{senior}}|}{|D|} \times \text{Entropy}(D_{\text{senior}})$$

$$\text{Entropy}_{\text{age}}(D) = \frac{5}{14} [-\frac{2}{5} \log_2(2/5) - \frac{3}{5} \log_2(3/5)] + \frac{4}{14} [-\frac{4}{4} \log_2(4/4)] + \frac{5}{14} [-\frac{3}{5} \log_2(3/5) - \frac{2}{5} \log_2(2/5)]$$

$$\text{Entropy}_{\text{age}}(D) = 0.629$$

E_{x2}

③

$$\text{Entropy}_{\text{income}}(D) = \frac{|D_{\text{high}}|}{|D|} \times \text{Entropy}(D_{\text{high}}) + \frac{|D_{\text{medium}}|}{|D|} \times \text{Entropy}(D_{\text{medium}}) + \frac{|D_{\text{low}}|}{|D|} \times \text{Entropy}(D_{\text{low}})$$

$$\text{Entropy}_{\text{income}}(D) = \frac{4}{14} [-\frac{2}{4} \log_2(2/4) - \frac{2}{4} \log_2(2/4)] + \frac{6}{14} [-\frac{4}{6} \log_2(4/6) - \frac{2}{6} \log_2(2/6)] + \frac{4}{14} [-\frac{3}{4} \log_2(3/4) - \frac{1}{4} \log_2(1/4)]$$

$$\text{Entropy}_{\text{income}}(D) = 0.908$$

E_{x3}

④

$$\text{Entropy}_{\text{student}}(D) = \frac{|D_{\text{yes}}|}{|D|} \times \text{Entropy}(D_{\text{yes}}) + \frac{|D_{\text{no}}|}{|D|} \times \text{Entropy}(D_{\text{no}})$$

$$\text{Entropy}_{\text{student}}(D) = 0.786$$

E_{x4}

⑤

$$\text{Entropy}_{\text{credit_rating}}(D) = \frac{|D_{\text{fair}}|}{|D|} \times \text{Entropy}(D_{\text{fair}}) + \frac{|D_{\text{excellent}}|}{|D|} \times \text{Entropy}(D_{\text{excellent}})$$

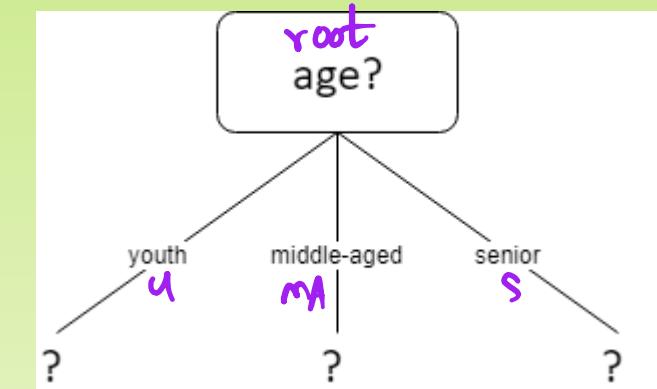
$$\text{Entropy}_{\text{credit_rating}}(D) = 0.89$$



Tabular View

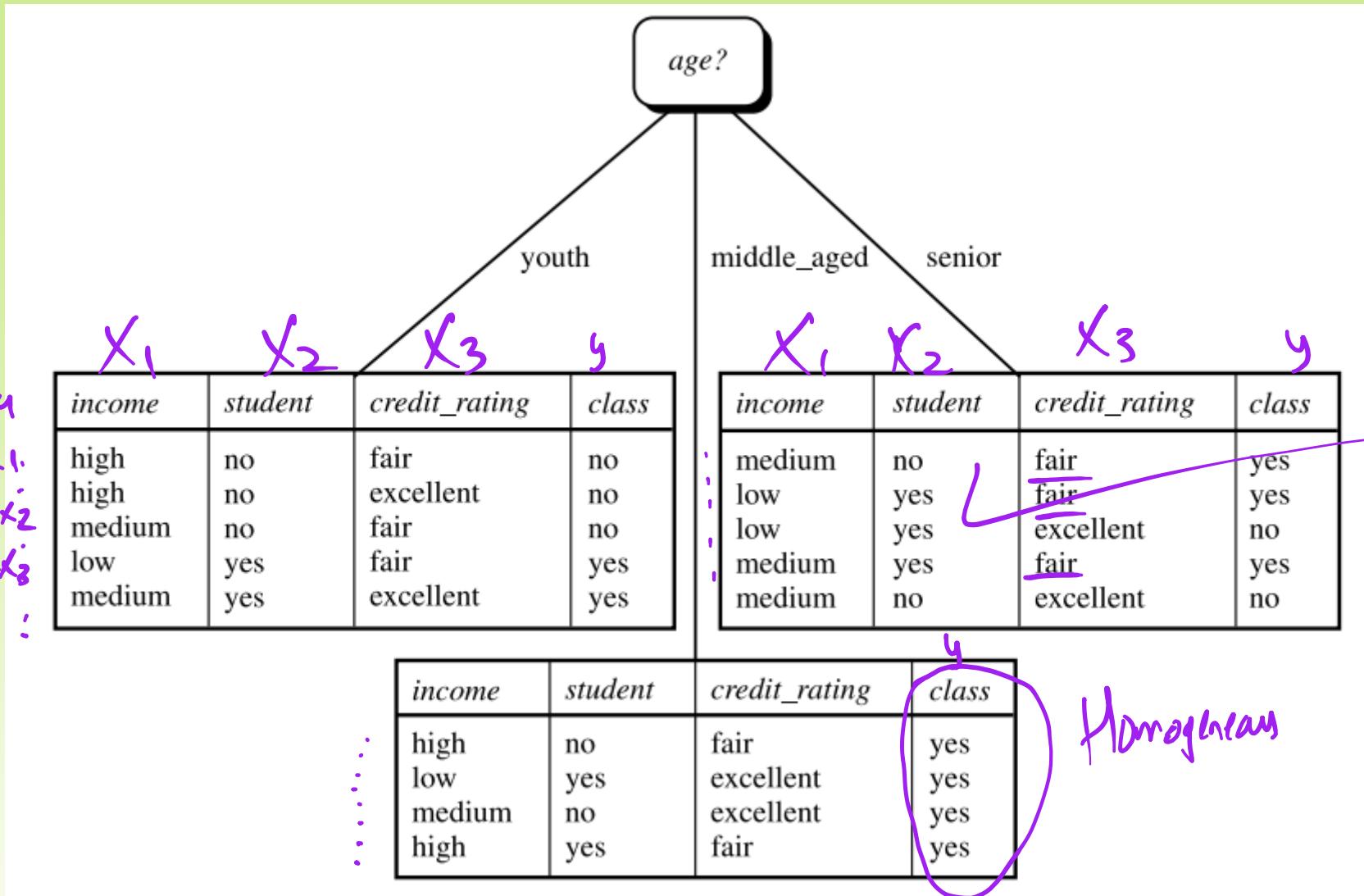
- The **age** attribute is **giving maximum information** gain. So the **root** node **will be age**.

	A	B	C
1	Data		
2	age	E_y	0
3	income	E_{x_1}	= 0.248 (Max)
4	student	E_{x_2}	= 0.032
5		E_{x_3}	= 0.154
6	credit_rating	E_{x_4}	= 0.05



- But how should I choose next attribute?
 - Repeat the step we've done so far on the subset of data.

Next View of Representation



Tabular view on subset data

Solution for **youth data**

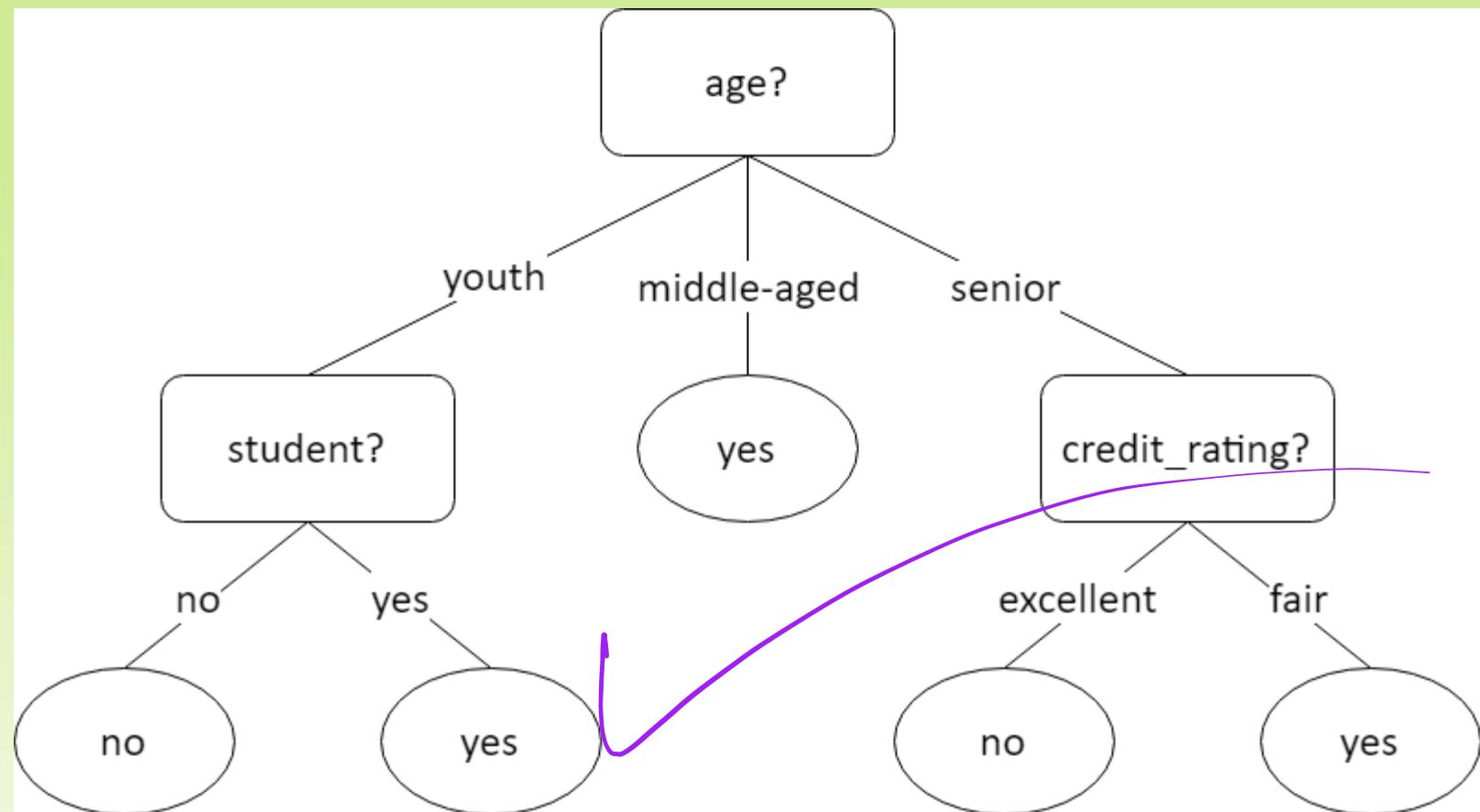
	Data	Entropy	Information Gain
12			
13	Data	0.97	0
14	income	0.399	0.571
15	student	0	0.97
16	credit_rating	0.948	0.022

Solution for **senior data**

	Data	Entropy	Information Gain
21			
22	Data	0.968	0
23	income	0.95	0.018
24	student	0.95	0.018
25	credit_rating	0	0.968

- Note: Entropy for middle-aged data is zero.

Decision Tree Complete View



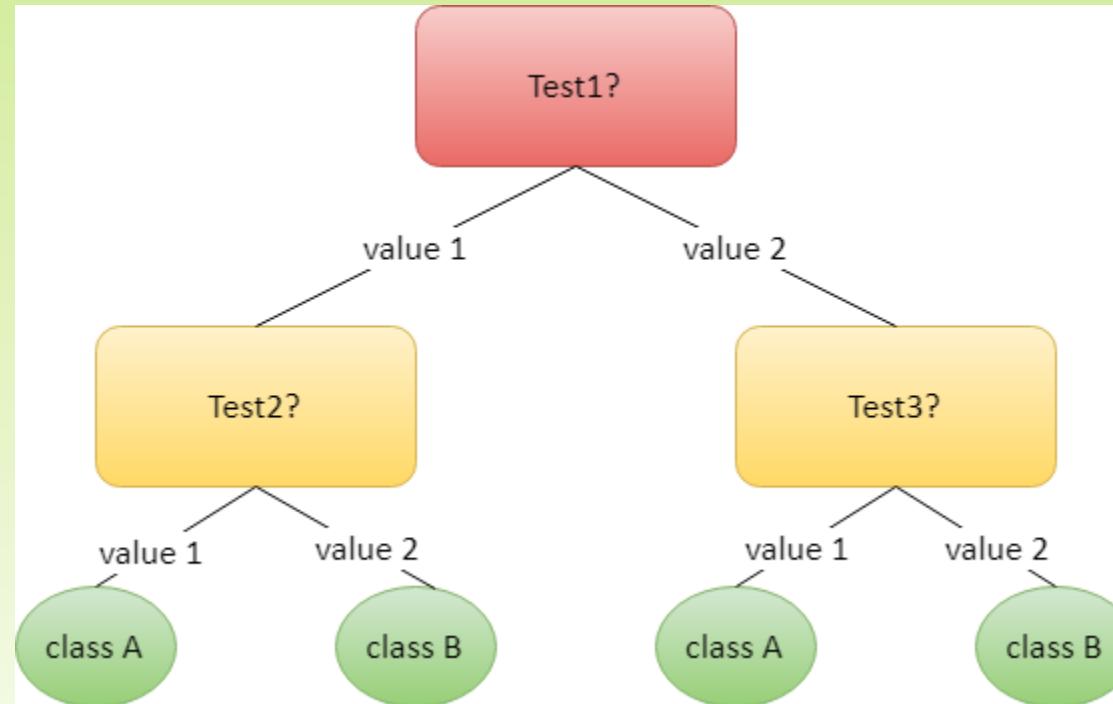
Agenda

- Terminology Related to Trees
- Decision Tree
- Decision Tree Algorithms
- Attribute Selection Measures
- ID3 Algorithm
- Entropy & Information Gain
- Steps to Estimate Entropy & Information Gain

- CART Algorithm**
- Gini Index
- Steps to estimate Gini Index
- CART – Regression Example
- Issues with Decision Trees
- Tree Pruning
- Decision Tree Applications

CART Algorithm

- CART is known as **Classification And Regression Trees**. It has binary representation of data.
- It uses **Gini index** to find best split on the attributes.



Agenda

- Terminology Related to Trees
- Decision Tree
- Decision Tree Algorithms
- Attribute Selection Measures
- ID3 Algorithm
- Entropy & Information Gain
- Steps to Estimate Entropy & Information Gain

- CART Algorithm
- Gini Index**
- Steps to estimate Gini Index
- CART – Regression Example
- Issues with Decision Trees
- Tree Pruning
- Decision Tree Applications

Gini Index

- Gini Index is an attribution selection measure to **find best split** on attributes.
- The concept is same as Information Gain except it helps to **find best split** for binary representation.
- It helps to **find maximum reduction in impurity** of data.



Remember: If samples are complete **homogeneous**, then their **Gini Index** will be **zero**.

Steps to estimate Gini Index

$$n=14 \quad \begin{cases} \text{Yes} = 9 \\ \text{No} = 5 \end{cases}$$

- Gini index measures the impurity of D, a data partition or set of training tuples, as

$$\text{Gini}(D) = 1 - \sum_{i=1}^m p_i^2 = 1 - \sum_{i=1}^2 p_i^2 = 1 - \left[p_{\text{Yes}}^2 + p_{\text{No}}^2 \right]$$

$$= 1 - \left[\left(\frac{9}{14} \right)^2 + \left(\frac{5}{14} \right)^2 \right]$$

- We will check how many tuples are **yes** and **no** in target variable.

	x_1	x_2	x_3	x_4	$y \in \{\text{Yes}, \text{No}\}$
1	A	B	C	D	E
1	age	income	student	credit_rating	buys_computer
2	youth	high	no	fair	no
3	youth	high	no	excellent	no
4	middle-aged	high	no	fair	yes
5	senior	medium	no	fair	yes
6	senior	low	yes	fair	yes
7	senior	low	yes	excellent	no
8	middle-aged	low	yes	excellent	yes
9	youth	medium	no	fair	no
10	youth	low	yes	fair	yes
11	senior	medium	yes	fair	yes
12	youth	medium	yes	excellent	yes
13	middle-aged	medium	no	excellent	yes
14	middle-aged	high	yes	fair	yes
15	senior	medium	no	excellent	no

$$\text{Gini}(D) = 1 - p(\text{yes})^2 - p(\text{no})^2$$

$$\text{Gini}(D) = 1 - (9/14)^2 - (5/14)^2$$

$$\text{Gini}(D) = 0.459$$

Steps to estimate Gini Index Cont.

Age: $\{U, MA, S\} \rightarrow$ One vs Rest

U vs $\{MA, S\}$ → (a)
 MA vs $\{U, S\}$ → (b)
 S vs $\{U, MA\}$ → (c) *

3 subsets

- When considering a binary split, we compute a weighted sum of the impurity of each resulting partition, the Gini index of D given that partitioning is

$$Gini = 1 - \sum P_i^2 = 1 - \left[P_{U_{yes}}^2 + P_{S_{no}}^2 \right] = 1 - \left[\left(\frac{3}{5}\right)^2 + \left(\frac{2}{5}\right)^2 \right]$$

$$Gini_A(D) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2)$$

- We will check how many tuples are **yes** and **no** in target variable along with that particular **predictor variable**.

A	B	C	D	E
1 age	income	student	credit_rating	buys_computer
2 youth	high	no	fair	no
3 youth	high	no	excellent	no
4 middle-aged	high	no	fair	yes
5 senior	medium	no	fair	yes
6 senior	low	yes	fair	yes
7 senior	low	yes	excellent	no
8 middle-aged	low	yes	excellent	yes
9 youth	medium	no	fair	no
10 youth	low	yes	fair	yes
11 senior	medium	yes	fair	yes
12 youth	medium	yes	excellent	yes
13 middle-aged	medium	no	excellent	yes
14 middle-aged	high	yes	fair	yes
15 senior	medium	no	excellent	no

(This value is computed only for subset of split)

$$Gini_{age}(D) = \frac{|D_{youth, middle-aged}|}{|D|} \times Gini(D_{youth, middle-aged}) + \frac{|D_{senior}|}{|D|} \times Gini(D_{senior})$$

$$Gini_{age}(D) = \frac{9/14}{14} \times [1 - \left(\frac{6}{9}\right)^2 - \left(\frac{3}{9}\right)^2] + \frac{5/14}{14} \times [1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2]$$

$$Gini_{age}(D) = 0.455$$

	Yes	No	Total
U	2	3	5
MA	4	6	10
S.	5	2	7

$$G = 1 - \sum P_i^2 = 1 - \left[P_{U,yes}^2 + P_{U,no}^2 \right]$$

$$n=14$$

$$= 1 - \left[\left(\frac{6}{9}\right)^2 + \left(\frac{3}{9}\right)^2 \right]$$

Steps to estimate Gini Index Cont.

- The reduction in impurity that would be incurred by a binary split on a discrete or continuous-valued attribute A is



$$\Delta \text{Gini}(A) = \text{Gini}(D) - \text{Gini}_A(D)$$

$$\Delta \text{Gini}(\text{age}) = 0.459 - 0.455 = 0.004$$

(This value is computed only for subset of split)

Example on Buys_Computer Data set

	A	B	C	D	E
1	age	income	student	credit_rating	buys_computer
2	youth	high	no	fair	no
3	youth	high	no	excellent	no
4	middle-aged	high	no	fair	yes
5	senior	medium	no	fair	yes
6	senior	low	yes	fair	yes
7	senior	low	yes	excellent	no
8	middle-aged	low	yes	excellent	yes
9	youth	medium	no	fair	no
10	youth	low	yes	fair	yes
11	senior	medium	yes	fair	yes
12	youth	medium	yes	excellent	yes
13	middle-aged	medium	no	excellent	yes
14	middle-aged	high	yes	fair	yes
15	senior	medium	no	excellent	no



- **Predictors** are: age, income, student, credit_rating. **Target variable** is buys_computer.

Calculation Work

G_{D_y} ①
~~X₁~~
Age ②

$$Gini(D) = 1 - (9/14)^2 - (5/14)^2 = 0.459$$

A $Gini_{age} \in \{\text{youth, middle-aged or senior}, \text{youth, senior or middle-aged}, \text{middle-aged, senior or youth}\}$

c b a

$Gini_{age} \in \{\text{youth, middle-aged or senior}\}$

$$Gini_{age} = 9/14 \times Gini(D1) + 5/14 \times Gini(D2)$$

$$Gini_{age} = 9/14 \times [1 - (6/9)^2 - (3/9)^2] + 5/14 \times [1 - (3/5)^2 - (2/5)^2]$$

$$Gini_{age} = 0.455$$

$$\Delta Gini(\text{age}) = 0.459 - 0.455 = 0.004$$

b {
c }
 $Gini_{age} \in \{\text{youth, senior or middle-aged}\}$

$$Gini_{age} = 10/14 \times Gini(D1) + 4/14 \times Gini(D2)$$

$$Gini_{age} = 10/14 \times [1 - (5/10)^2 - (5/10)^2] + 4/14 \times [1 - (4/4)^2]$$

$$Gini_{age} = 0.357$$

(Optimal for binary split)

$$\Delta Gini(\text{age}) = 0.459 - 0.357 = 0.102$$

highest reduction

a {
b }
 $Gini_{age} \in \{\text{middle-aged, senior or youth}\}$

$$Gini_{age} = 9/14 \times Gini(D1) + 5/14 \times Gini(D2)$$

$$Gini_{age} = 9/14 \times [1 - (7/9)^2 - (2/9)^2] + 5/14 \times [1 - (2/5)^2 - (3/5)^2]$$

$$Gini_{age} = 0.393$$

$$\Delta Gini(\text{age}) = 0.459 - 0.393 = 0.066$$

Calculation Work Cont.

x_2

B $\text{Gini}_{\text{income}} \in \{\{\text{low}, \text{medium}\} \text{ or } \{\text{high}\}, \{\text{low}, \text{high}\} \text{ or } \{\text{medium}\}, \{\text{medium}, \text{high}\} \text{ or } \{\text{low}\}\}$

C b a

$\text{Gini}_{\text{income}} \in \{\text{low}, \text{medium}\} \text{ or } \{\text{high}\}$

$$\text{Gini}_{\text{income}} = 10/14 \times \text{Gini}(D1) + 4/14 \times \text{Gini}(D2)$$

$$\text{Gini}_{\text{income}} = 10/14 \times [1 - (7/10)^2 - (3/10)^2] + 4/14 \times [1 - (2/4)^2 - (2/4)^2] \quad (\text{Optimal for binary split})$$

$$\text{Gini}_{\text{income}} = 0.443$$

$$\Delta \text{Gini}(\text{income}) = 0.459 - 0.443 = 0.016$$

c

b

a

$\text{Gini}_{\text{income}} \in \{\text{low}, \text{high}\} \text{ or } \{\text{medium}\}$

$$\text{Gini}_{\text{income}} = 8/14 \times \text{Gini}(D1) + 6/14 \times \text{Gini}(D2)$$

$$\text{Gini}_{\text{income}} = 8/14 \times [1 - (5/8)^2 - (3/8)^2] + 6/14 \times [1 - (4/6)^2 - (2/6)^2]$$

$$\text{Gini}_{\text{income}} = 0.458$$

$$\Delta \text{Gini}(\text{income}) = 0.459 - 0.458 = 0.001$$

$\text{Gini}_{\text{income}} \in \{\text{medium}, \text{high}\} \text{ or } \{\text{low}\}$

$$\text{Gini}_{\text{income}} = 8/14 \times \text{Gini}(D1) + 6/14 \times \text{Gini}(D2)$$

$$\text{Gini}_{\text{income}} = 10/14 \times [1 - (5/10)^2 - (5/10)^2] + 4/14 \times [1 - (3/4)^2 - (1/4)^2]$$

$$\text{Gini}_{\text{income}} = 0.450$$

$$\Delta \text{Gini}(\text{income}) = 0.459 - 0.450 = 0.009$$

Calculation Work Cont.

x_3

C $Gini_{student} \in \{\text{yes}\} \text{ or } \{\text{no}\}$

$$Gini_{student} = \underline{7/14 \times Gini(D1) + 7/14 \times Gini(D2)}$$

$$Gini_{student} = 7/14 \times [1 - (6/7)^2 - (1/7)^2] + 7/14 \times [1 - (3/7)^2 - (4/7)^2]$$

$$Gini_{student} = 0.368$$

$$\Delta Gini(\text{student}) = 0.459 - 0.368 = 0.091$$

x_u

D $Gini_{credit_rating} \in \{\text{excellent}\} \text{ or } \{\text{fair}\}$

$$Gini_{credit_rating} = \underline{8/14 \times Gini(D1) + 6/14 \times Gini(D2)}$$

$$Gini_{credit_rating} = 8/14 \times [1 - (6/8)^2 - (2/8)^2] + 6/14 \times [1 - (3/6)^2 - (3/6)^2]$$

$$Gini_{credit_rating} = 0.428$$

$$\Delta Gini(\text{credit_rating}) = 0.459 - 0.428 = 0.031$$

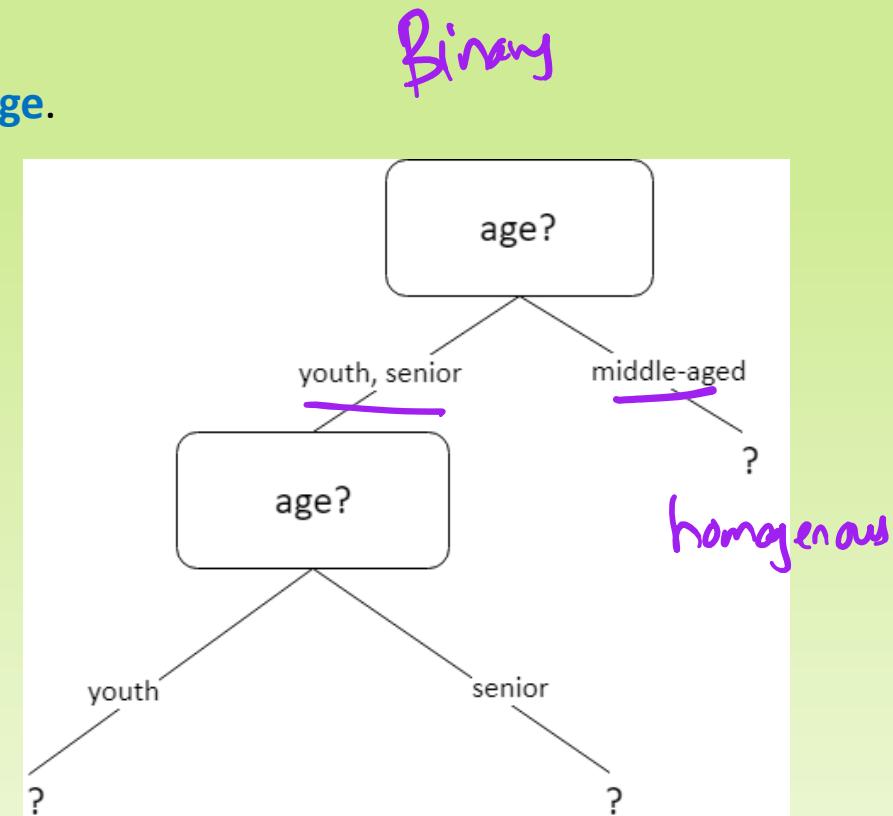
Tabular Representation

- The **age** attribute is **giving Reduction in Impurity**. So the **root node will be age**.

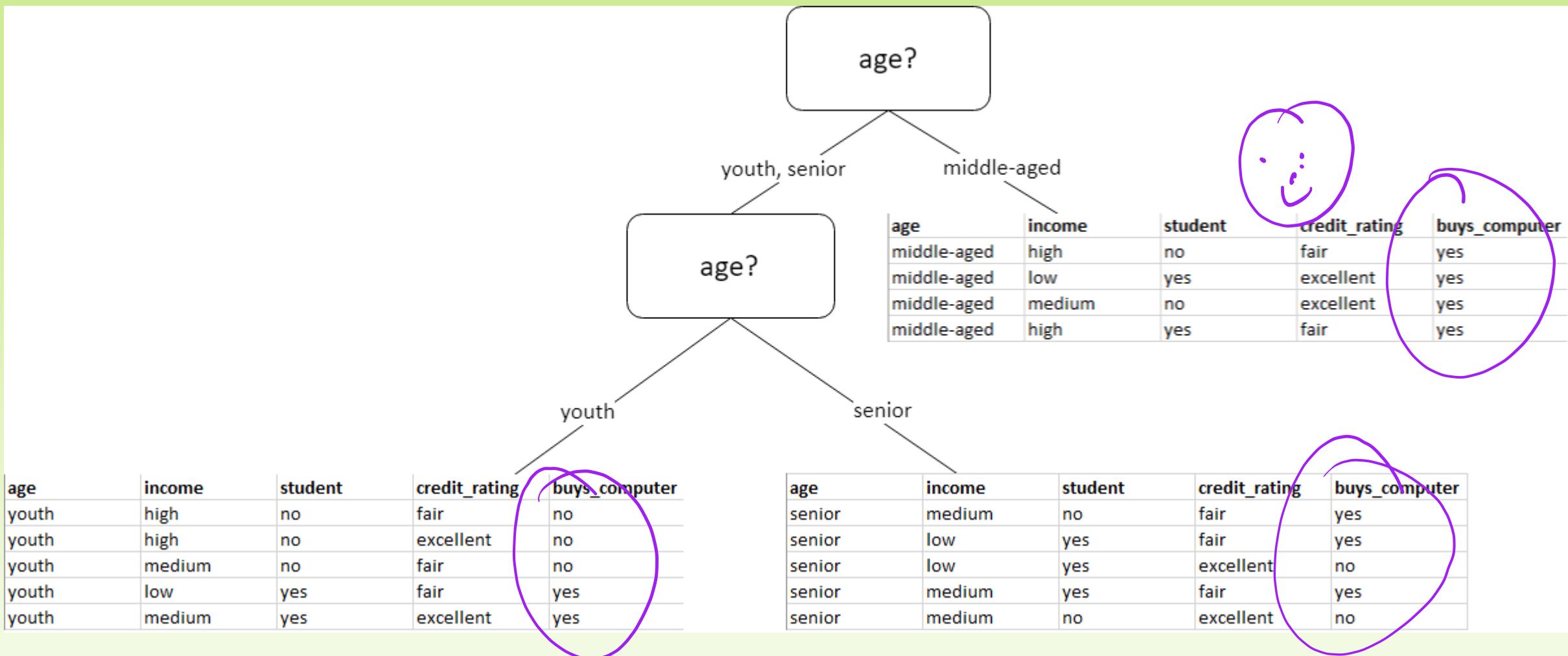
	A	B	C	D
1				
2	Data			
3	age	✓ 0.459 ✓ 0.455 ✓ 0.357 ✓ 0.393	0 0.004 0.102 0.066	(Highest Reduction)
4	income			
5		0.443	0.016	
6		0.458	0.001	
7		0.45	0.009	
8	student	0.091	0.091	
9				
10	credit_rating	0.428	0.031	



- But how should I choose next attribute?
 - Repeat the step we've done so far on the subset of data.



Next View of Representation



Tabular view on subset data

Solution for **youth** data

Data	Gini Index	Reduction in Impurity
income	0.266	0.214
student	0.466	0.014
credit_rating	0.3	0.18
student	0	0.48
credit_rating	0.466	0.014

Solution for **senior** data

Data	Gini Index	Reduction in Impurity
income	0.48	0
student	0.466	0.014
student	0.466	0.014
credit_rating	0	0.48

- **Note:** Gini Index for middle-aged data is zero.

CART Complete View

