

# Decision Tree



# Course Overview

You are here...

Term	CDF	GCD	GCDAI	PGPDSAI
Term 1	Data Analytics with Python	Data Analytics with Python	Data Analytics with Python	Data Analytics with Python
Term 2	Data Visualization Techniques	Data Visualization Techniques	Data Visualization Techniques	Data Visualization Techniques
Term 3	EDA & Data Storytelling	EDA & Data Storytelling	EDA & Data Storytelling	EDA & Data Storytelling
		Minor Project	Minor Project	Minor Project
Term 4		Machine Learning Foundation	Machine Learning Foundation	Machine Learning Foundation
Term 5		Machine Learning Intermediate	Machine Learning Intermediate	Machine Learning Intermediate
Term 6		Machine Learning Advanced (Mandatory)	Machine Learning Advanced (Mandatory)	Machine Learning Advanced (Mandatory)
		Data Visualization with Tableau (Elective - I)	Data Visualization with Tableau (Elective - I)	Data Visualization with Tableau (Elective - I)
		Data Analytics with R (Elective - II)	Data Analytics with R (Elective - II)	Data Analytics with R (Elective - II)
		Capstone Project	Capstone Project	Capstone Project
Term 7		Bonus: Industrial ML (ML – 4 & 5)	Basics of AI, TensorFlow, and Keras	Basics of AI, TensorFlow, and Keras
Term 8			Deep Learning Foundation	Deep Learning Foundation
Term 9			NPL – I/CV – I	CV – I
Term 10			NLP – II/CV – II	NLP – I
			Capstone Project	Capstone Project
Term 11				CV – II
Term 12				NLP – II
				NLP – III + CV – III
				AutoVision & AutoNLP
				Building AI product

# Term Context

- Decision Tree ← You are here...
- Random Forest
- Principal Component Analysis
- Naïve Bayes Classifier

# 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

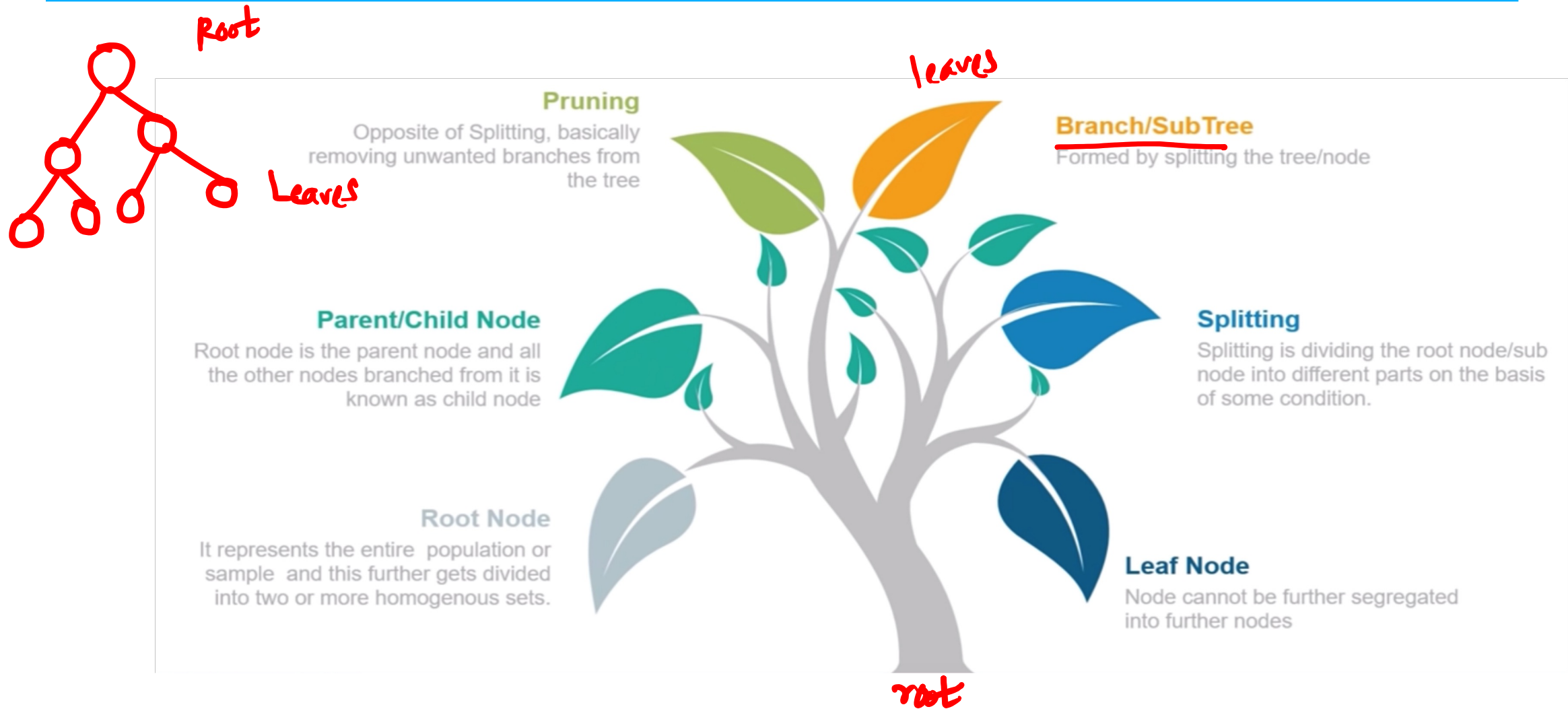


Tree Pruning



Decision Tree Applications

# Terminology Related to Trees



# Agenda

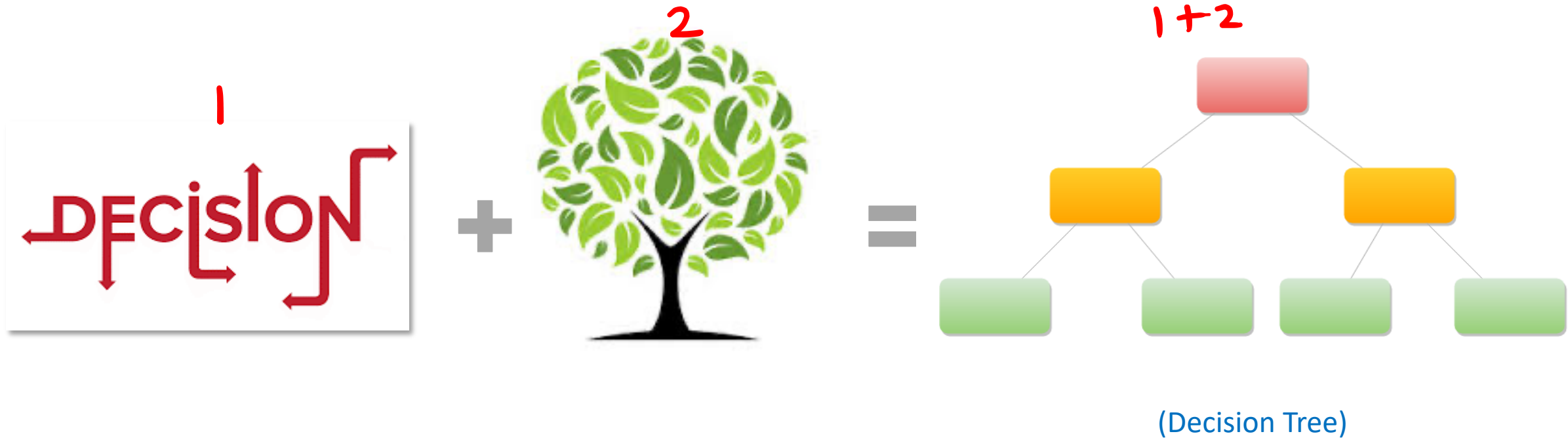
---

- ☐ 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

# Decision Tree

$$X_1 + X_2 + \dots + X_n \approx y$$

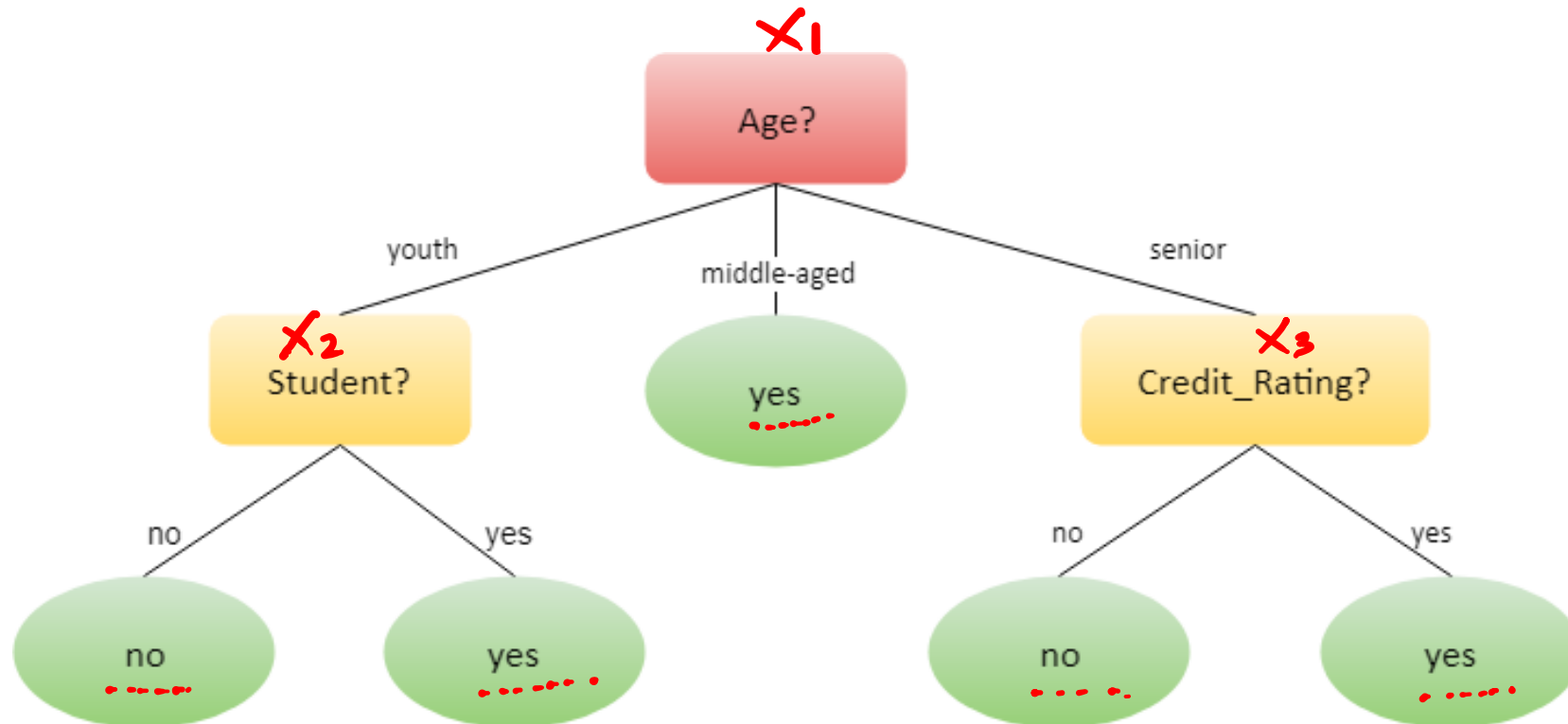
- It is a type of Supervised Learning algorithm which works for both Categorical as well as Continuous data.
- A **Decision Tree** is a graphical representation of all the possible solutions based on certain conditions.
- These solutions can be seen as IF-THEN rules.



# Decision Tree Example 1

$x_1$  Age,  $x_2$  Student,  $x_3$  Credit\_Rating  $\sim$  Buys Comp  $y$

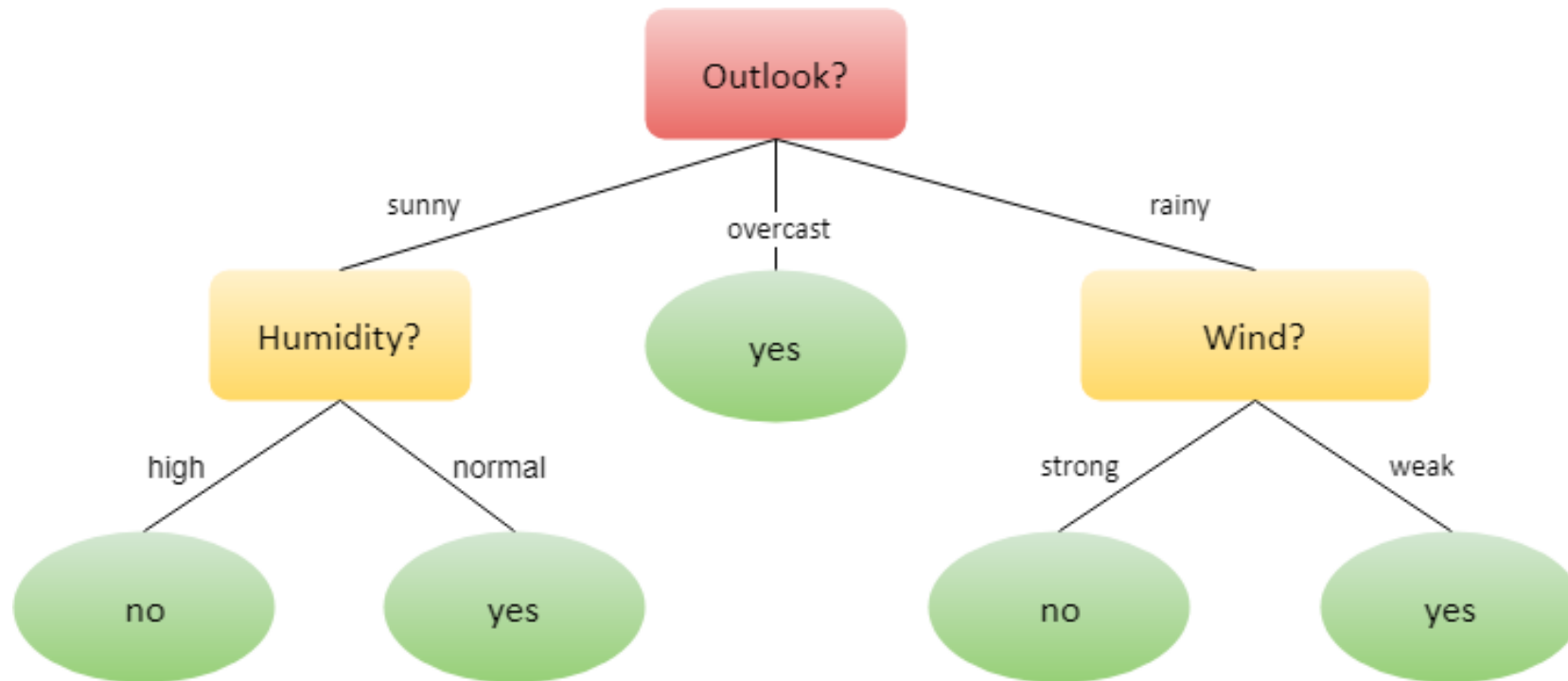
- Buys-Computer Classification: Based on certain factors whether a person will buy computer or not.





# Decision Tree Example 2

- **Play Golf Classification:** Whether a person will **play** or **not** based on environmental factors.



# Agenda

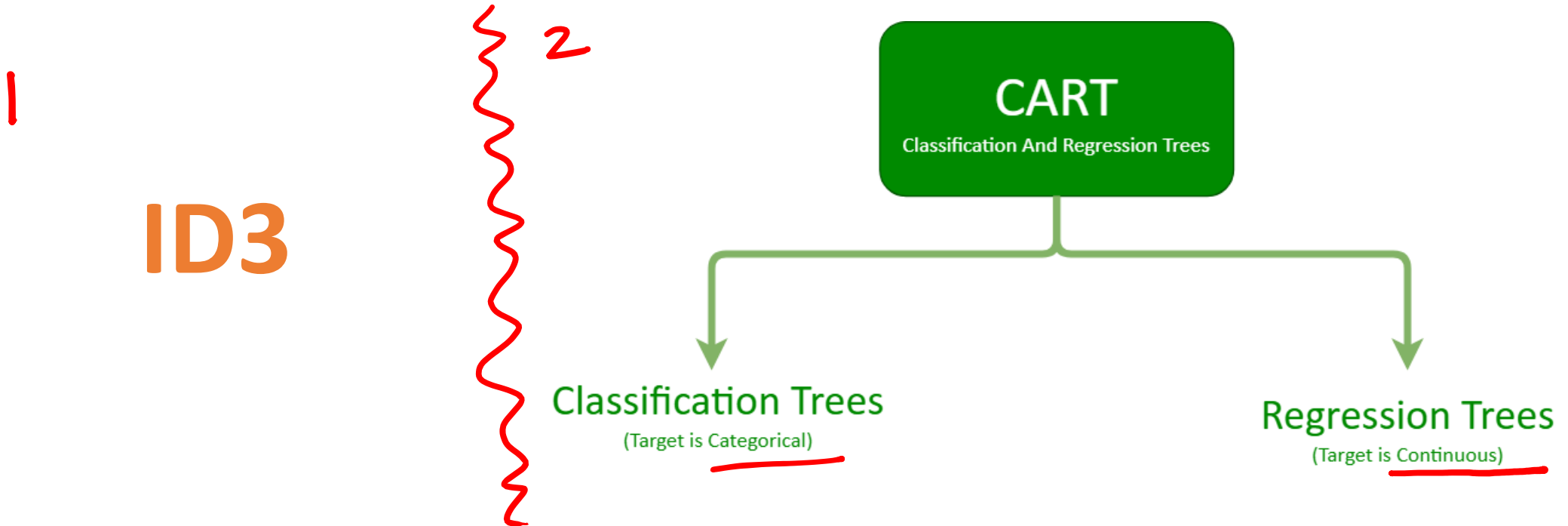
---

- ☐ 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

# Decision Tree Algorithms

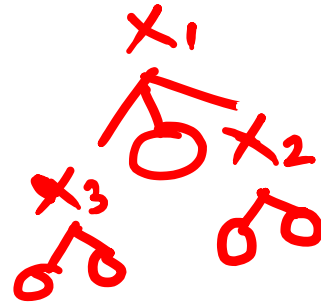
---

- 1 • ID3 – Also known as Iterative Dichotomiser 3.
- 2 • CART – Also known as Classification and Regression Trees.



# Agenda

- ☐ Terminology Related to Trees
- ☐ Decision Tree
- ☐ Decision Tree Algorithms  $X_1 + X_2 + X_3 + X_4 + N$
- ☒ Column 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



# Attribute Selection Measures

---



You can use –

- Information Gain
- Gini Index



# 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

# ID3 Algorithm

---

- ID3 uses Information Gain and Entropy as its attribute selection measure.



# 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



# 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  $\neq 0$

Heterogeneous Data



# Agenda

---

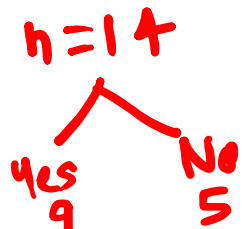
- ☐ Terminology Related to Trees
- ☐ Decision Tree
- ☐ Decision Tree Algorithms
- ☐ Attribute Selection Measures
- ☐ ID3 Algorithm
- ☐ Entropy & Information Gain
- ☒ Steps to Estimate Entropy & Information Gain
- ☐ CART Algorithm
- ☐ Gini Index
- ☐ Steps to estimate Gini Index
- ☐ CART – Regression Example
- ☐ Issues with Decision Trees
- ☐ Tree Pruning
- ☐ Decision Tree Applications

Math begins. Focus: Understanding formulae & how are the values put in <sup>numerator</sup> that formulae.

# Steps to estimate Entropy & Information Gain

$$E = - \sum_{i=1}^n p_i \cdot \log_2 p_i$$

- ④ Calculate the expected information needed to classify a tuple in data set (D) is given by –  $= - \sum_{i=1}^n p_i \cdot \log_2 p_i$



$y \rightarrow \text{Unique values} : \{ \text{yes}, \text{no} \}$

$$\text{Entropy}(D) = - \sum_{i=1}^m p_i \log_2(p_i) = - [p_y \cdot \log_2 p_y + p_n \cdot \log_2 p_n]$$

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

	$x_1$	$x_2$	$x_3$	$x_4$	$y$
	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

$$= - p_y \cdot \log_2 p_y - p_n \cdot \log_2 p_n$$

$$= - \frac{9}{14} \cdot \log_2 \frac{9}{14} - \frac{5}{14} \cdot \log_2 \frac{5}{14}$$

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

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

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

# Steps to estimate Entropy & Information Gain Cont.

$n=14$

4 youth: 2 Unique values {4 yes, 0 no}

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

This amount is measured by

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

$$E_{\text{youth}} = - \sum_{i=1}^2 p_i \cdot \log_2 p_i$$

$$= - [p_y \cdot \log_2 p_y + p_n \cdot \log_2 p_n]$$

$$= - p_y \cdot \log_2 p_y - p_n \cdot \log_2 p_n$$

$$= - \frac{2}{5} \cdot \log_2 \frac{2}{5} - \frac{3}{5} \cdot \log_2 \frac{3}{5}$$

- 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

Ex,

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

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

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

# Steps to estimate Entropy & Information Gain Cont.

- Final step is to calculate Information Gain –



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

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

# Example on Buys\_Computer Data set

	A	B <sup>X<sub>2</sub></sup>	C <sup>X<sub>3</sub></sup>	D <sup>X<sub>4</sub></sup>	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

$E_y$  [ Entropy(D) =  $-9/14 \times \log_2(9/14) - 5/14 \times \log_2(5/14) = 0.94$

$E_{x1}$  { Entropy<sub>age</sub>(D) =  $\frac{|D_{youth}|}{|D|} \times \text{Entropy}(D_{youth}) + \frac{|D_{middle}|}{|D|} \times \text{Entropy}(D_{middle}) + \frac{|D_{senior}|}{|D|} \times \text{Entropy}(D_{senior})$

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

Entropy<sub>age</sub>(D) = 0.629

$E_{x2}$  { Entropy<sub>income</sub>(D) =  $\frac{|D_{high}|}{|D|} \times \text{Entropy}(D_{high}) + \frac{|D_{medium}|}{|D|} \times \text{Entropy}(D_{medium}) + \frac{|D_{low}|}{|D|} \times \text{Entropy}(D_{low})$

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

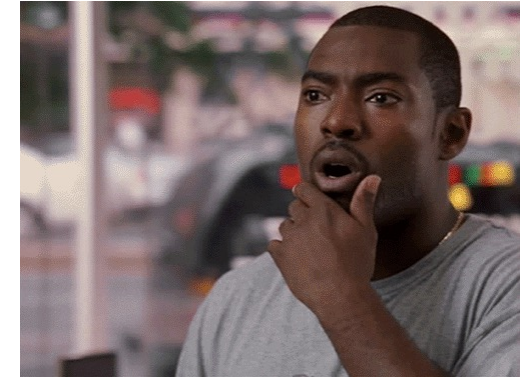
Entropy<sub>income</sub>(D) = 0.908

$E_{x3}$  { Entropy<sub>student</sub>(D) =  $\frac{|D_{yes}|}{|D|} \times \text{Entropy}(D_{yes}) + \frac{|D_{no}|}{|D|} \times \text{Entropy}(D_{no})$

Entropy<sub>student</sub>(D) = 0.786

$E_{x4}$  { Entropy<sub>credit\_rating</sub>(D) =  $\frac{|D_{fair}|}{|D|} \times \text{Entropy}(D_{fair}) + \frac{|D_{excellent}|}{|D|} \times \text{Entropy}(D_{excellent})$

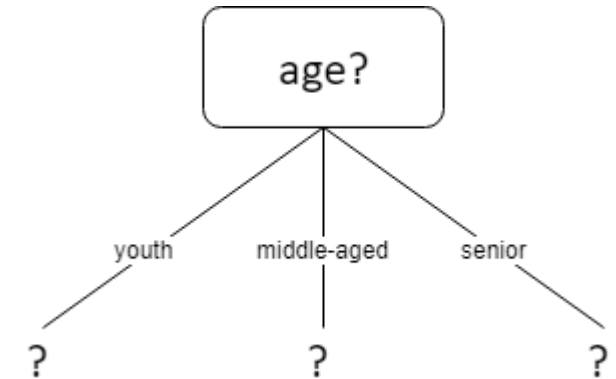
Entropy<sub>credit\_rating</sub>(D) = 0.89



# Tabular View

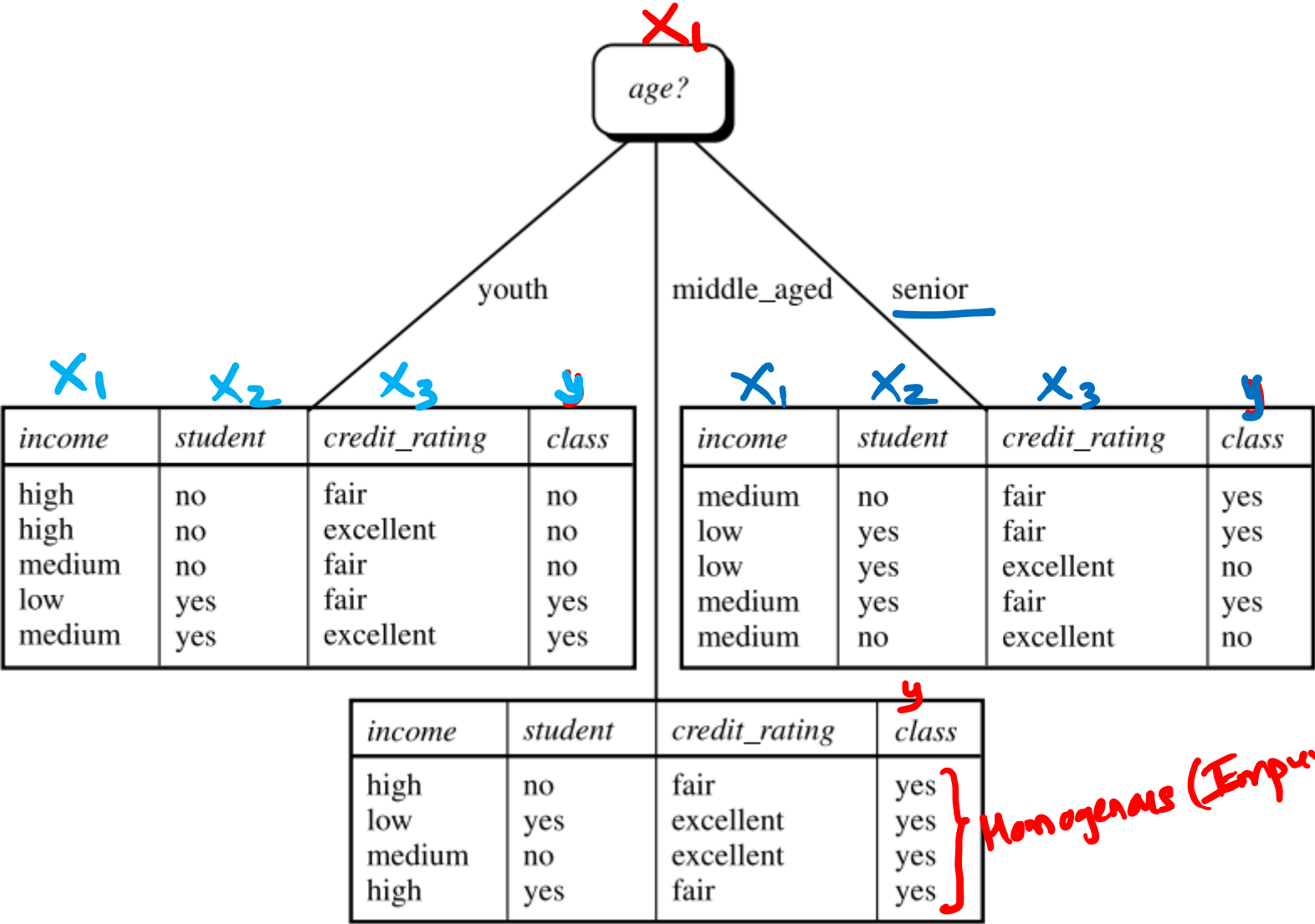
- The **age** attribute is **giving maximum information** gain. So the **root** node **will be age**.

	A	B	C
1		Entropy	Information Gain
2	Data	0.94	0
3	age (root)	0.629	= 0.248 (max)
4	income	0.908	= 0.032
5	student	0.786	= 0.154
6	credit_rating	0.89	= 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

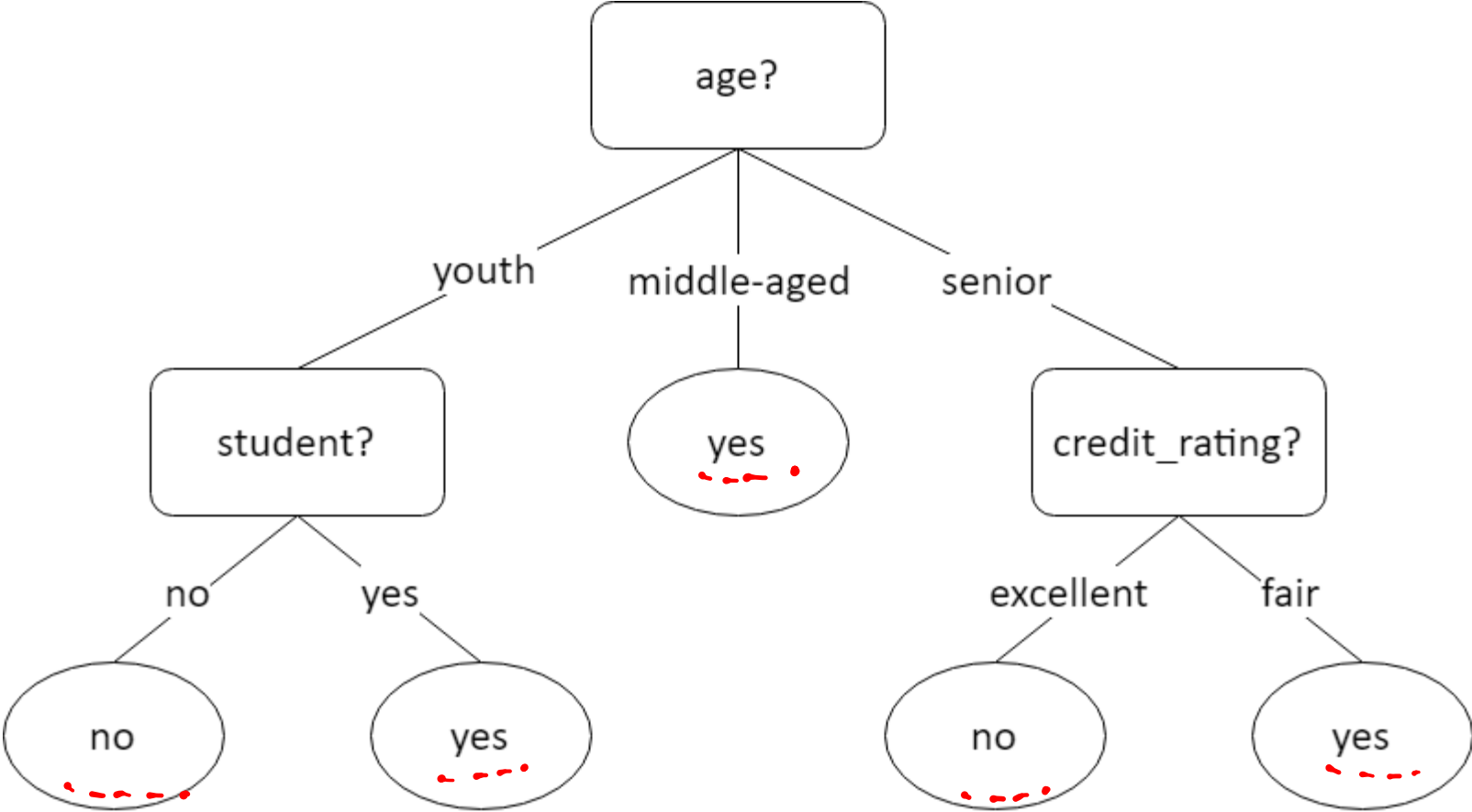
12		Entropy	Information Gain
13	Data $y$	0.97	0
14	income $x_1$	0.399	= 0.571
15	student $x_2$	0	= 0.97 (Max)
16	credit_rating $x_3$	0.948	= 0.022

Solution for senior data

21		Entropy	Information Gain
22	Data $y$	0.968	0
23	income $x_1$	0.95	= 0.018
24	student $x_2$	0.95	= 0.018
25	credit_rating $x_3$	0	= 0.968 (Max)

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

# Decision Tree Complete View





Thanks for watching