DATA SCIENCE

DATE : 30.07.2024
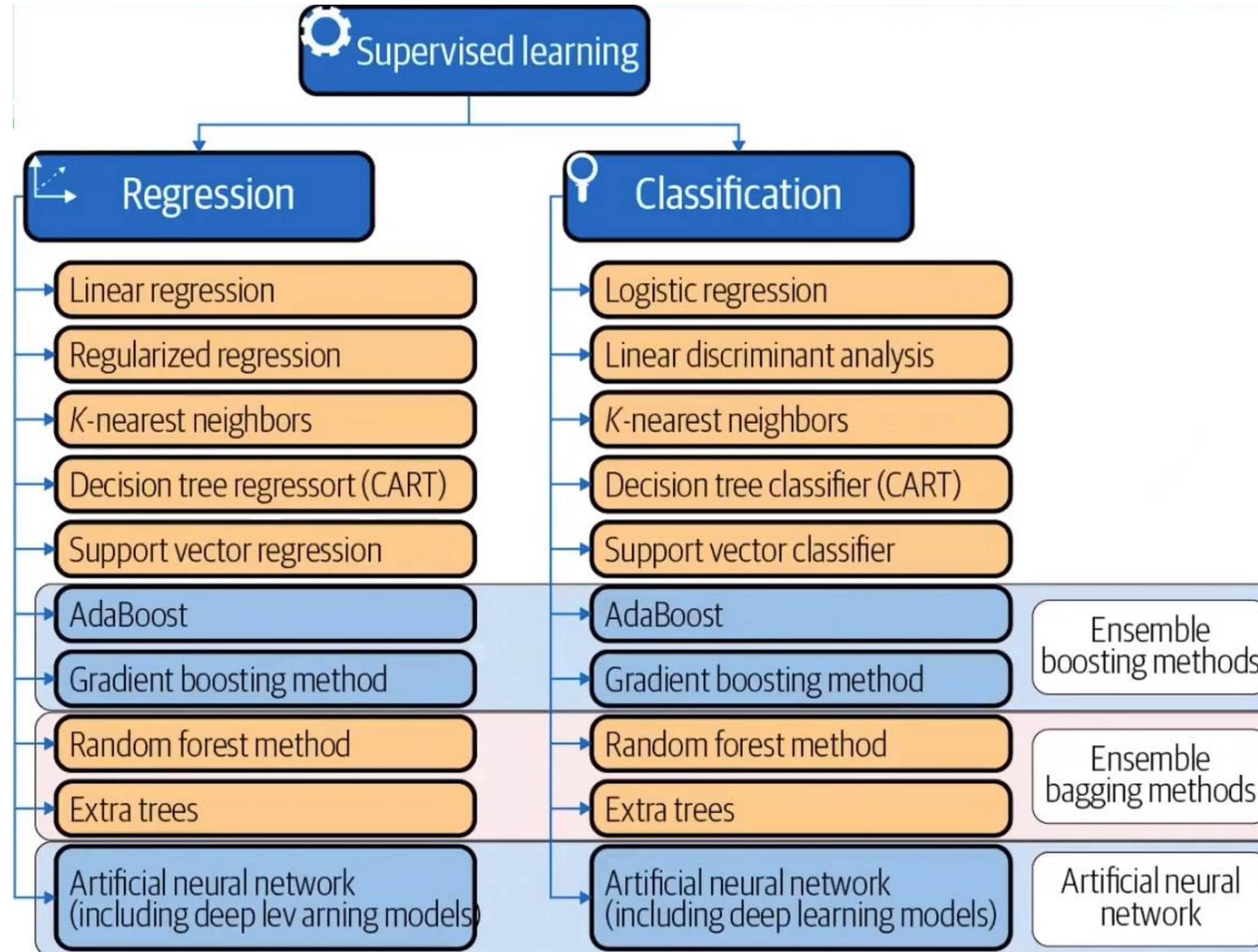
DT/NT :

LESSON : MACHINE LEARNING
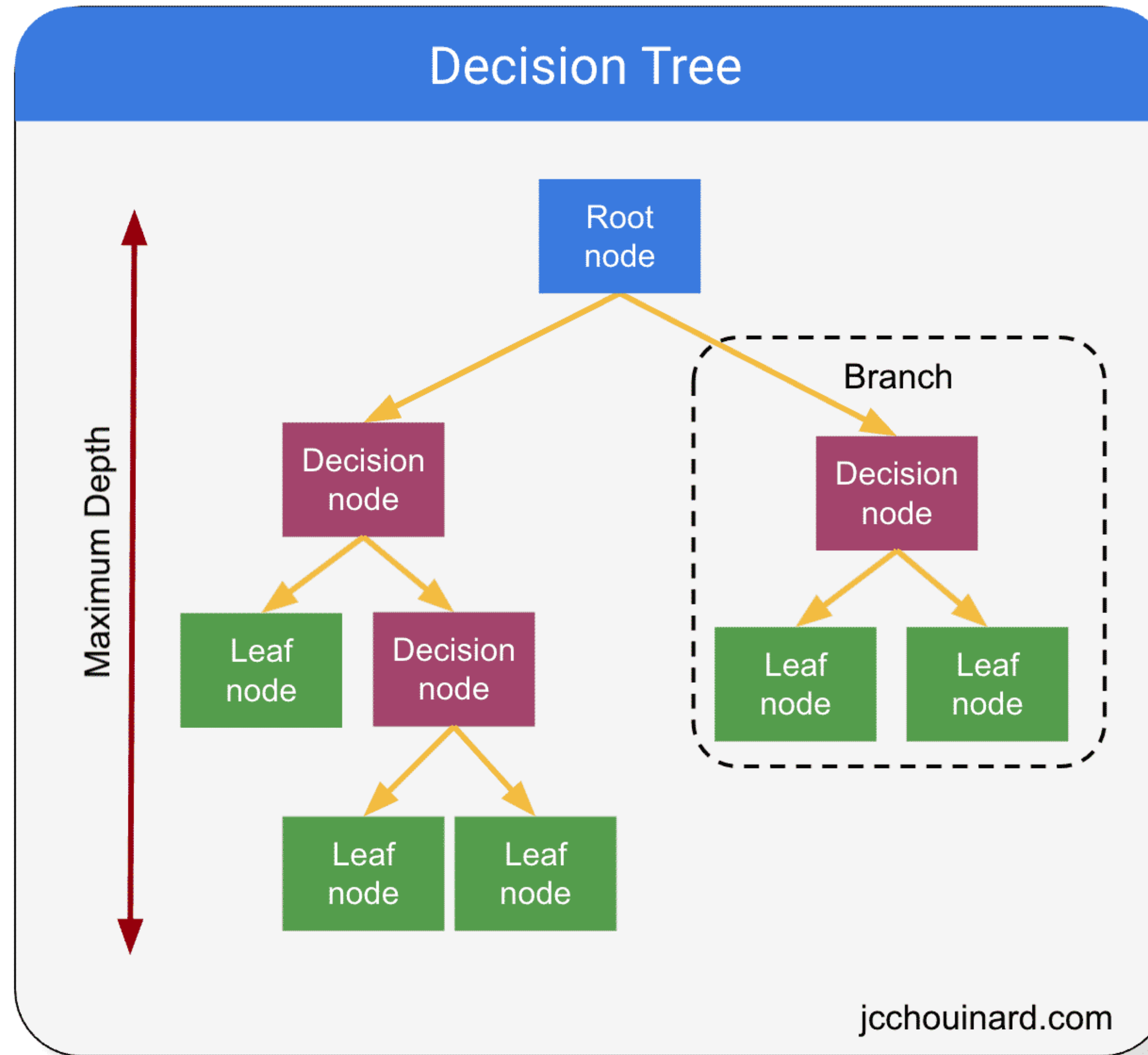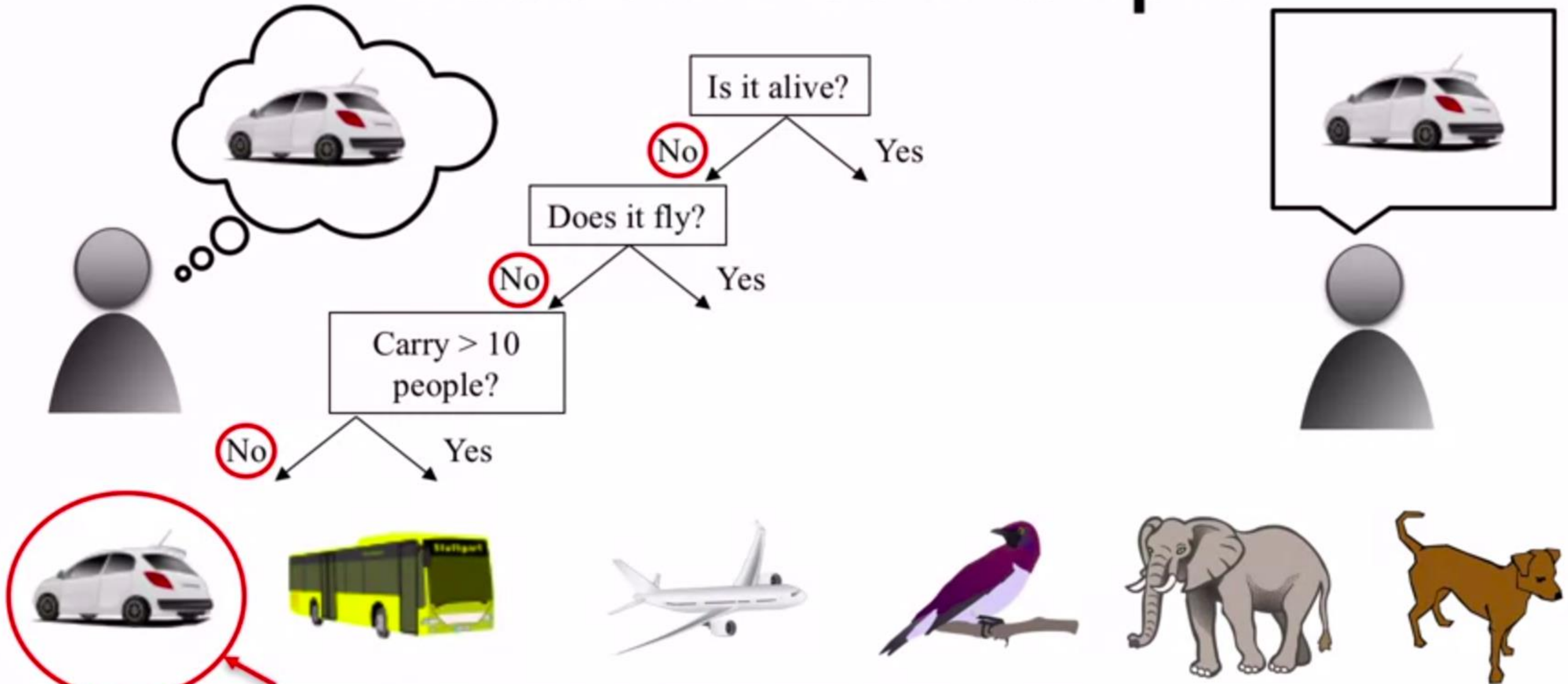
SUBJECT : DECISION TREE (CART)
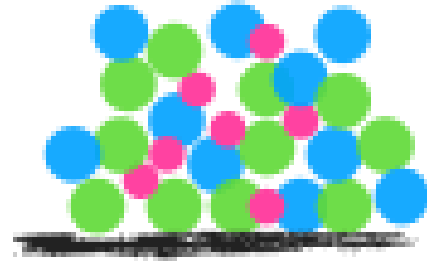
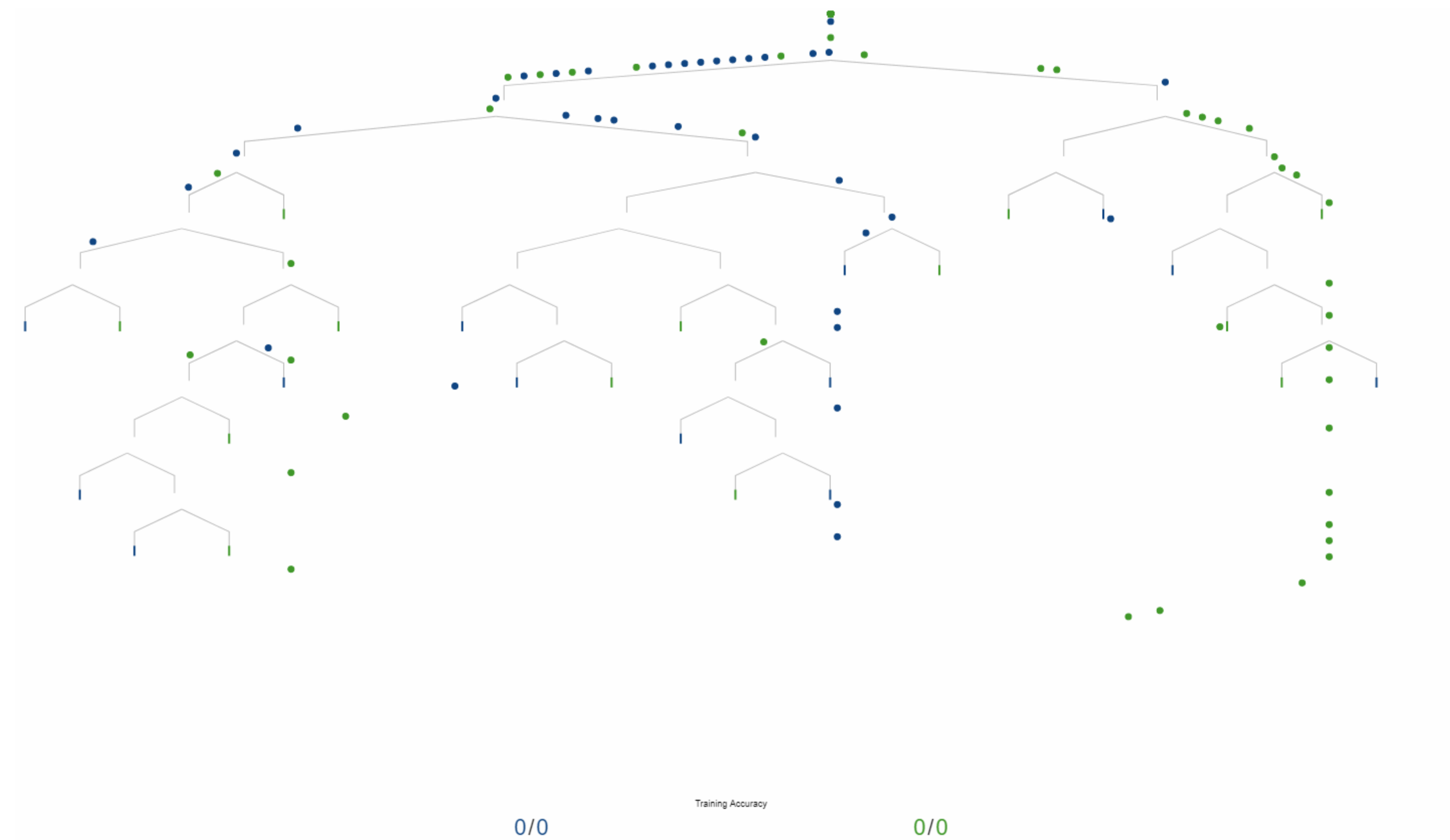BATCH : 247

# DECISION TREE

# Where We Are?

# Decision Tree

# Decision Tree Example

Color == yellow?

True — False

Height=10?

Height<10?

True — False

True — False

diameter size

colour

Training Accuracy

0/0                    0/0

# TEST DATA

# DECISION TREE



node #0
petal width (cm) ≤ 0.8
gini = 0.667
samples = 144
value = [48, 48, 48]
class = setosa

node #1
gini = 0.0
samples = 48
value = [48, 0, 0]
class = setosa

node #2
petal width (cm) ≤ 1.75
gini = 0.5
samples = 96
value = [0, 48, 48]
class = versicolor

node #3
petal length (cm) ≤ 4.95
gini = 0.174
samples = 52
value = [0, 47, 5]
class = versicolor

node #12
petal length (cm) ≤ 4.85
gini = 0.044
samples = 44
value = [0, 1, 43]
class = virginica

node #4
petal width (cm) ≤ 1.65
gini = 0.043
samples = 46
value = [0, 45, 1]
class = versicolor

node #7
petal width (cm) ≤ 1.55
gini = 0.444
samples = 6
value = [0, 2, 4]
class = virginica

node #13
sepal length (cm) ≤ 5.95
gini = 0.444
samples = 3
value = [0, 1, 2]
class = virginica

node #16
gini = 0.0
samples = 41
value = [0, 0, 41]
class = virginica

node #5
gini = 0.0
samples = 45
value = [0, 45, 0]
class = versicolor

node #6
gini = 0.0
samples = 1
value = [0, 0, 1]
class = virginica

node #8
gini = 0.0
samples = 3
value = [0, 0, 3]
class = virginica

node #9
sepal length (cm) ≤ 6.95
gini = 0.444
samples = 3
value = [0, 2, 1]
class = versicolor

node #14
gini = 0.0
samples = 1
value = [0, 1, 0]
class = versicolor

node #15
gini = 0.0
samples = 2
value = [0, 0, 2]
class = virginica

node #10
gini = 0.0
samples = 2
value = [0, 2, 0]
class = versicolor

node #11
gini = 0.0
samples = 1
value = [0, 0, 1]
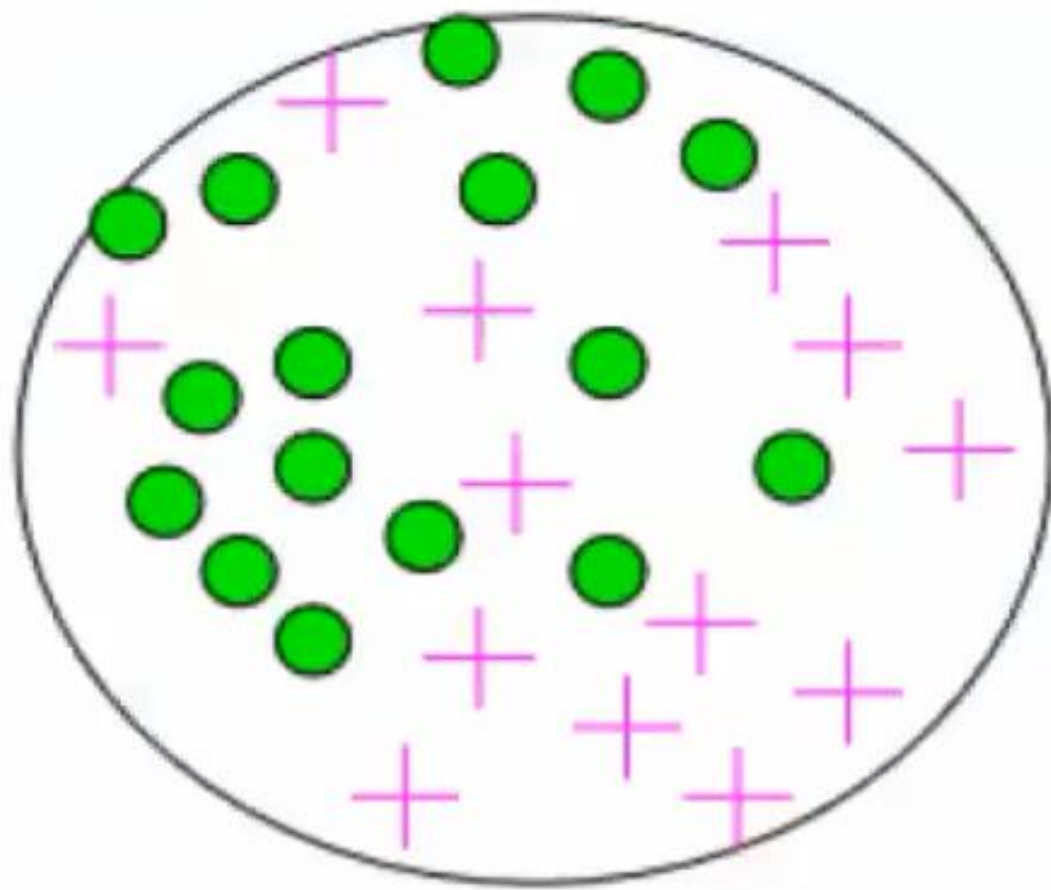class = virginica

# Gini Impurity



Gini Impurity

Max Gini Impurity = 1 - 1/n
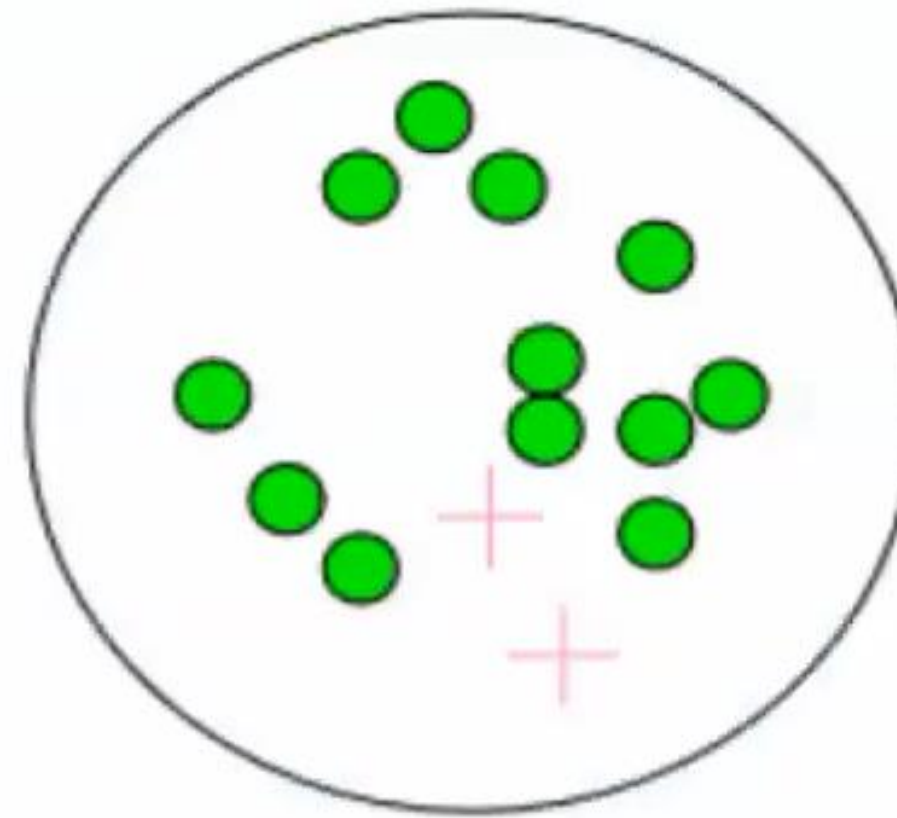n = number of classes
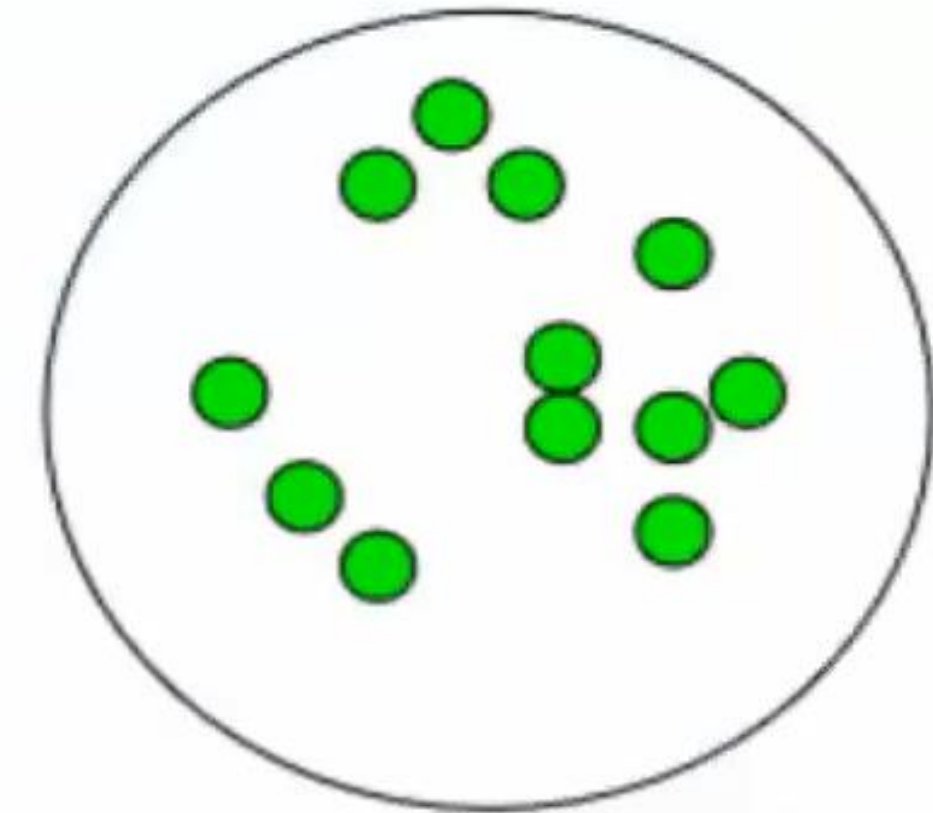Max gini impurity is **0.5** for binary class

**Very impure group**

$$1 - [(16/29)^2 + (13/29)^2] = \textbf{0.4946}$$

**Less impure**

$$1 - [(12/14)^2 + (2/14)^2] = \textbf{0.244}$$

**Minimum impurity**

$$1 - [(12/12)^2 + (0/12)^2] = \textbf{0}$$

# Gini Impurity



**Gini Index Formula**

$$Gini = 1 - \sum_{i=1}^{n} P^2(x_i)$$

# of classes

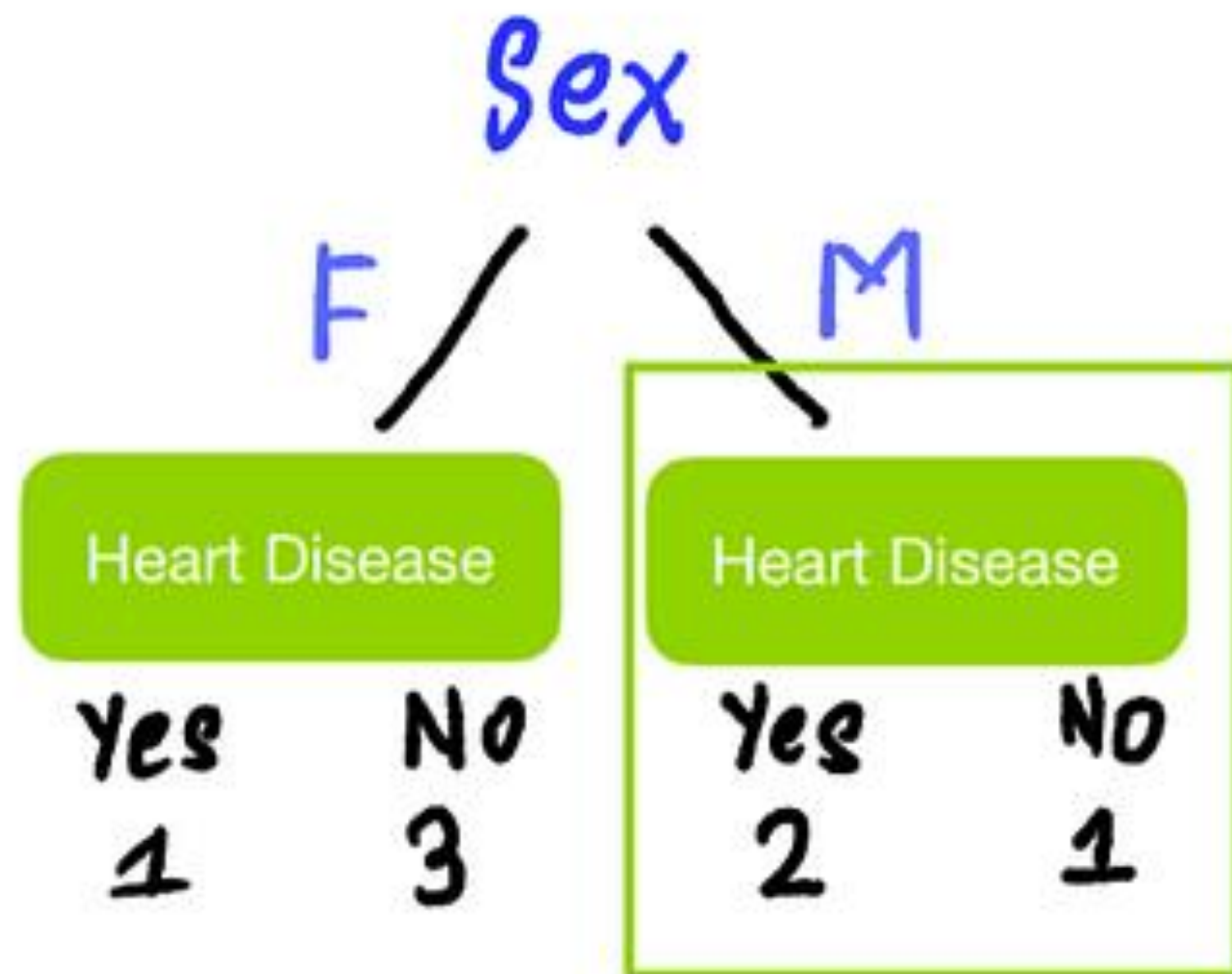Probability of class "i"

Another commonly used formula is:

$$Gini\ impurity = 1 - \Sigma\left(p(i) * \left(1 - p(i)\right)\right)$$

# Gini Impurity

## Sex

F / M

**Heart Disease**

Yes: 1  No: 3

**Heart Disease**

Yes: 2  No: 1

$$Gini = 1 - \sum_{i=1}^{C} (p_i)^2$$

Gini Impurity of the right leaf

$$= 1 - \left(\frac{2}{2+1}\right)^2 - \left(\frac{1}{2+1}\right)^2$$
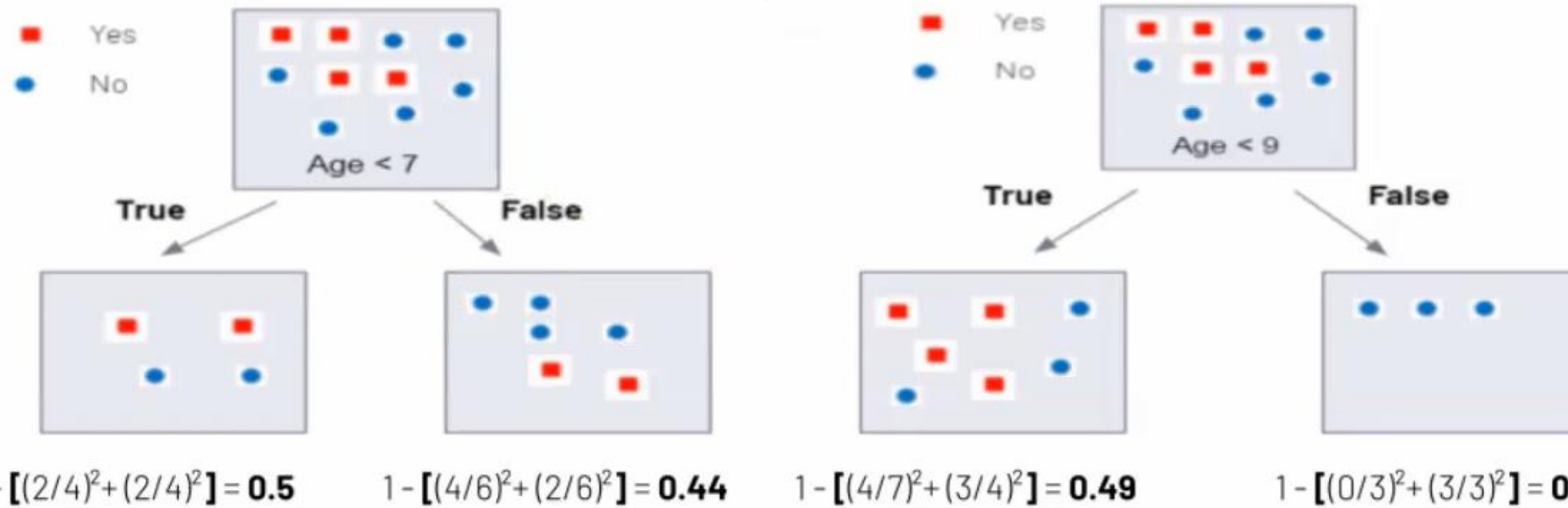
$$= 0.444$$

# Gini Impurity

| Age | Ability | Footballer (tiny stars) |
|-----|---------|------------------------|
| 10 | Yes | No |
| 8 | No | No |
| 6 | Yes | Yes |
| 6 | No | No |
| 8 | Yes | Yes |
| 8 | Yes | Yes |
| 6 | Yes | Yes |
| 10 | Yes | No |
| 10 | No | No |
| 6 | No | No |

## Selection Feature For Root Node With Gini Impurity

| Age | Ability | Footballer (tiny stars) |
|-----|---------|------------------------|
| 10 | 1 | No |
| 8 | 0 | No |
| 6 | 1 | Yes |
| 6 | 0 | No |
| 8 | 1 | Yes |
| 8 | 1 | Yes |
| 6 | 1 | Yes |
| 10 | 1 | No |
| 10 | 0 | No |
| 6 | 0 | No |

**Sorting Age** : 6 - 8 - 10
**Sorting Ability** : 0 - 1

## For Age Feature



Age < 7

True — $1 - [(2/4)^2 + (2/4)^2] = 0.5$

False — $1 - [(4/6)^2 + (2/6)^2] = 0.44$

$4/10 \times 0.5 + 6/10 \times 0.44 = 0.464$

Age < 9

True — $1 - [(4/7)^2 + (3/4)^2] = 0.49$

False — $1 - [(0/3)^2 + (3/3)^2] = 0$

$7/10 \times 0.49 + 3/10 \times 0 = 0.34$

## For Ability Feature

Ability < 0.5

True — $1 - [(0/4)^2 + (4/4)^2] = 0$

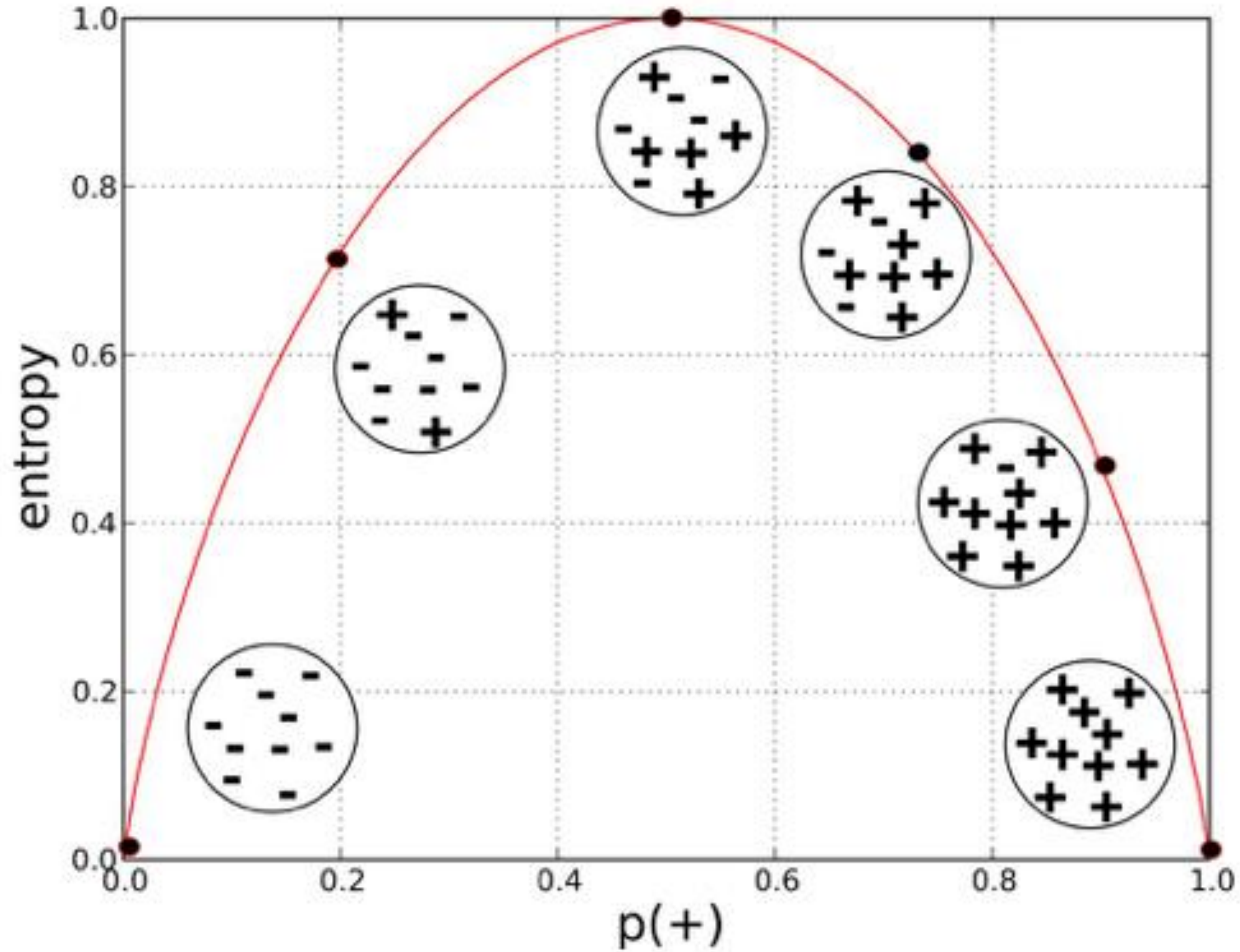False — $1 - [(4/6)^2 + (2/6)^2] = 0.44$

$4/10 \times 0 + 6/10 \times 0.44 = 0.26$

**Root Node**

# Entropy

# Entropy

The formula for Entropy is shown below:

$$E(S) = -p_{(+)} \log p_{(+)} - p_{(-)} \log p_{(-)}$$

Here,

- $p_{+}$ is the probability of positive class

- $p_{-}$ is the probability of negative class

- S is the subset of the training example

$$E(S) = \sum_{i=1}^{c} -p_i \log_2 p_i$$

TECHPRO EDUCATION

# Entropy

$$E(S) = \sum_{i=1}^{c} -p_i \log_2 p_i$$

| Play Golf | |
|---|---|
| Yes | No |
| 9 | 5 |

Entropy(PlayGolf) = Entropy (5,9)

= Entropy (0.36, 0.64)

= - (0.36 $\log_2$ 0.36) - (0.64 $\log_2$ 0.64)

= 0.94

# Entropy

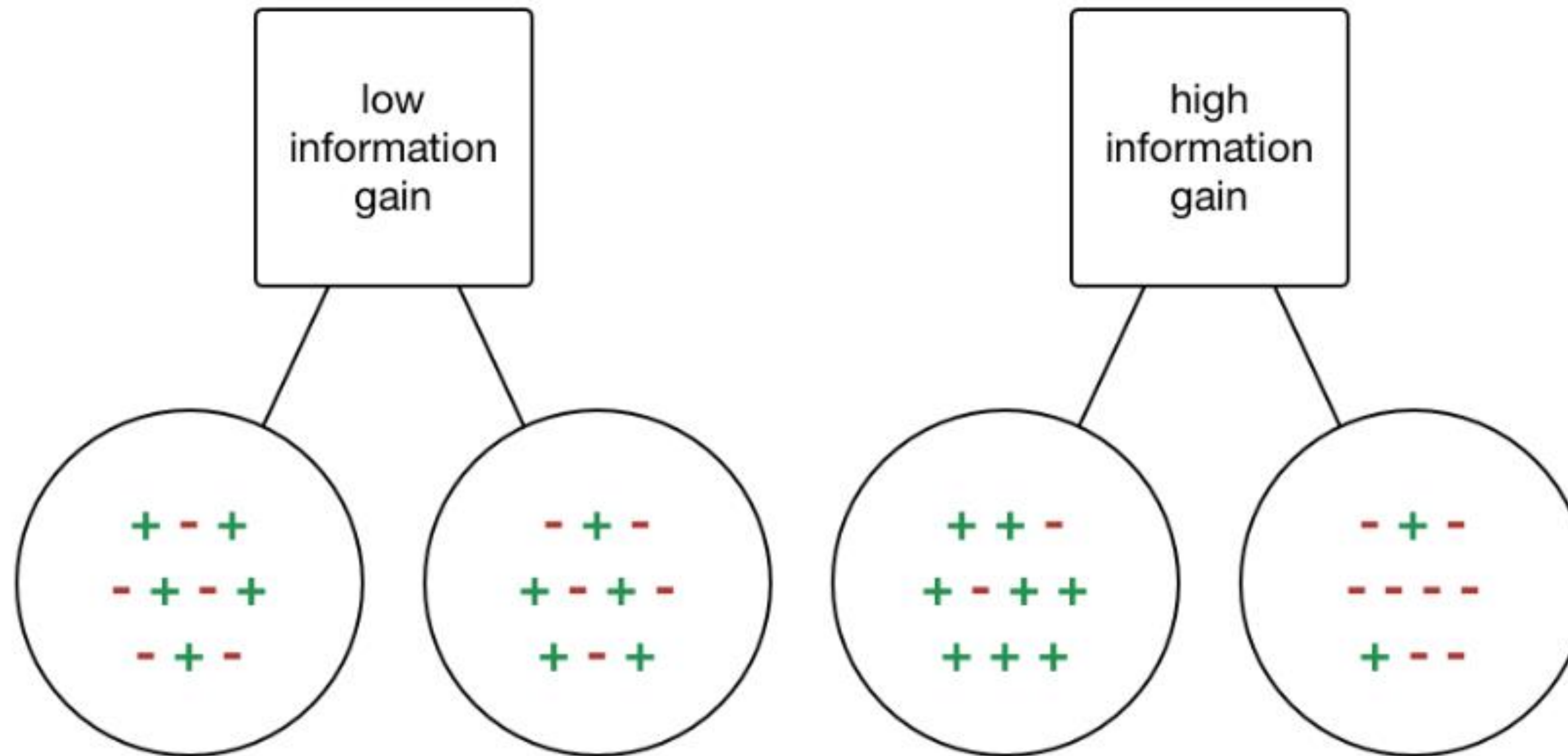$$E(T, X) = \sum_{c \in X} P(c)E(c)$$

| | | Play Golf | | |
|---|---|---|---|---|
| | | Yes | No | |
| | Sunny | 3 | 2 | 5 |
| Outlook | Overcast | 4 | 0 | 4 |
| | Rainy | 2 | 3 | 5 |
| | | | | 14 |

**E**(PlayGolf, Outlook) = **P**(Sunny)*E(3,2) + **P**(Overcast)*E(4,0) + **P**(Rainy)*E(2,3)

= (5/14)*0.971 + (4/14)*0.0 + (5/14)*0.971

= 0.693

# Information Gain

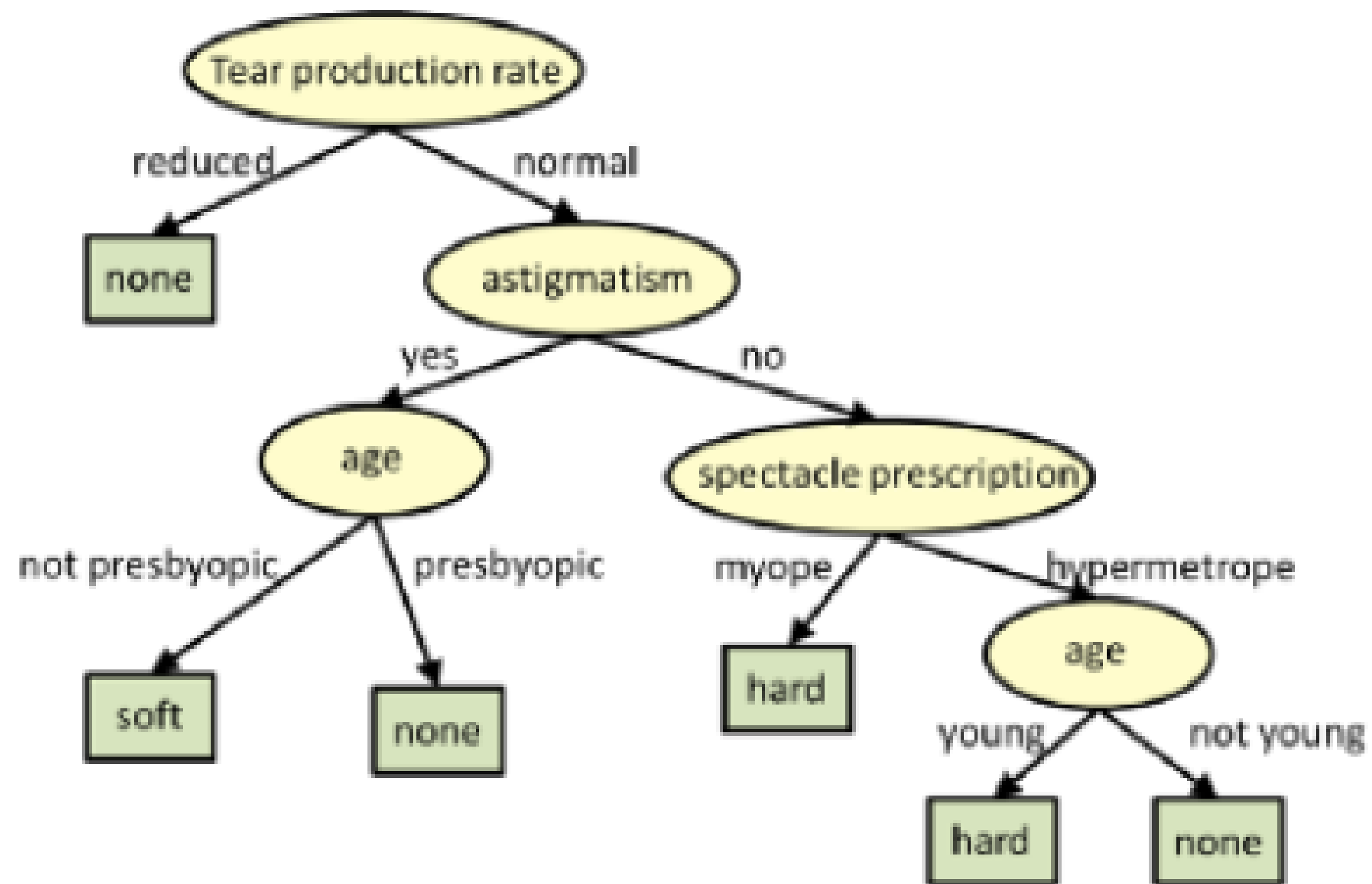# Information Gain

$$\text{Information Gain}(T, X) = \text{Entropy}(T) - \text{Entropy}(T, X)$$

$$
\begin{aligned}
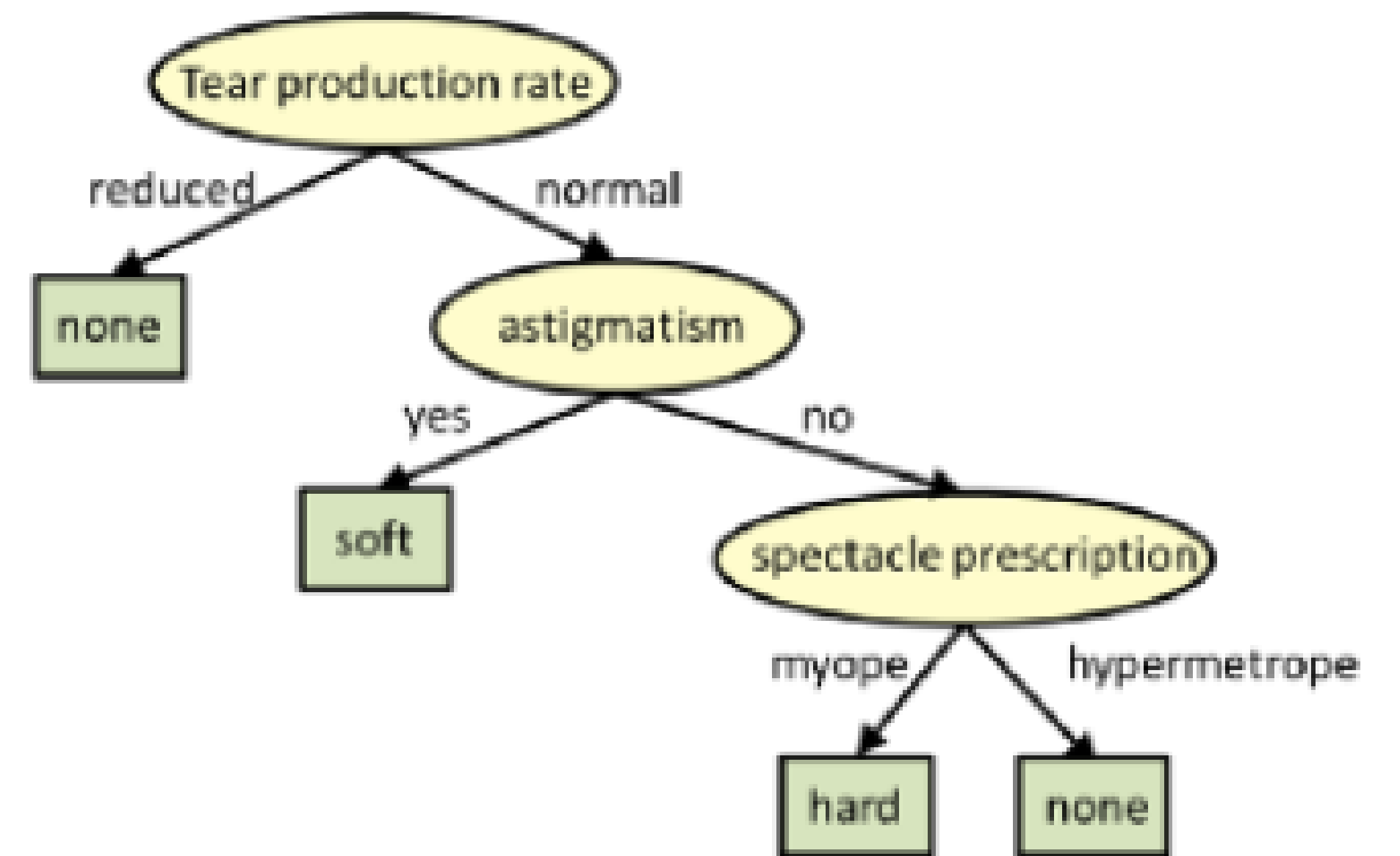IG(PlayGolf, Outlook) &= E(PlayGolf) - E(PlayGolf, Outlook) \\
&= 0.940 - 0.693 \\
&= 0.247
\end{aligned}
$$

# Pruned Tree



Original Tree

Pruned Tree

# Karar ağacının nasıl bölüneceğini nasıl belirleriz?

- Karar ağaçlarında bölünme (veya dal ayrımı), veri setindeki homojenliği maksimize etmek için yapılır.

- Yani, bir ağacın her bir dalında, sonucun mümkün olduğunca bir sınıfa özgü olması hedeflenir.

- Bu, bilgi kazancı (Information Gain), Gini saflığı (Gini Impurity), veya entropi gibi ölçümler kullanılarak belirlenir.

Tea break...

# 10:00