



NANYANG  
TECHNOLOGICAL  
UNIVERSITY  
SINGAPORE

SC1015 : Review Lecture

# Classification

Dr Smitha K G



Sample  
COLLECTION



Practical  
MOTIVATION

Data  
PREPARATION



Problem  
FORMULATION

Exploratory  
ANALYSIS



Statistical  
DESCRIPTION

Analytic  
VISUALIZATION



Pattern  
RECOGNITION

Algorithmic  
OPTIMIZATION



Machine  
LEARNING

Information  
PRESENTATION



Statistical  
INFERENCE

Ethical  
CONSIDERATION



Intelligent  
DECISION

SC1015

## Admin Announcements

1. Mini-Project details posted on NTULearn. Check the FAQs and take inspiration from the datasets.
2. Talk to your Lab TA to form your Project Teams. Deadline for team formation March 1<sup>st</sup> 5pm.

### LAMS Completion Status

Module 1 Parts 1, 2 : Above 800 – Quiz solutions posted

Module 2 Parts 1, 2 : Above 750 – Quiz solutions posted

Module 3 : Above 600 – quiz solutions will be posted by next Monday

Module 4 : Above 250 – Complete by Exercise 5 (W7)

**LAMS DS deadline:** 3<sup>rd</sup> March 11.59 pm

**DS Theory Quiz in Recess Week : 8 March, Friday.** Slots :  
12:30 pm – 2:00 pm and 2:30 pm to 4:00 pm. Lab allocations  
and FAQs posted.

Let's touch upon the basic ideas of ...

# **CLASSIFICATION**

SC1015

## Let's clarify some of these ...

Connection between Data Partitions and Decision Tree

55 (15%)

The intuition of Gini Index, and how Partitioning works

74 (20%)

How to predict Binary Classes using the Decision Tree

69 (18%)

The concept of Classification Accuracy and the Errors

75 (20%)

The idea of Multi-Variate Decision Tree and Partitions

101 (27%)

*Which part of this Lesson will you like me to review in t...*



Connection between Data Partitions and Decision Tree

59 (15%)

The intuition of Gini Index, and how Partitioning works

81 (20%)

How to predict Binary Classes using the Decision Tree

72 (18%)

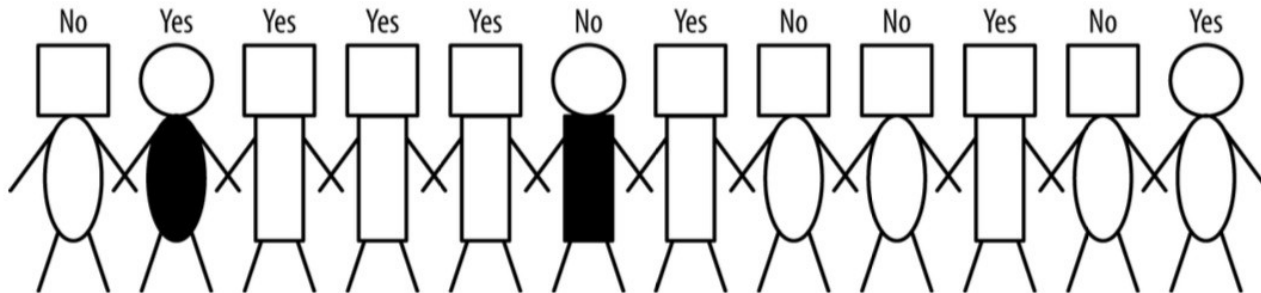
The concept of Classification Accuracy and the Errors

74 (19%)

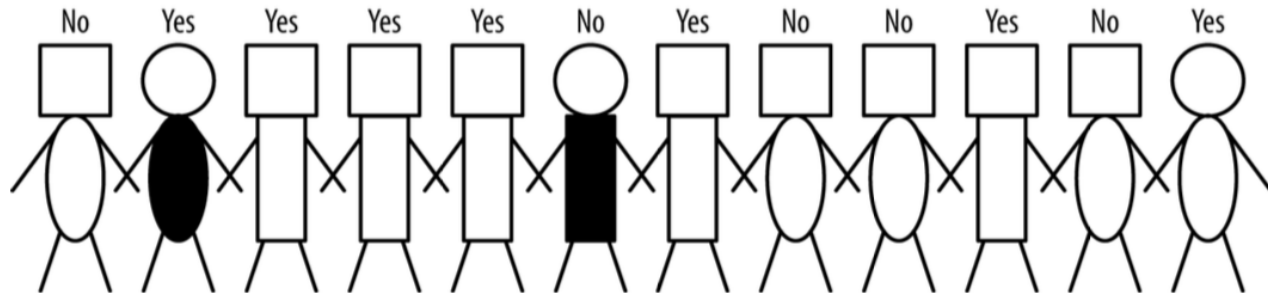
The idea of Multi-Variate Decision Tree and Partitions

113 (28%)



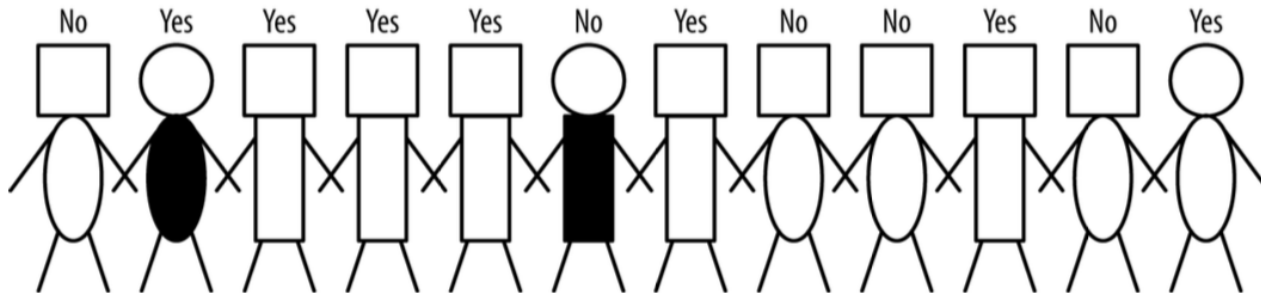


Body Shape	Body Color	Head Shape	Head Color	Funny?
Oval	White	Square	White	No
Oval	Black	Circle	White	Yes
Rectangle	White	Square	White	Yes
Rectangle	White	Square	White	Yes
Rectangle	White	Square	White	Yes
Rectangle	Black	Circle	White	No
Rectangle	White	Square	White	Yes
Oval	White	Square	White	No
Oval	White	Square	White	No
Rectangle	White	Square	White	Yes
Oval	White	Square	White	No
Oval	White	Circle	White	Yes



YES : NO

7 : 5

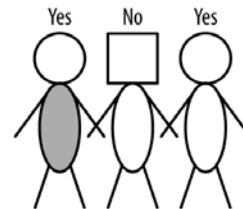
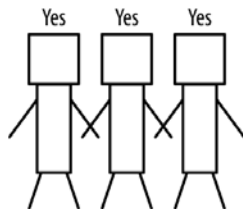


Rectangular Bodies

Oval Bodies

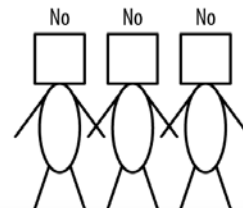
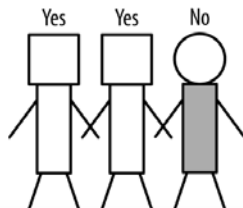
YES : NO

5 : 1

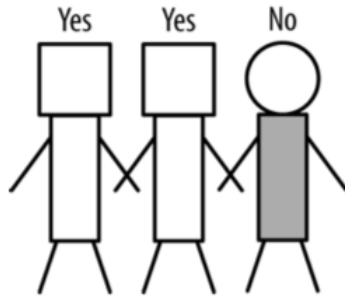
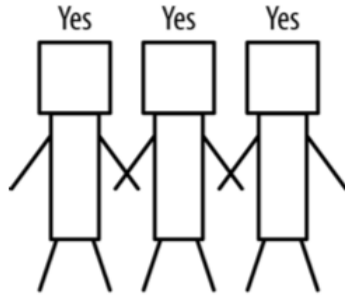


YES : NO

2 : 4



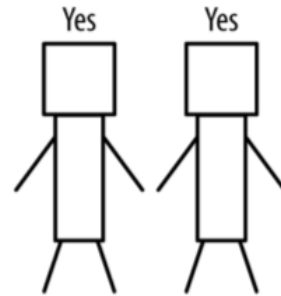
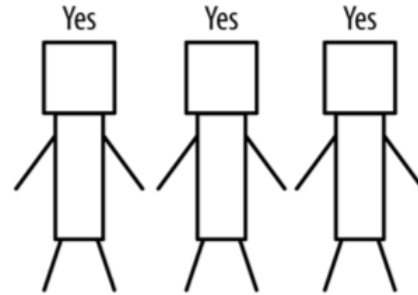
## Rectangular Bodies



YES : NO

5 : 1

## Rectangular Body and White



YES : NO

5 : 0

## Rectangular Body and Gray

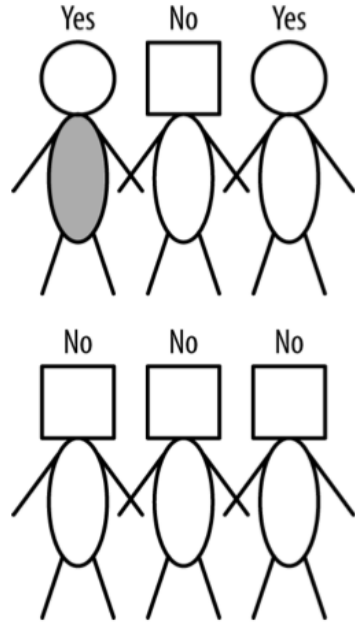


YES : NO

0 : 1



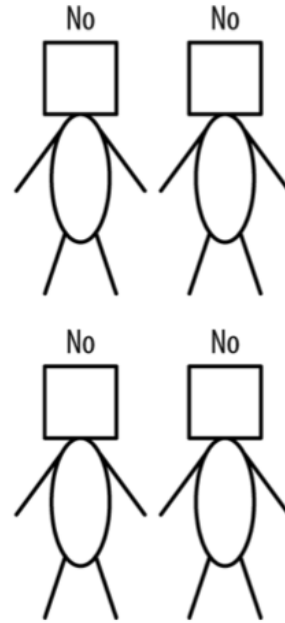
## Oval Bodies



YES : NO

**2 : 4**

## Oval Body and Square Head



YES : NO

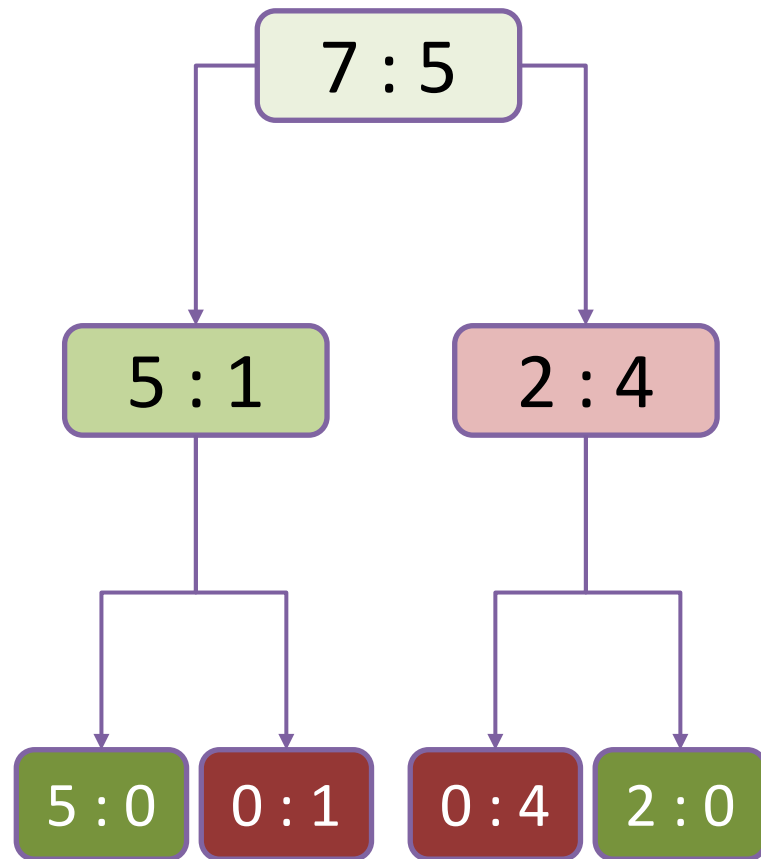
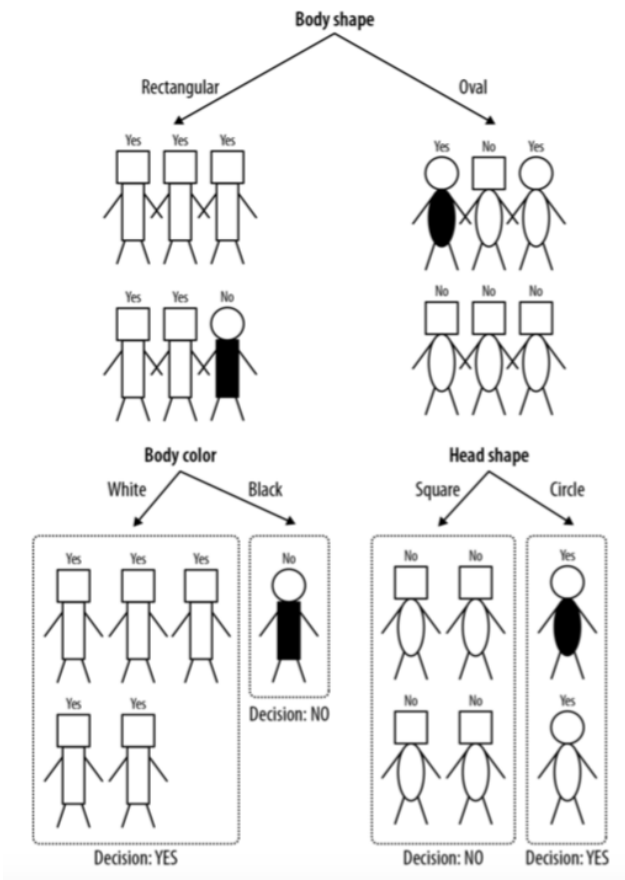
**0 : 4**

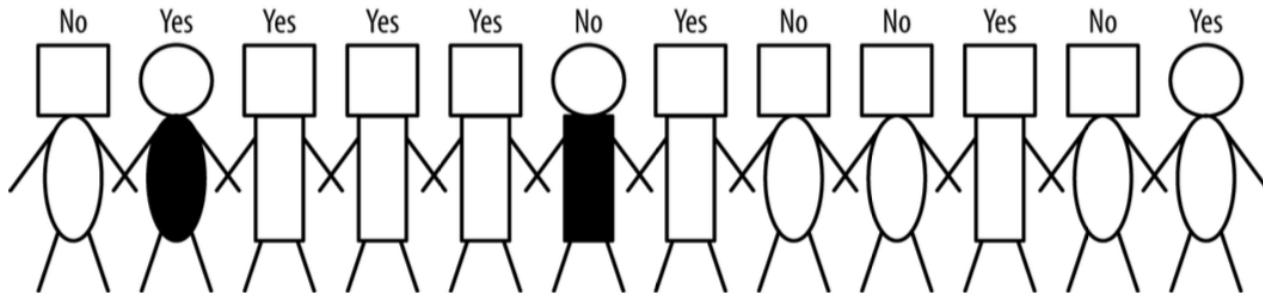
## Oval Body and Circular Head



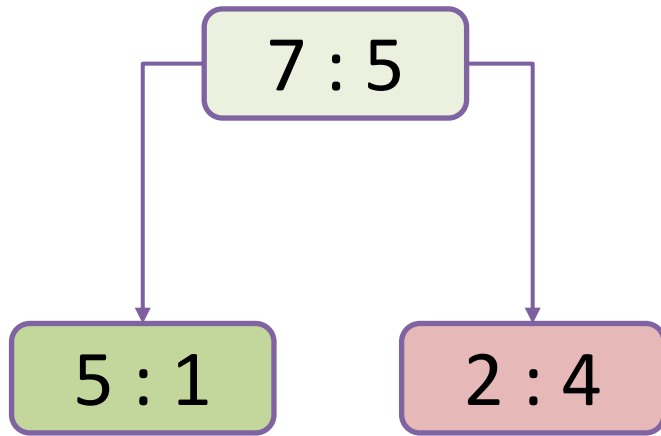
YES : NO

**2 : 0**

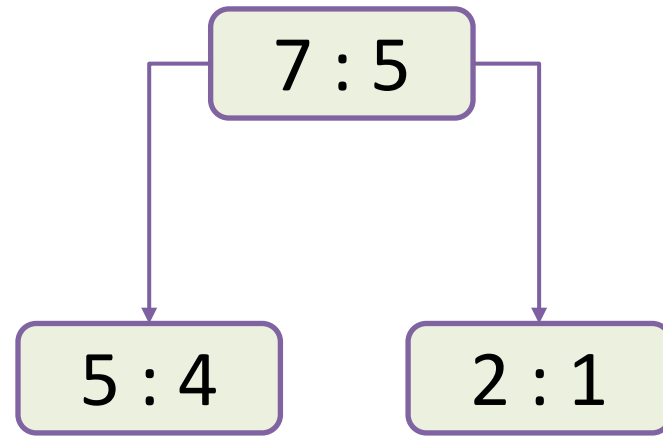




**Body : Rectangle or Oval**



**Head : Square or Circle**



**Parent** : Gini =  $1 - (7/12)^2 - (5/12)^2 = \mathbf{0.486}$

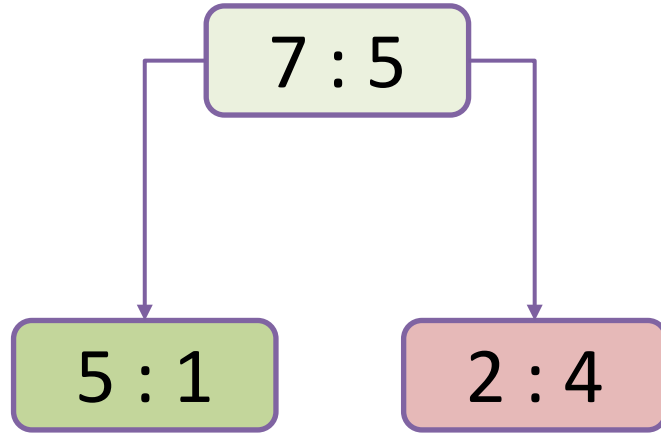
**LChild** :  $1 - (5/6)^2 - (1/6)^2 = 0.278$  | 6 samples

**RChild** :  $1 - (2/6)^2 - (4/6)^2 = 0.444$  | 6 samples

**Children** :  $0.278 \times (6/12) + 0.444 \times (6/12) = \mathbf{0.361}$

**Improvement** =  $0.486 - 0.361 = 0.125$

**Body : Rectangle or Oval**



**Parent** : Gini =  $1 - (7/12)^2 - (5/12)^2 = \mathbf{0.486}$

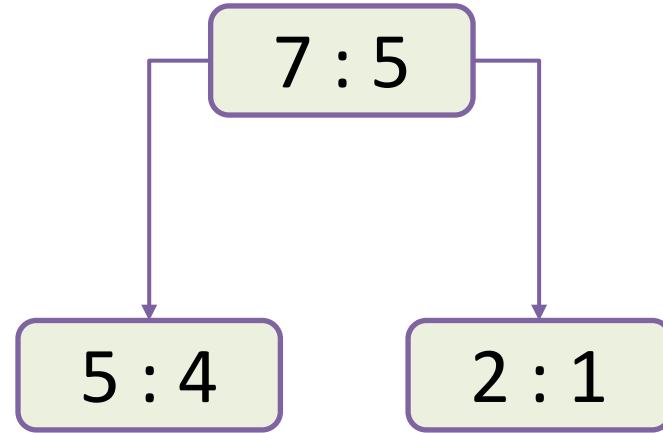
**LChild** :  $1 - (5/9)^2 - (4/9)^2 = 0.494$  | 9 samples

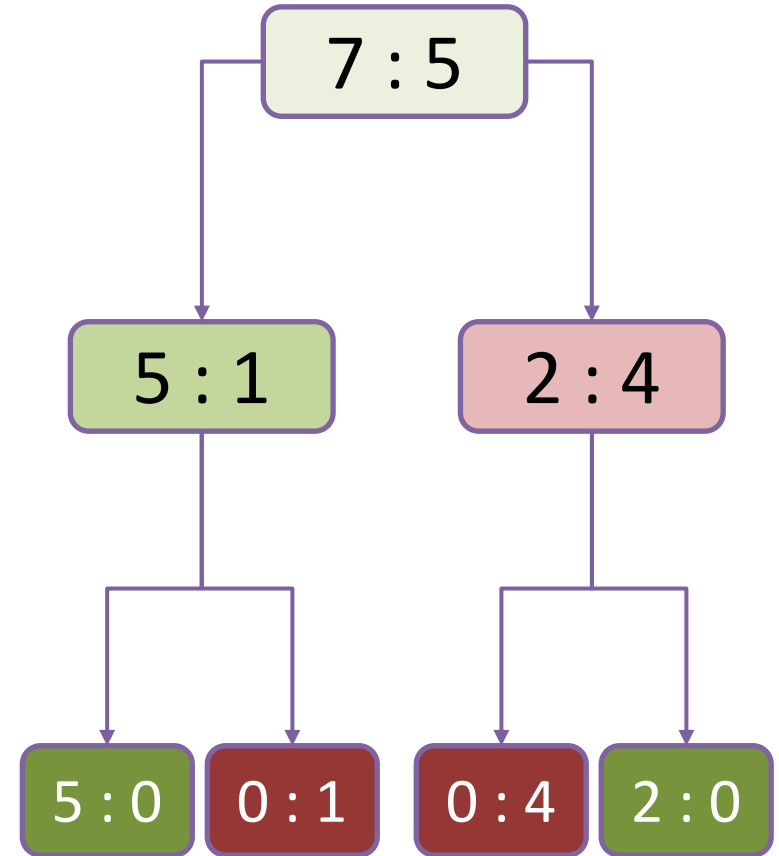
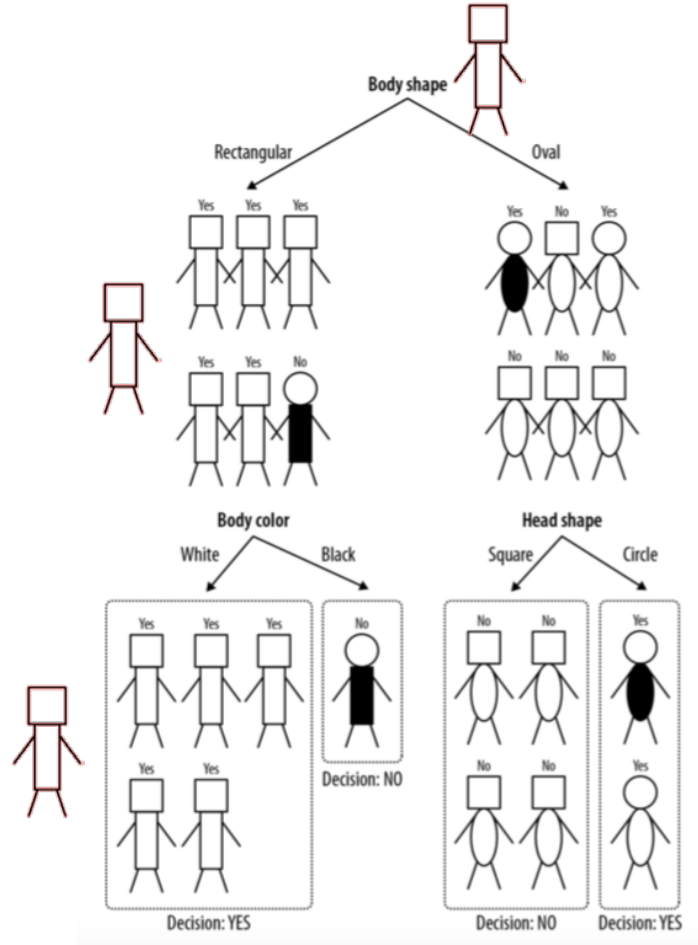
**RChild** :  $1 - (2/3)^2 - (1/3)^2 = 0.444$  | 3 samples

**Children** :  $0.494 \times (9/12) + 0.444 \times (3/12) = \mathbf{0.482}$

**Improvement** =  $0.488 - 0.482 = 0.006$

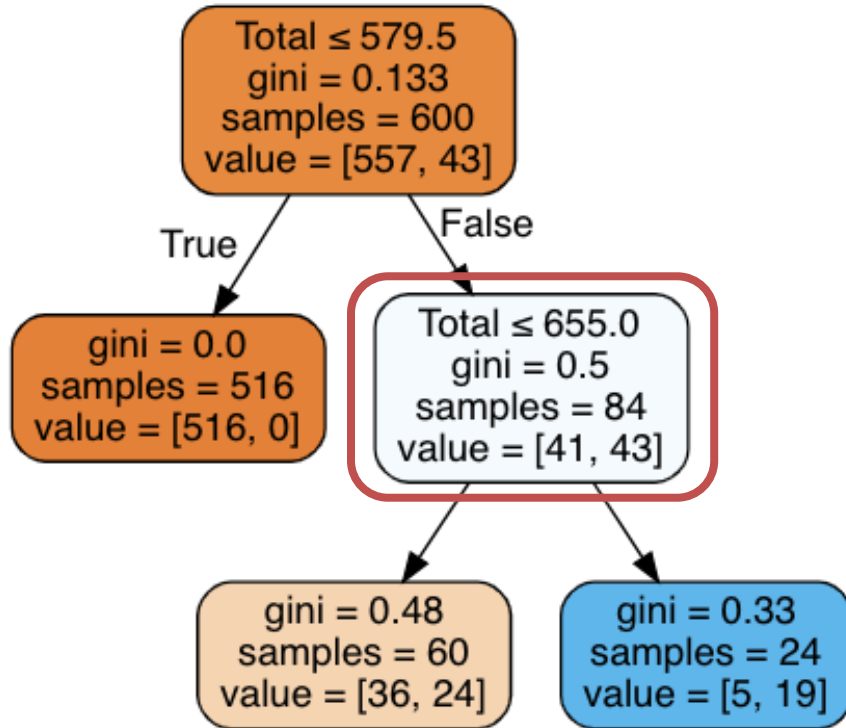
**Head : Square or Circle**





# Binary Classification

## How to “read” a Decision Tree?



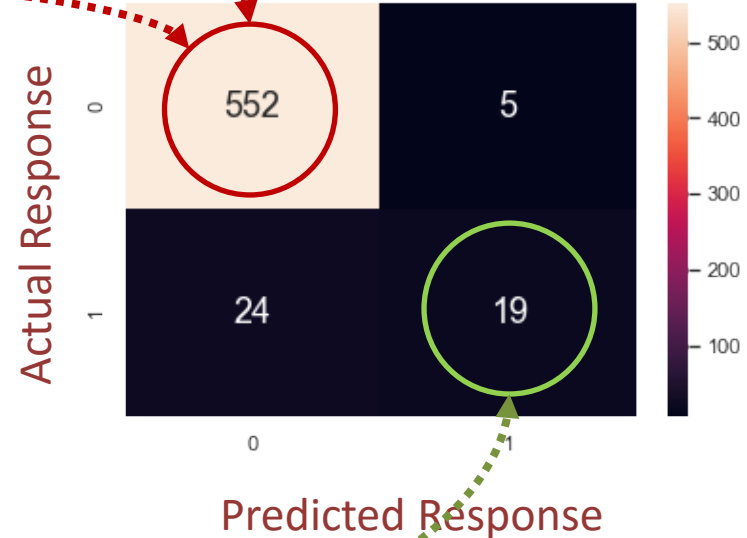
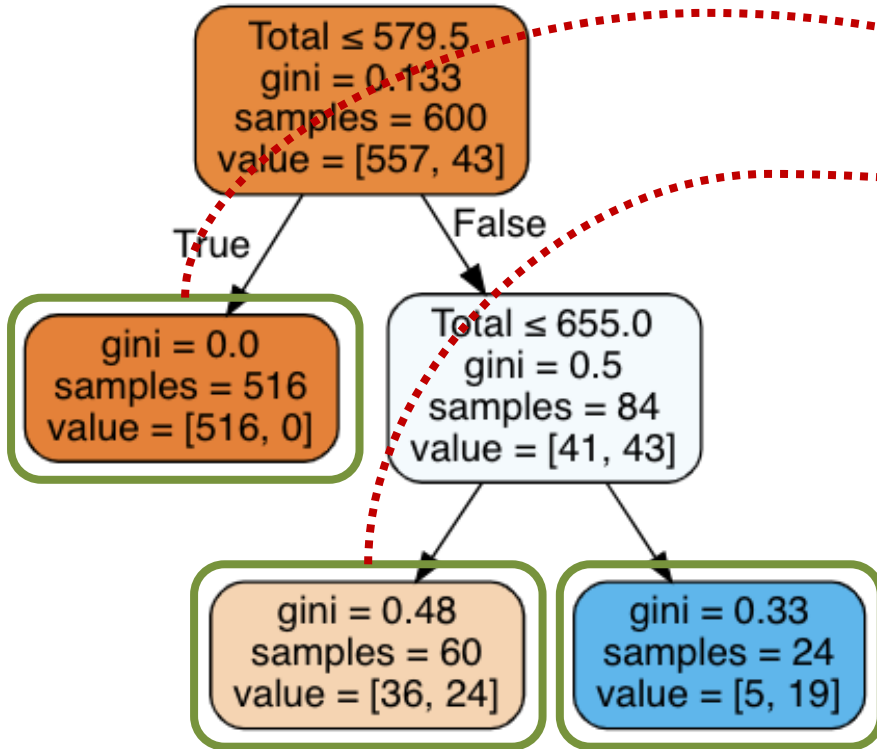
If **Total**  $\leq 655.0$ , go to Left Child  
Otherwise, go to Right Child

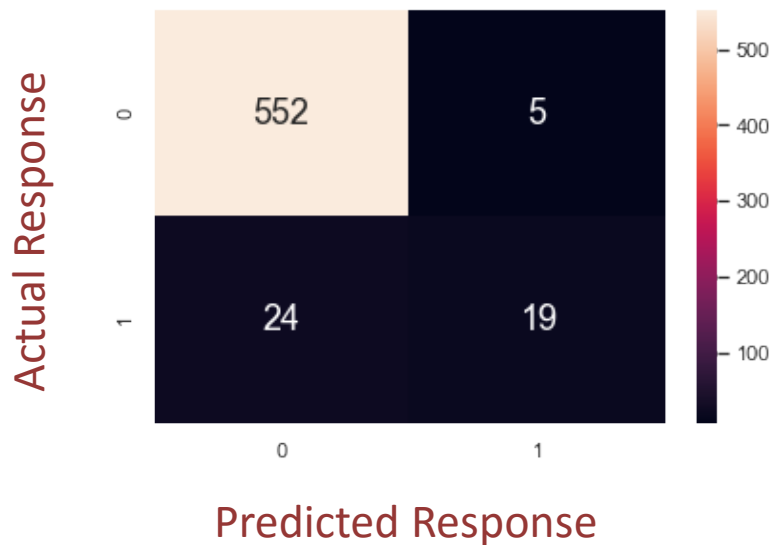
**Gini = 0.5** at this node

Total samples at node = **84**  
Proportion **NL : L = 41 : 43**

## Binary Classification

## Decision Tree to Confusion Matrix





Data Science

# Binary Classification

## Goodness of Fit of the Model

TP : True predicted as True	19
TN : False predicted as False	552
FN : True predicted as False	24
FP : False predicted as True	5

$$accuracy = \frac{552 + 19}{552 + 19 + 5 + 24}$$

“Positive” : 1

“Negative” : 0

Legendary  
Non-Legendary

$$tpr = \frac{19}{19 + 24}, \quad fnr = \frac{24}{24 + 19}$$

$$fpr = \frac{5}{5 + 552}, \quad tnr = \frac{552}{552 + 5}$$

Confusion Matrix : [https://en.wikipedia.org/wiki/Confusion\\_matrix](https://en.wikipedia.org/wiki/Confusion_matrix)



Actual N	TN	FP
Actual P	FN	TP
	Predicted N	Predicted P

When will you be happy?

**Ideal**      $TPR = 1, FPR = 0$

**Bad?**      $TPR = 1, FPR = 1$

**Bad?**      $TPR = 0, FPR = 0$

**Trash**      $TPR = 0, FPR = 1$

25	0
0	75

0	25
0	75

25	0
75	0

0	25
75	0

Balancing classes to achieve the desired TPR and FPR is a tricky thing to do. 😊

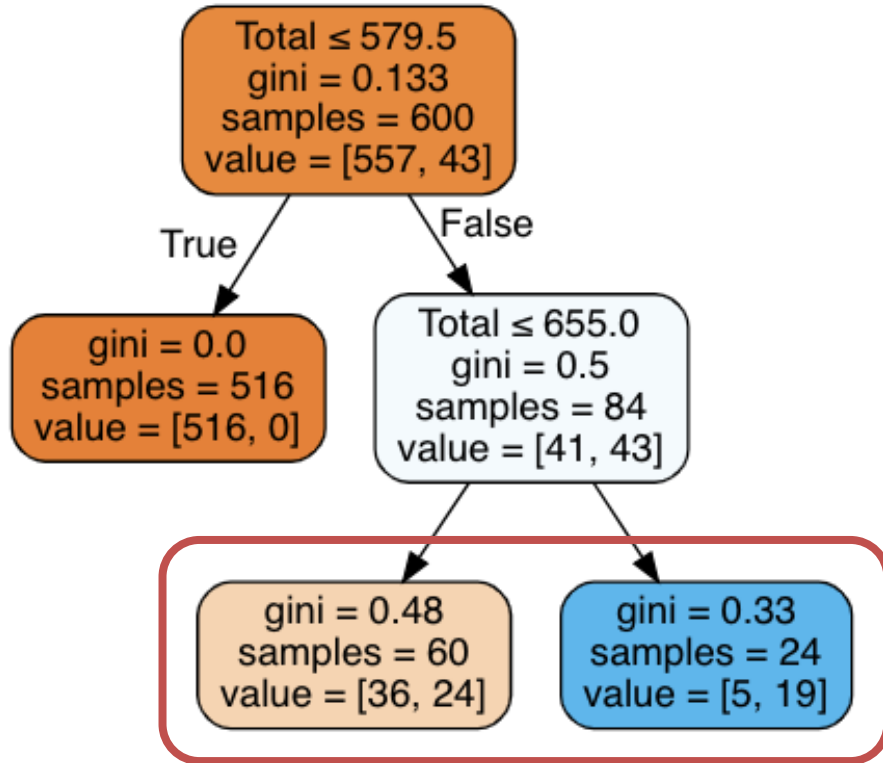
## Binary Classification

### How does a Tree “decide” classes?

The tree doesn't! **You** decide it on your own by choosing **Decision Threshold**.

If Proportion > T, you call it Positive, and else, you call it Negative class.

**Default Threshold for Trees = 0.5**



## Experiment with the Decision Threshold!

Use your tree to find Leaf Nodes.

Vary your decision threshold  $T$  in steps from 0 to 1 and note the TPR and FPR.

$T = 0$  : Everyone P  
TPR = 1, FPR = 1

$T = 1$  : Everyone N  
TPR = 0, FPR = 0

(0,1) Best

Bad (1,1)

