

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

## Decision Trees

Gini Impurity

Led by : Shreyas Shukla

# Gini Impurity

A mathematical measurement of how “pure” the information in a data set is.

We can think of this as a measurement of class uniformity.

Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

## Gini Impurity for Classification:

- For a set of classes **C** for a given dataset **Q**:

$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$

Led by : Shreyas Shukla

## Gini Impurity for Classification:

- For a set of classes  $\mathbf{C}$  for a given dataset  $\mathbf{Q}$ ,  $p_c$  is probability of class  $\mathbf{c}$ .

$$p_c = \frac{1}{N_Q} \sum_{x \in Q} \mathbb{1}(y_{class} = c) \quad G(Q) = \sum_{c \in C} p_c(1 - p_c)$$

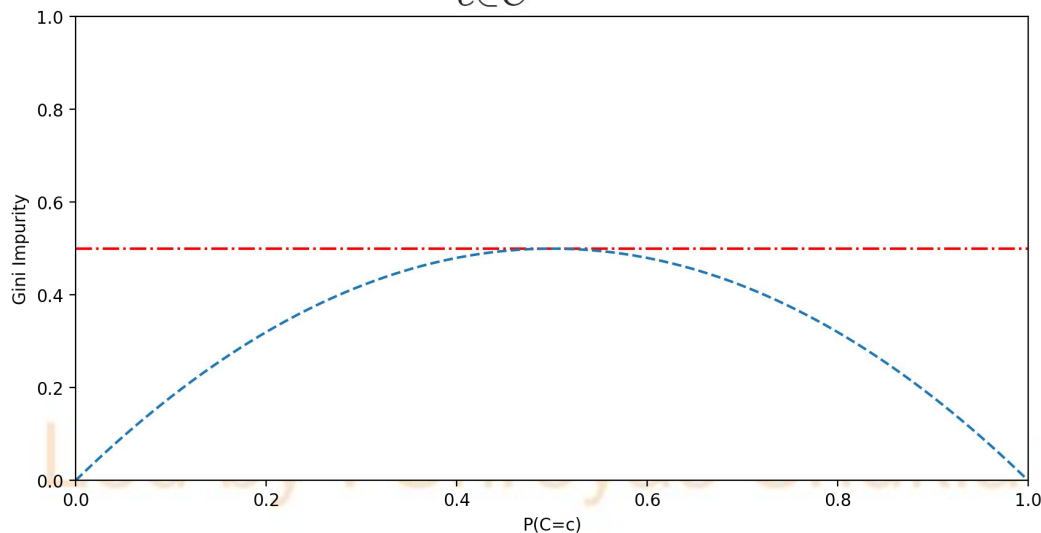
Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

## Gini Impurity for Classification:

$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$



# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

## Gini Impurity for Classification:

$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$



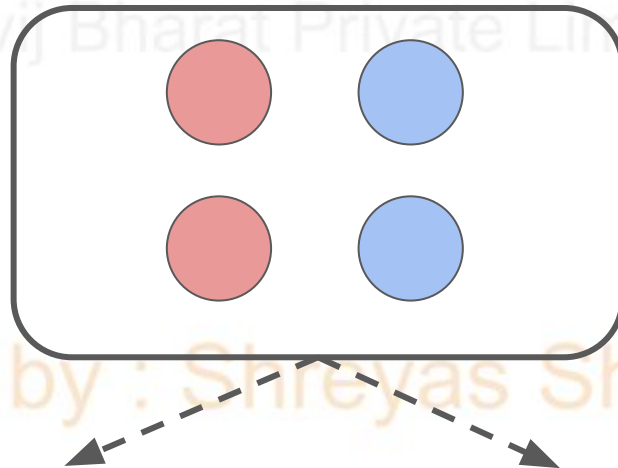
Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

## Gini Impurity for Classification:

$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$

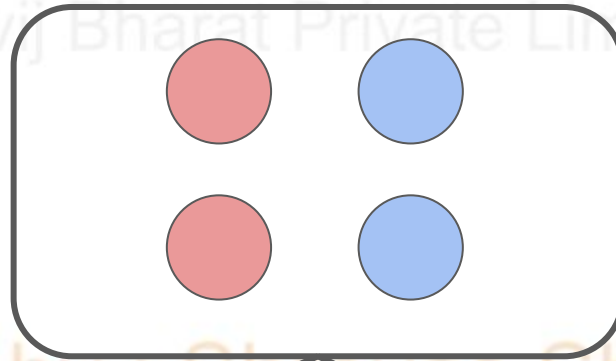


# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$

Class Red  
 $(2/4)(1 - 2/4) = 0.25$

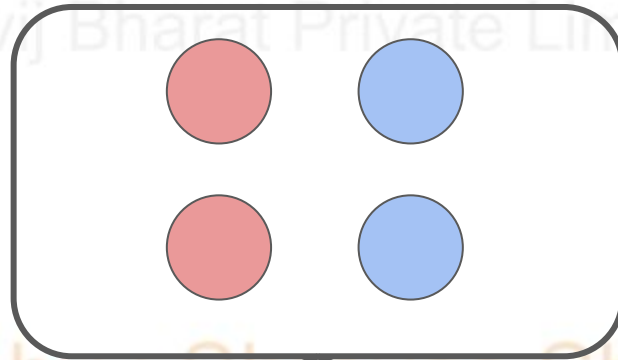




# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$



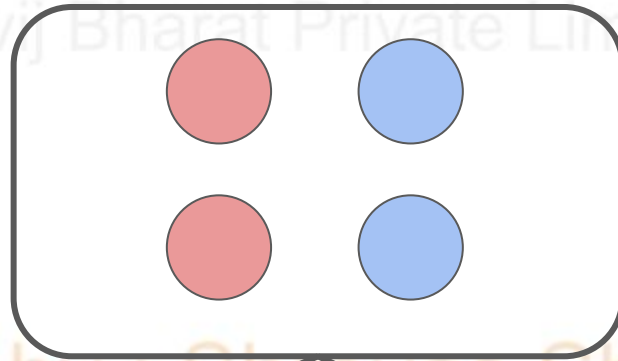
Class Red  
 $(2/4)(1 - 2/4) = 0.25$

Class Blue  
 $(2/4)(1 - 2/4) = 0.25$

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$



Class Red  
 $(2/4)(1 - 2/4) = 0.25$



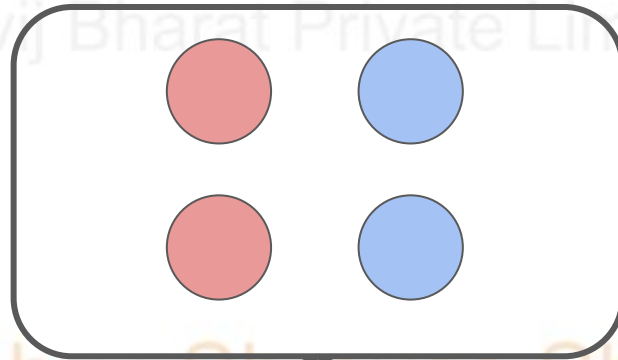
Class Blue  
 $(2/4)(1 - 2/4) = 0.25$



Gini Impurity  
 $0.25 + 0.25 = 0.5$

## “Maximum” Impurity Possible

$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$



Class Red  
 $(2/4)(1 - 2/4) = 0.25$



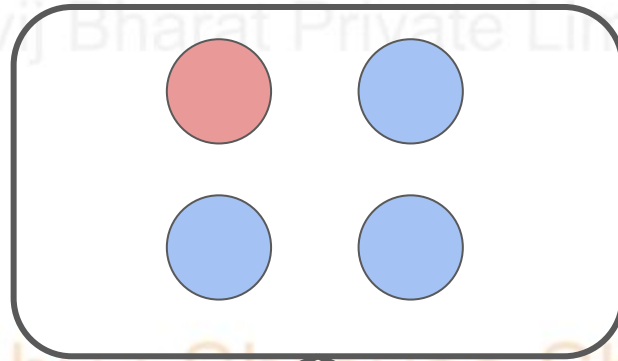
Class Blue  
 $(2/4)(1 - 2/4) = 0.25$



Gini Impurity  
 $0.25 + 0.25 = 0.5$

Data is more “pure” (less impurity)

$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$



Class Red  
 $(1/4)(1 - 1/4) = 0.1875$



Class Blue  
 $(3/4)(1 - 3/4) = 0.1875$



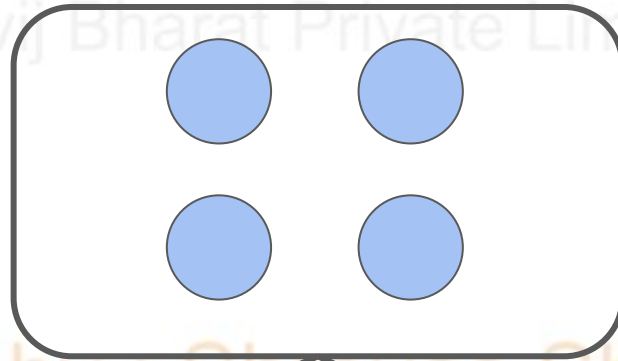
Gini Impurity  
 $0.1875 + 0.1875 = 0.375$

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Data is completely “pure” (no impurity)

$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$



Class Red  
 $(0/4)(1 - 0/4) = 0$



Class Blue  
 $(4/4)(1 - 4/4) = 0$



Gini Impurity  
 $0 + 0 = 0$

If the goal of a decision tree is to separate out classes, we can use gini impurity to decide on data split values.

We want to minimize the gini impurity at leaf nodes.

Minimized impurity at leaf nodes means we are separating classes effectively

Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

## Decision Trees

Gini Impurity in Trees

Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

For constructing a tree, we have to decide what feature will be root node.

Use gini impurity to compare the information contained within features for the training data.

Led by : Shreyas Shukla



## Gini Impurity for Classification:

- For a set of classes  $\mathbf{C}$  for a given dataset  $\mathbf{Q}$ ,  $p_c$  is probability of class  $\mathbf{c}$ .

$$p_c = \frac{1}{N_Q} \sum_{x \in Q} \mathbb{1}(y_{class} = c)$$

$$G(Q) = \sum_{c \in C} p_c (1 - p_c)$$

Led by : Shreyas Shukla

# Mastering Machine Learning with Python

Create a decision tree to predict spam.

X - URL Link	Y-Spam
Yes	Yes
Yes	Yes
No	No
No	No
No	Yes
No	No
Yes	No

Led by : Shreyas Shukla

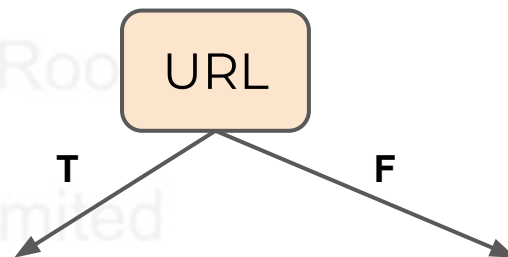
# Only one X feature to use for a node.

X - URL Link	Y-Spam
Yes	Yes
Yes	Yes
No	No
No	No
No	Yes
No	No
Yes	No

URL

Predict if email is spam if it contains a URL:

X - URL Link	Y-Spam
Yes	Yes
Yes	Yes
No	No
No	No
No	Yes
No	No
Yes	No

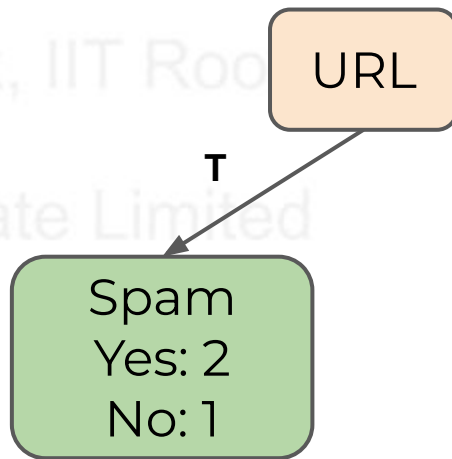


Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

X - URL Link	Y-Spam
Yes	Yes
Yes	Yes
No	No
No	No
No	Yes
No	No
Yes	No

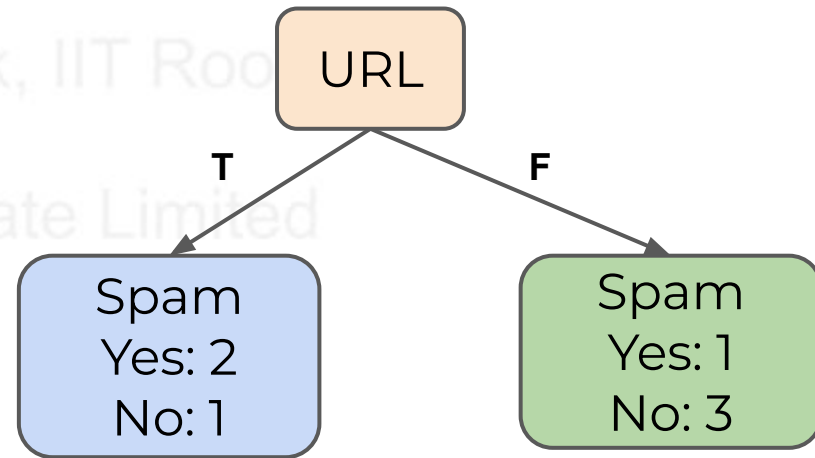


Preyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

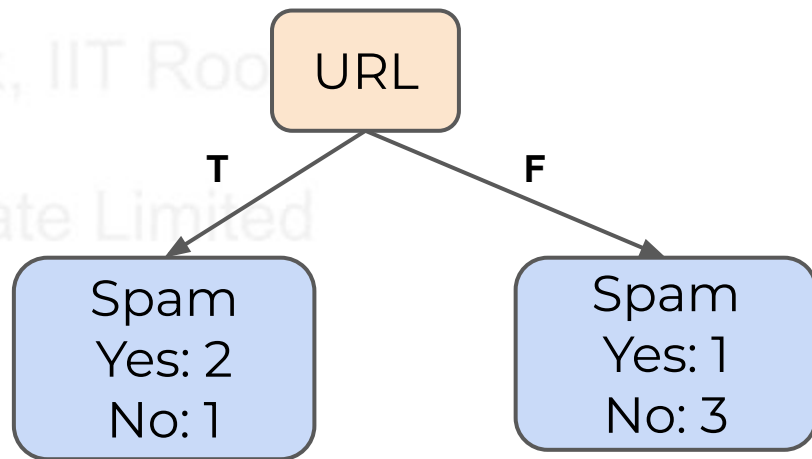
X - URL Link	Y-Spam
Yes	Yes
Yes	Yes
No	No
No	No
No	Yes
No	No
Yes	No



Led by : Shreyas Shukla

Predict if email is spam if it contains a URL:

X - URL Link	Y-Spam
Yes	Yes
Yes	Yes
No	No
No	No
No	Yes
No	No
Yes	No

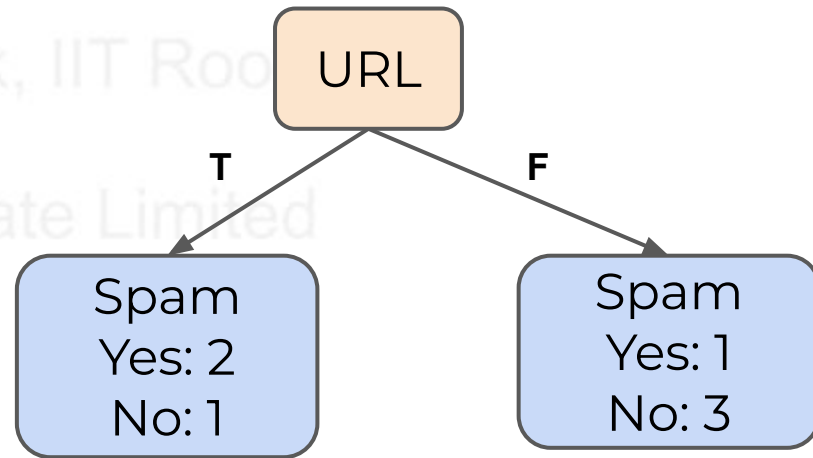


Led by : Shreyas Shukla

# Recall the gini impurity formula:

(27th Aug 2024 - 18th Oct 2024)

X - URL Link	Y-Spam
Yes	Yes
Yes	Yes
No	No
No	No
No	Yes
No	No
Yes	No



$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$



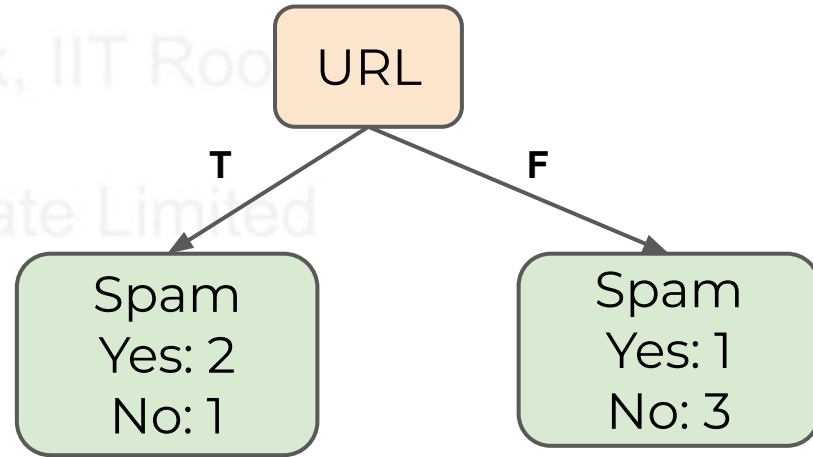
# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Treat Yes Spam and No Spam as **c** classes:

Left Leaf Node:

$$\left(\frac{2}{3}\right)\left(1-\frac{2}{3}\right) + \left(\frac{1}{3}\right)\left(1-\frac{1}{3}\right)$$



Led by : Shreya

$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$

Treat Yes Spam and No Spam as **c** classes:

Left Leaf Node:

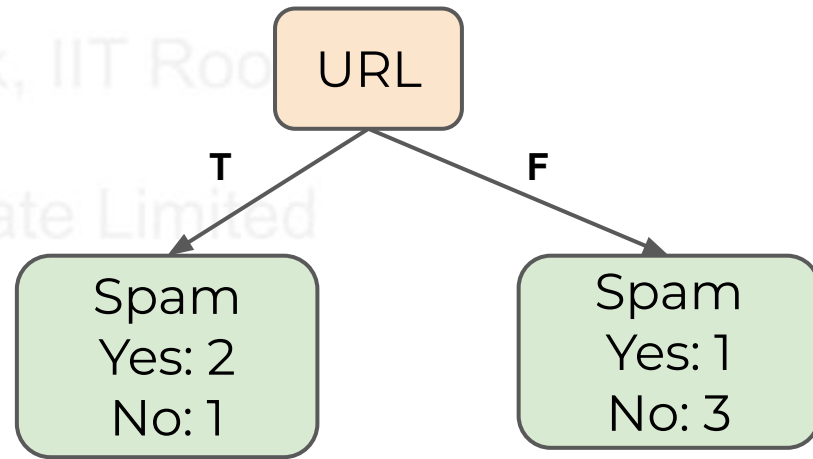
$$\left(\frac{2}{3}\right)\left(1-\frac{2}{3}\right) + \left(\frac{1}{3}\right)\left(1-\frac{1}{3}\right)$$

Left Leaf Gini=0.44

Right Leaf Node:

$$\left(\frac{1}{4}\right)\left(1-\frac{1}{4}\right) + \left(\frac{3}{4}\right)\left(1-\frac{3}{4}\right)$$

Right Leaf Gini=0.375



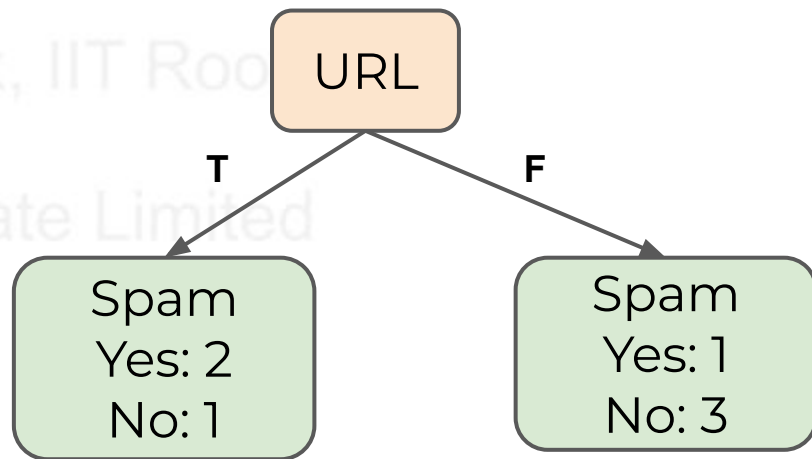
$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$

Calculate gini impurity of URL feature.

Weighted Average of both:

Left Leaf Gini=0.44

Right Leaf Gini=0.375



Led by : Shreya

$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

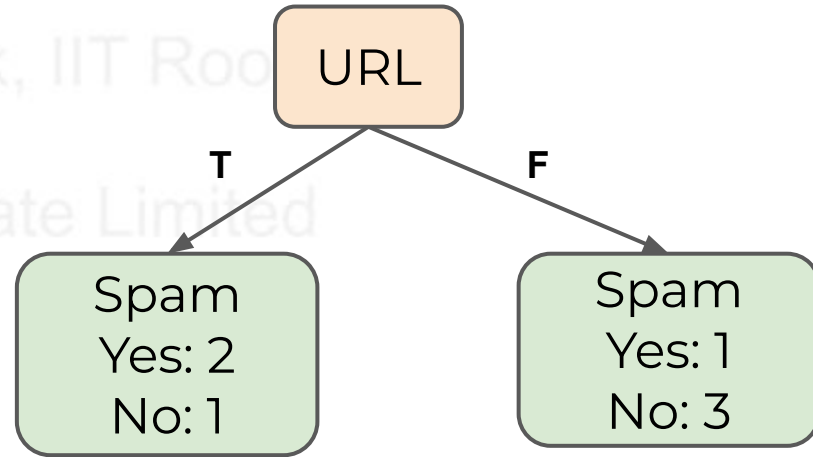
Total Emails:  $(2+1) + (1+3) = 7$

Left Leaf Gini=0.44

Right Leaf Gini=0.375

Left Emails: 3

Right Emails: 4



Led by : Shreya

$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Total Emails:  $(2+1) + (1+3) = 7$

Left Leaf Gini=0.44

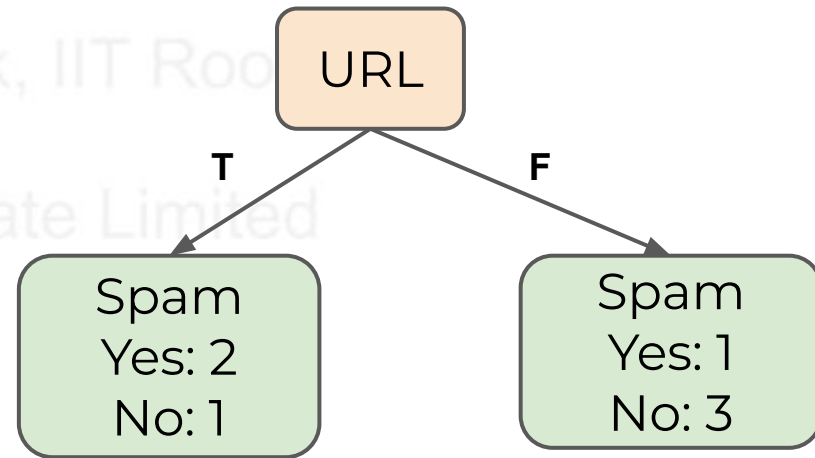
Right Leaf Gini=0.375

Left Emails: 3

Right Emails: 4

$(3/7)*0.44 + (4/7)*0.375$

Gini Impurity: 0.403



Led by : Shreya

$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

More issues to consider:

- Multiple Features
- Continuous Features
- Multi-categorical Features

We use the gini impurity to each of these issues to solve for best root nodes and best split parameters for leaves.

Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

## Decision Trees

### Gini Impurity Part Two

Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Let's explore:

- Continuous numeric features
- Multi-categorical features ( $N > 2$ )
- Choosing a root node feature

Led by : Shreyas Shukla



# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Imagine a continuous feature.

Calculate the feature gini impurity:

X - Words in Email	Y-Spam
10	Yes
40	No
20	Yes
50	No
30	No

Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Sort data:

X - Words in Email	Y-Spam
10	Yes
40	No
20	Yes
50	No
30	No

Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

## Calculate potential split values for node

Words  $\leq N$

X - Words in Email	Y-Spam
10	Yes
20	Yes
30	No
40	No
50	No

Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Use averages between rows as values:

Words  $\leq N$

X - Words in Email		Y-Spam
15	10	Yes
25	20	Yes
	30	No
35	40	No
45	50	No

Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Perform all the potential split:

Words  $\leq 15$

X - Words in Email		Y-Spam
15	10	Yes
25	20	Yes
	30	No
35	40	No
45	50	No

Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Calculate gini impurity for each split:

Words  $\leq 15$

X - Words in Email		Y-Spam
15	10	Yes
	20	Yes
	30	No
	40	No
	50	No

Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Calculate gini impurity for each split:

Words  $\leq 15$

X - Words in Email	Y-Spam
10	Yes
20	Yes
30	No
40	No
50	No

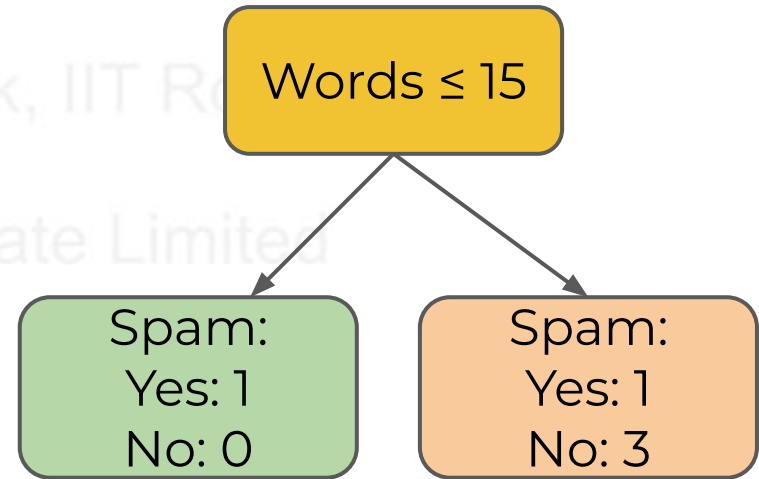
Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Calculate gini impurity for each split:

X - Words in Email	Y-Spam
10	Yes
20	Yes
30	No
40	No
50	No



$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$

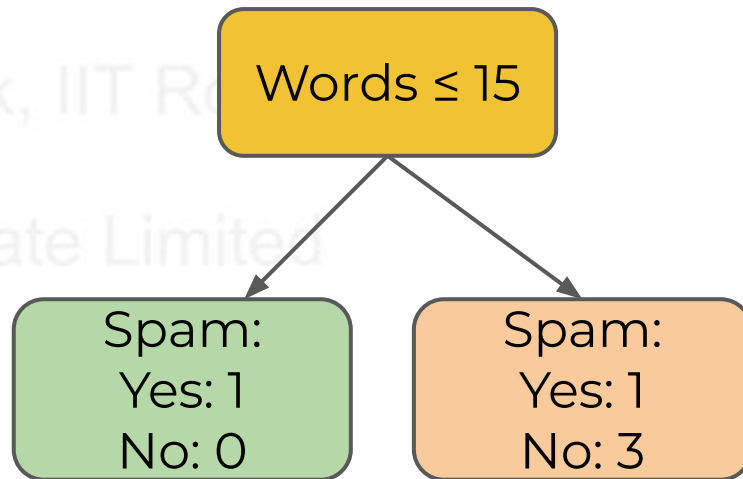


# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Calculate gini impurity for each split:

X - Words in Email	Y-Spam
10	Yes
20	Yes
30	No
40	No
50	No



$$\begin{aligned} G(Q) &= \left(\frac{1}{5}\right)(0+0) + \left(\frac{4}{5}\right)\left(\left(\frac{1}{4}\right)\left(1-\frac{1}{4}\right) + \left(\frac{3}{4}\right)\left(1-\frac{3}{4}\right)\right) \\ &= 0.3 \end{aligned}$$

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Do it for all possible splits:

X - Words in Email		Y-Spam	
15	10	Yes	→ Gini=0.3
25	20	Yes	
35	30	No	→ Gini=0
45	40	No	→ Gini=0.26
	50	No	→ Gini=0.4

Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Choose lowest impurity split value

X - Words in Email	Y-Spam
10	Yes
20	Yes
25	No
30	No
40	No
50	No

Gini=0

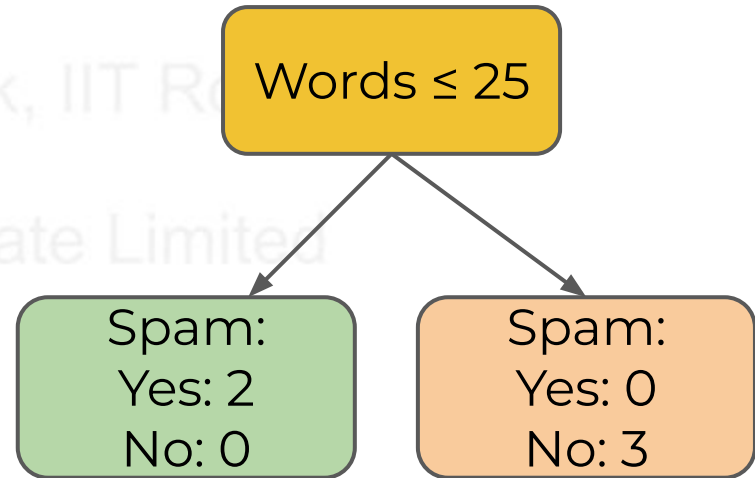
Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Choose this as split value for node

X - Words in Email	Y-Spam
10	Yes
20	Yes
25	No
30	No
40	No
50	No



$$G(Q) = 0$$

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

## Multicategorical feature

Calculate gini impurity for all combinations:

X - Sender	Y-Spam
Abe	Yes
Bob	Yes
Claire	No
Abe	No
Bob	No

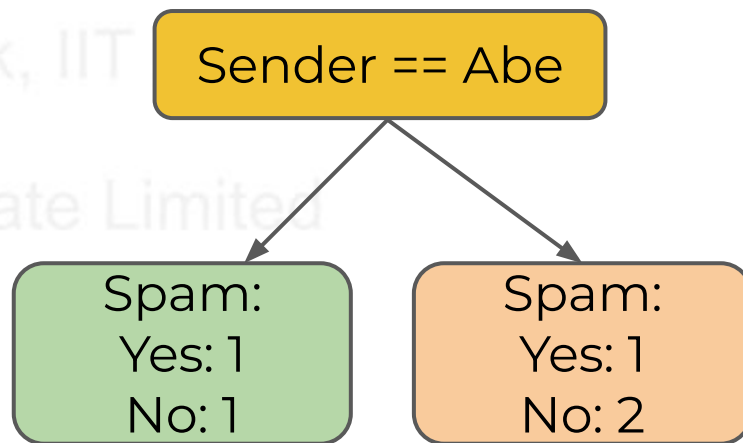
Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

Calculate gini impurity for all combinations:

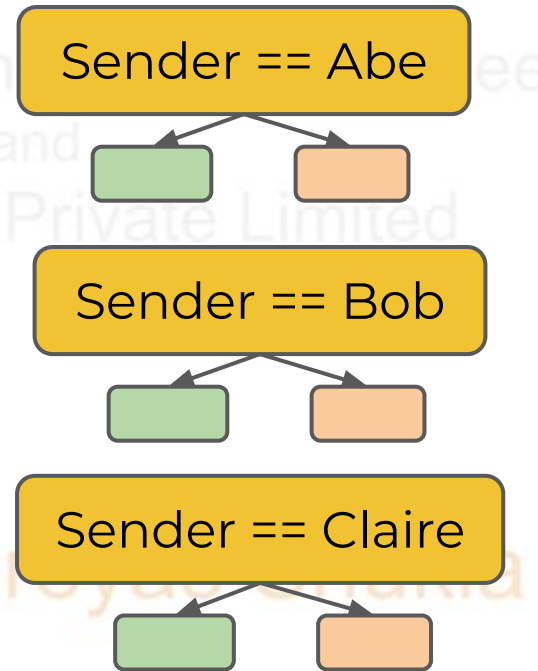
X - Sender	Y-Spam
Abe	Yes
Bob	Yes
Claire	No
Abe	No
Bob	No



Led by : Shreyas Shukla

## Calculate gini impurity for all combinations

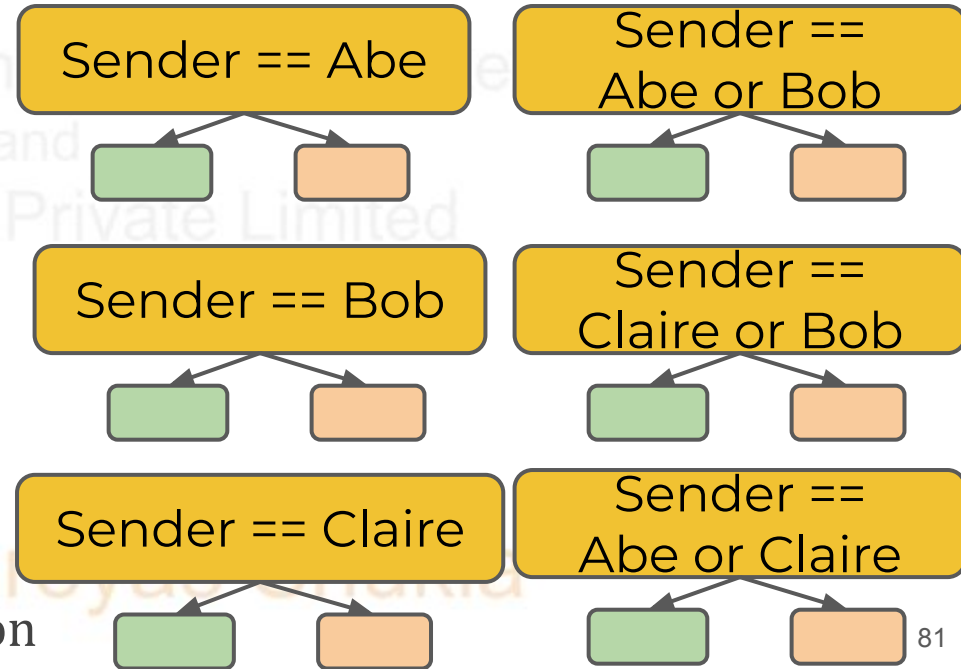
X - Sender	Y-Spam
Abe	Yes
Bob	Yes
Claire	No
Abe	No
Bob	No



# Calculate gini impurity for all combinations

(27th Aug 2024 - 18th Oct 2024)

X - Sender	Y-Spam
Abe	Yes
Bob	Yes
Claire	No
Abe	No
Bob	No



Choose lowest impurity split combination



Now we can split any type of feature.

(27th Aug 2024 - 18th Oct 2024)

## **How does the decision tree decide on the root node of a multi-feature dataset?**

Calculate the gini impurity values of each feature and choose the lowest impurity value to split on first.

Ritvij Bharat Private Limited

By choosing the feature with the lowest resulting gini impurity in its leaf nodes, we are choosing the feature that best splits the data into “pure” classes.

Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

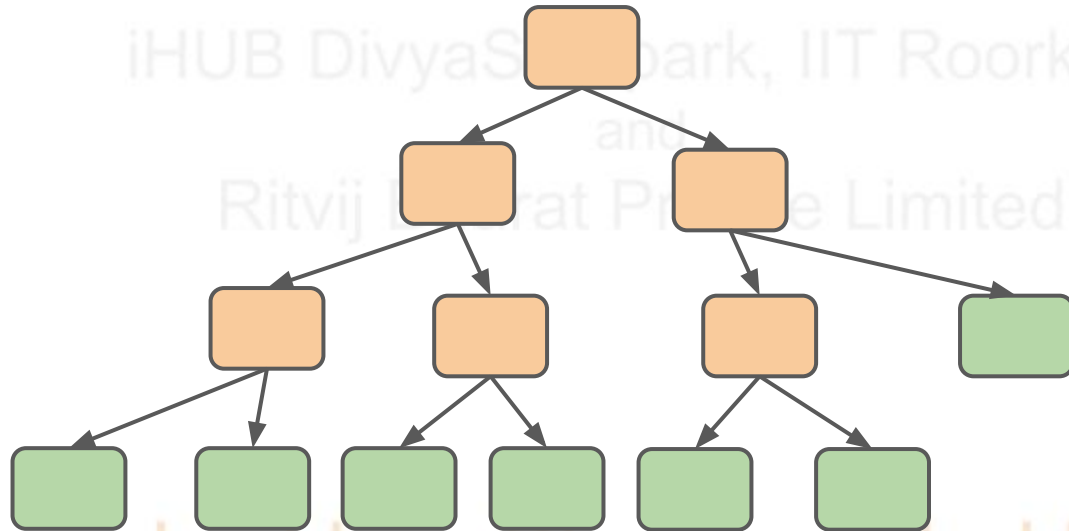
By using gini impurity as a measurement of the effectiveness of a node split, we can perform automatic feature selection by mandating an impurity threshold for an additional feature based split to occur.

Led by : Shreyas Shukla

# A large overfitted tree.

(27th Aug 2024 - 18th Oct 2024)

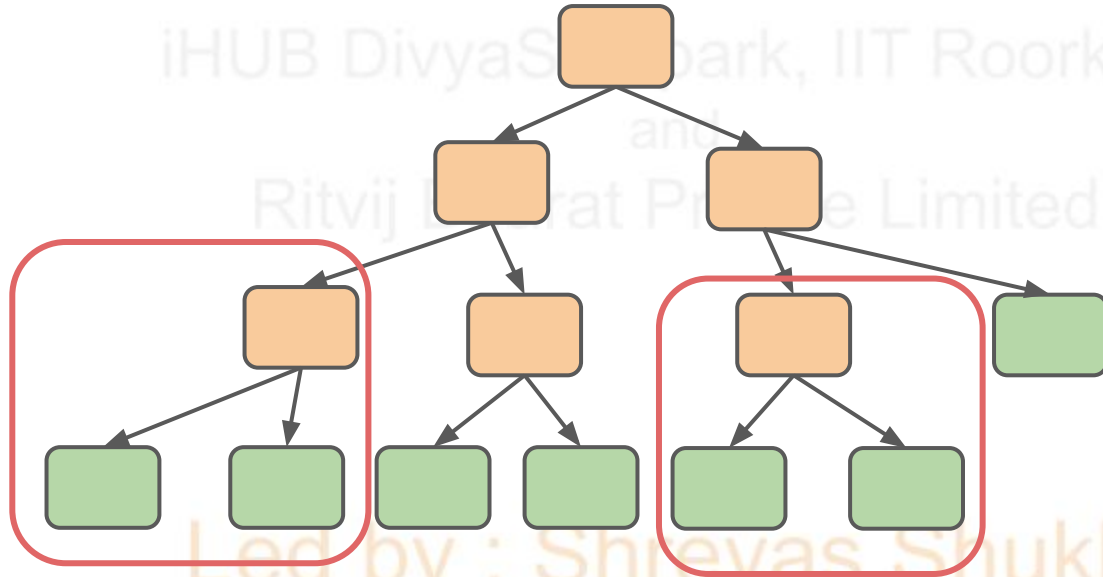
## Add minimum gini impurity decrease



Led by : Shreyas Shukla

# Mastering Machine Learning with Python

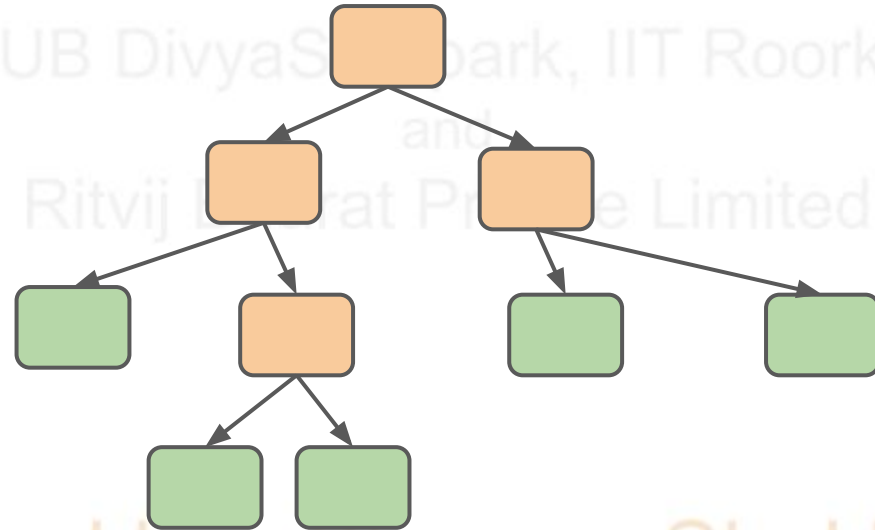
(27th Aug 2024 - 18th Oct 2024)



Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

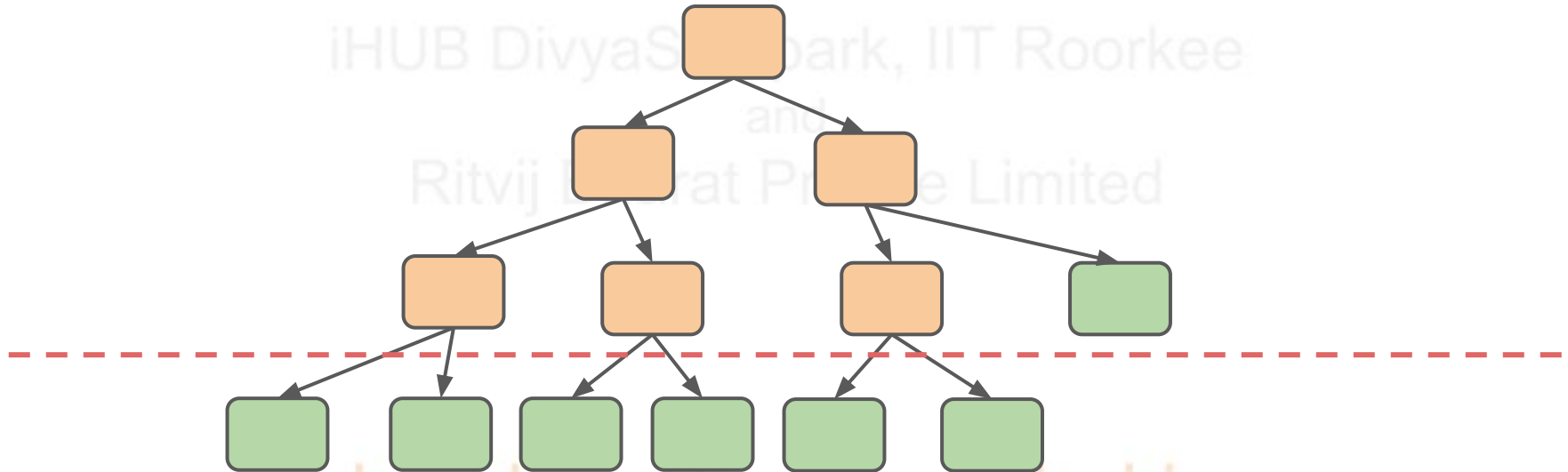


Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

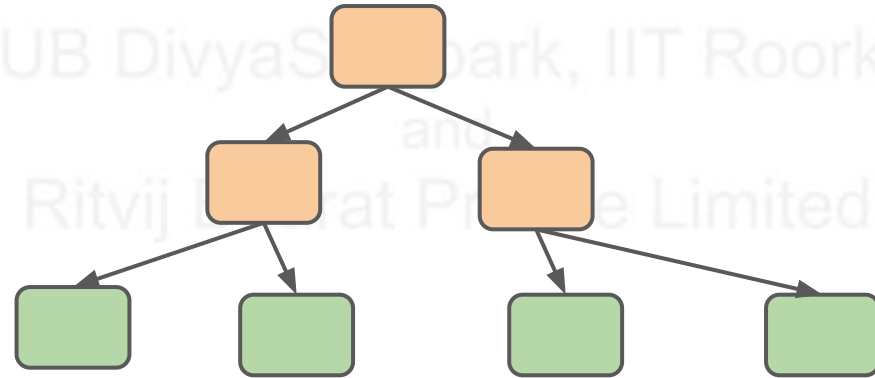
We can also mandate a max depth



Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)



Led by : Shreyas Shukla

# Mastering Machine Learning with Python

(27th Aug 2024 - 18th Oct 2024)

iHUB DivyaSampark, IIT Roorkee  
and  
Ritvij Bharat Private Limited

**Let's code !!**

**Led by : Shreyas Shukla**