

Decision Tree

①

Introduction

- Non-parametric \rightarrow Non-linear data
- white box model \rightarrow Inference
- Mother of all tree based algo
- Work with both classification and Regression.

Intuition behind DT

Gender	Occupation	Suggestion
F	Student	PUBG
F	Programmer	Github
M	Programmer	Whatsapp
F	Programmer	Github
M	Student	PUBG
M	Student	PUBG

Giant Nested if - else structure

ex:-

```
if Occupation = student  
    print ("PUBG")
```

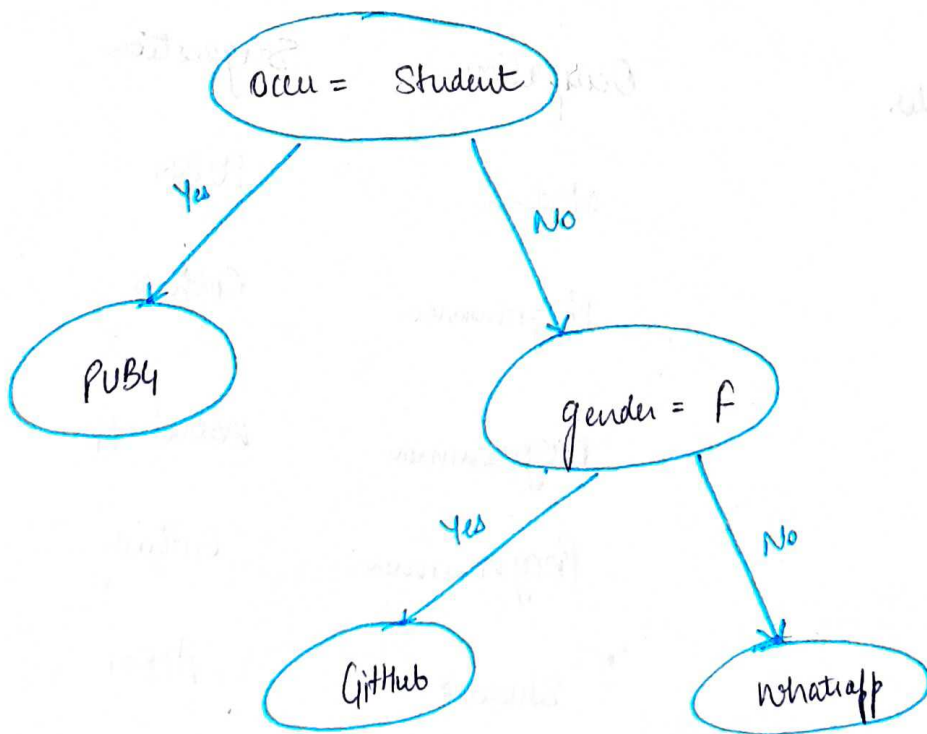
else :

```
    if gender = f  
        print ("Github")
```

else :

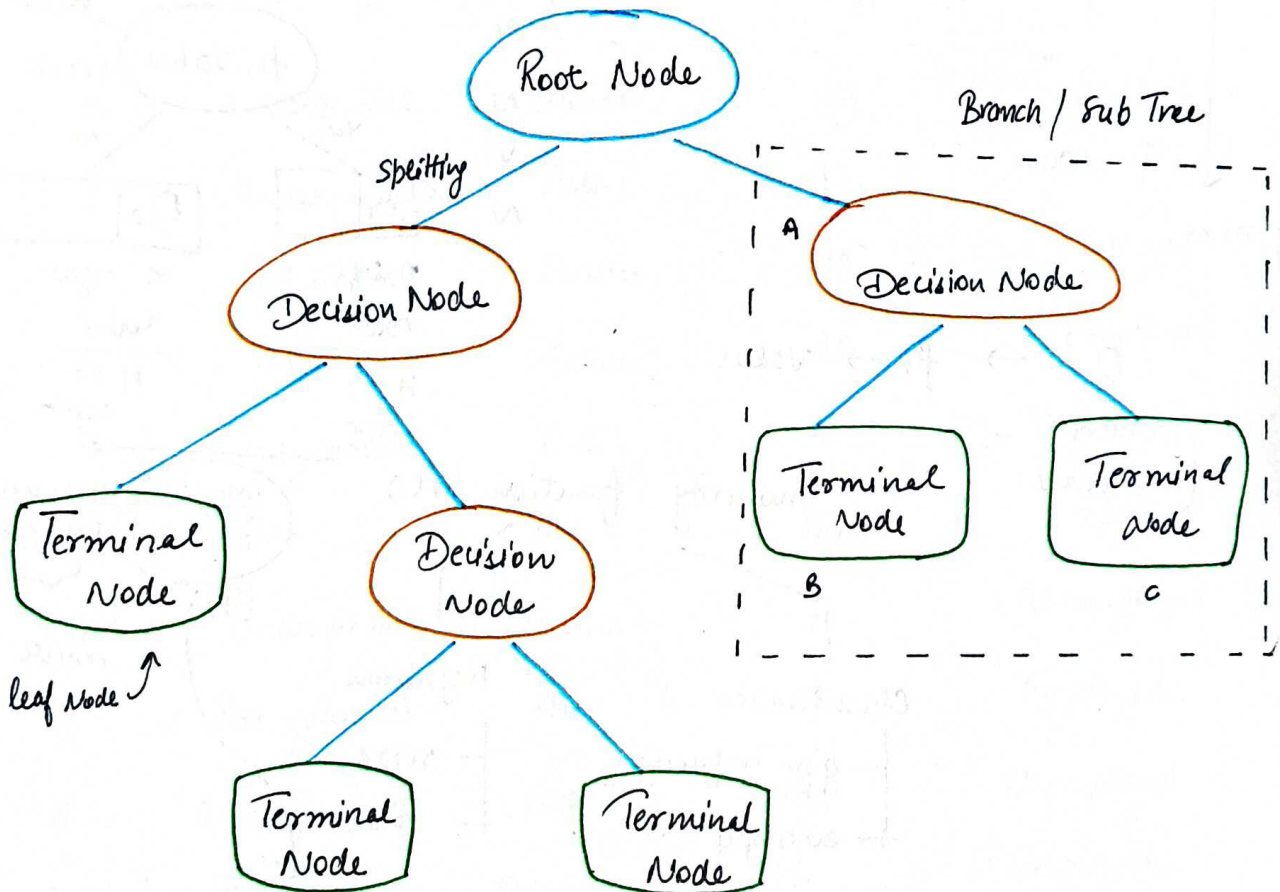
```
        print ("whatsapp")
```

basic
decision
Tree



Vocab

②



Note:- A is parent node of B and C

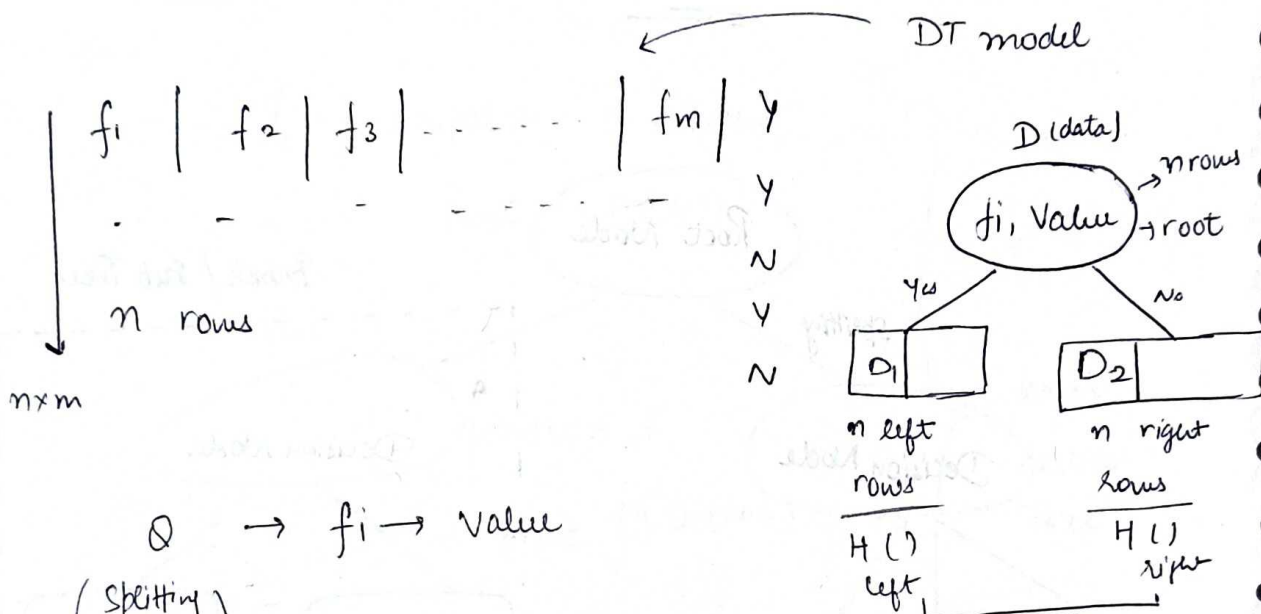
CART → Classification and Regression Tree

- ID3
- C4.5
- CART

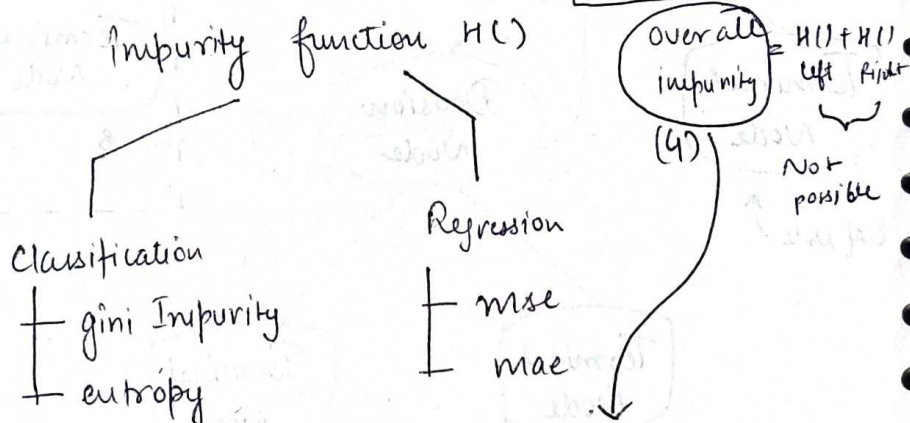
Algorithm to grow decision Tree

Sklearn use internally

Classification



$\emptyset \rightarrow f_i \rightarrow \text{value}$
(splitting criteria)



Overall impurity $(G) = \frac{n_{\text{left}} H(\text{left})}{n} + \frac{n_{\text{right}} H(\text{right})}{n}$

$f_1, \text{value} \rightarrow G()$
 $f_2, \text{value} \rightarrow G()$
 \vdots
 \vdots

}

Pick those combination which have minimum $G()$ value

impurity value

Splitting Categorical Features

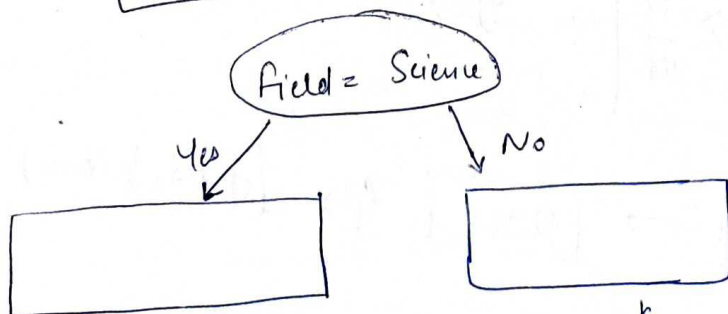
(3)

	$f_1 \downarrow$ Degree-type	$f_2 \downarrow$ Field	$f_3 \downarrow$ Average-Grade	Job-Outcome
0	Undergraduate	Science	89	Employed
1	Undergraduate	Arts	92	Unemployed
2	Post graduate	Science	95	Employed
3	PhD	Science	85	Employed
4	Postgraduate	Arts	98	Unemployed
5	PhD	Arts	90	Employed
6	Undergraduate	Science	88	Unemployed
7	Post graduate	Arts	93	Employed
8	Undergraduate	Arts	94	Unemployed
9	PhD	Science	86	Employed

multi
class \uparrow

binary
class \uparrow

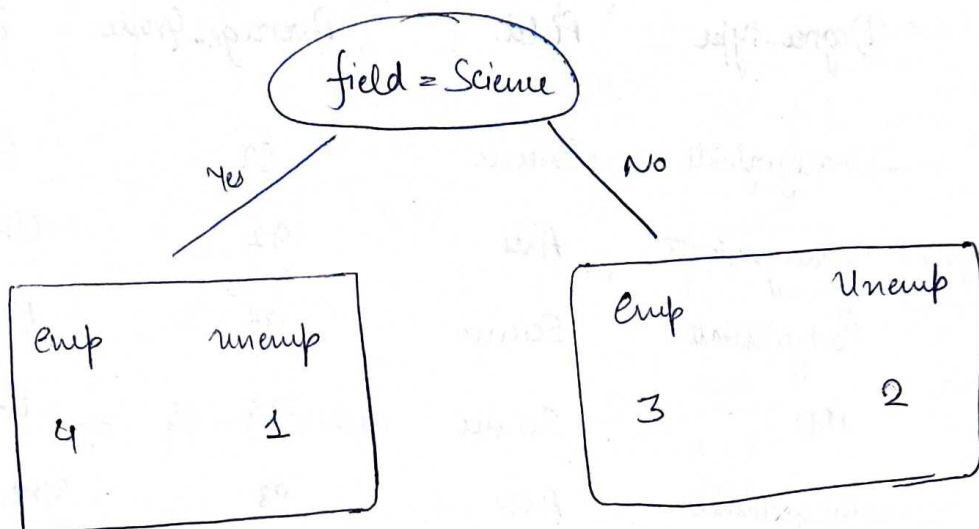
Numerical \uparrow



gini impurity $= 1 - \sum_{k=1}^k p_k^2$ \rightarrow $p_k =$ Proportion of each classes

$p_{\text{science}} = \left(\frac{5}{10}\right)$ $p_{\text{Arts}} = \left(\frac{5}{10}\right)$

$$= 1 - \left(\frac{5}{10}\right)^2 - \left(\frac{5}{10}\right)^2$$



$$P_{\text{emp}} = \frac{4}{5}$$

$$P_{\text{unemp}} = \frac{1}{5}$$

↓
gini

$$1 - \left(\frac{4}{5}\right)^2 - \left(\frac{1}{5}\right)^2$$

$$1 - \left(\frac{16}{25}\right) - \left(\frac{1}{25}\right) = \frac{8}{25}$$

$$P_{\text{emp}} = \frac{3}{5}$$

$$P_{\text{unemp}} = \frac{2}{5}$$

↓
gini

$$1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2$$

$$1 - \frac{9}{25} - \frac{4}{25} = \frac{12}{25}$$

$$\left[\frac{5}{10} \times \frac{8}{25} + \frac{5}{10} \times \frac{12}{25} \right]$$

g()

$$0.16 + 0.24 \rightarrow \boxed{0.40} \quad q \rightarrow (\text{field, Science})$$

multiclass

Value \rightarrow UG / PG / PhD

(4)

degree, UG

G1

degree, PG

G1

degree, PhD

G1

min

degree = UG

Yes

No

emp	Unemp
1	3

emp	Unemp
5	1

$$1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2$$

$$1 - \left(\frac{5}{6}\right)^2 - \left(\frac{1}{6}\right)^2$$

0.3

0.4

$$\frac{4}{10} \times 0.3 + \frac{6}{10} \times 0.4 = 0.5$$

degree = PG

Yes

No

emp	Unemp
2	1

emp	Unemp
4	3

$$1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 0.25$$

$$1 - \left(\frac{4}{7}\right)^2 - \left(\frac{3}{7}\right)^2 = 0.35$$

$$\frac{3}{10} \times 0.25 + \frac{7}{10} \times 0.35 = 0.65$$

degree = UG

<

degree = PG

↓

repeat Criteria

because greater value

Numerical column

↓ Sort

	Avg-Grade	job outcome
	85	emp
85.5	86	emp
87	88	unemp
88.5	89	emp
89.5	90	emp
91	92	unemp
92.5	93	emp
93.5	94	unemp
94.5	95	emp
	98	unemp

Grade > 85.5

yes

no

emp	unemp
5	4

emp	unemp
1	0

pure node
Unemp have 0

$$1 - \left(\frac{5}{9}\right)^2 - \left(\frac{4}{9}\right)^2$$

$$1 - \left(\frac{1}{1}\right)^2 - \left(\frac{0}{1}\right)^2 \quad (3)$$

$$1 - \frac{25}{81} - \frac{16}{81}$$

$$1 - 1 = 0$$

$$1 - \frac{41}{81} = \frac{40}{81}$$

$$G(1) = \frac{9}{10} \times \frac{40}{81} + \frac{1}{10} \times 0 = 0.25$$

Let assume,

0.4 → 85.5

0.3 → 87

0.2 → 88.5

0.1 → 89.5

0.8 → 91

0.6 → 92.5

0.7 → 93.5

94.5

grade > 89.5 → 0.1

this value compare with other value

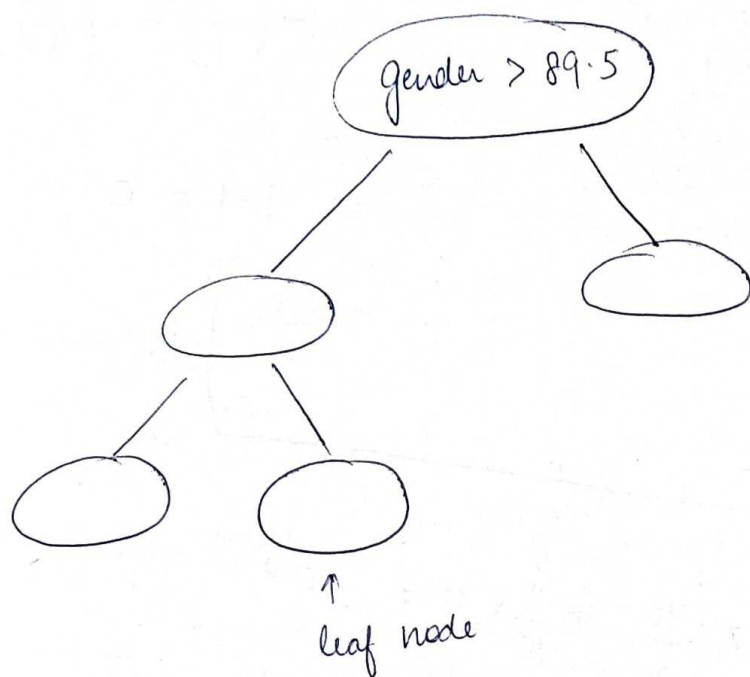
0.1 → Small
↑

Degree-Type

0.4
↑
Field

Average-Grade

Job-Occupation

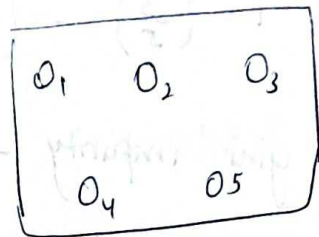
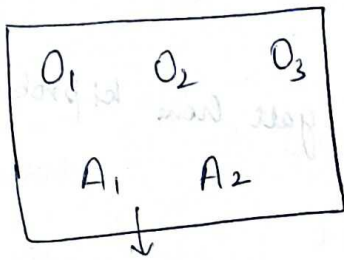


Understanding Gini Impurity?

The Gini impurity is a measure of how often a randomly chosen element from the set would be incorrectly labelled if it was randomly labelled according to the distribution of labels in the subset.

$$H(i) = 1 - \sum_{k=1}^k p_k^2$$

example:



3 Oranges

2 Apples

	O_1	O_2	O_3	A_1	A_2
O_1	✓				
O_2		✓			
O_3			✓		
A_1				✓	
A_2					✓

	O_1	O_2	O_3	A	A
O	✓	✓	✓	✗	✗
O	✓	✓	✓	✗	✗
O	✓	✓	✓	✗	✗
A	✗	✗	✗	✓	✓
A	✗	✗	✗	✓	✓

$$1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2$$

gini impurity \rightarrow guess galt hone ki prob.

Geometric Intuition of Decision Tree

Cgpa	iq	Placement
-	-	Yes
-	-	No
-	-	Yes

