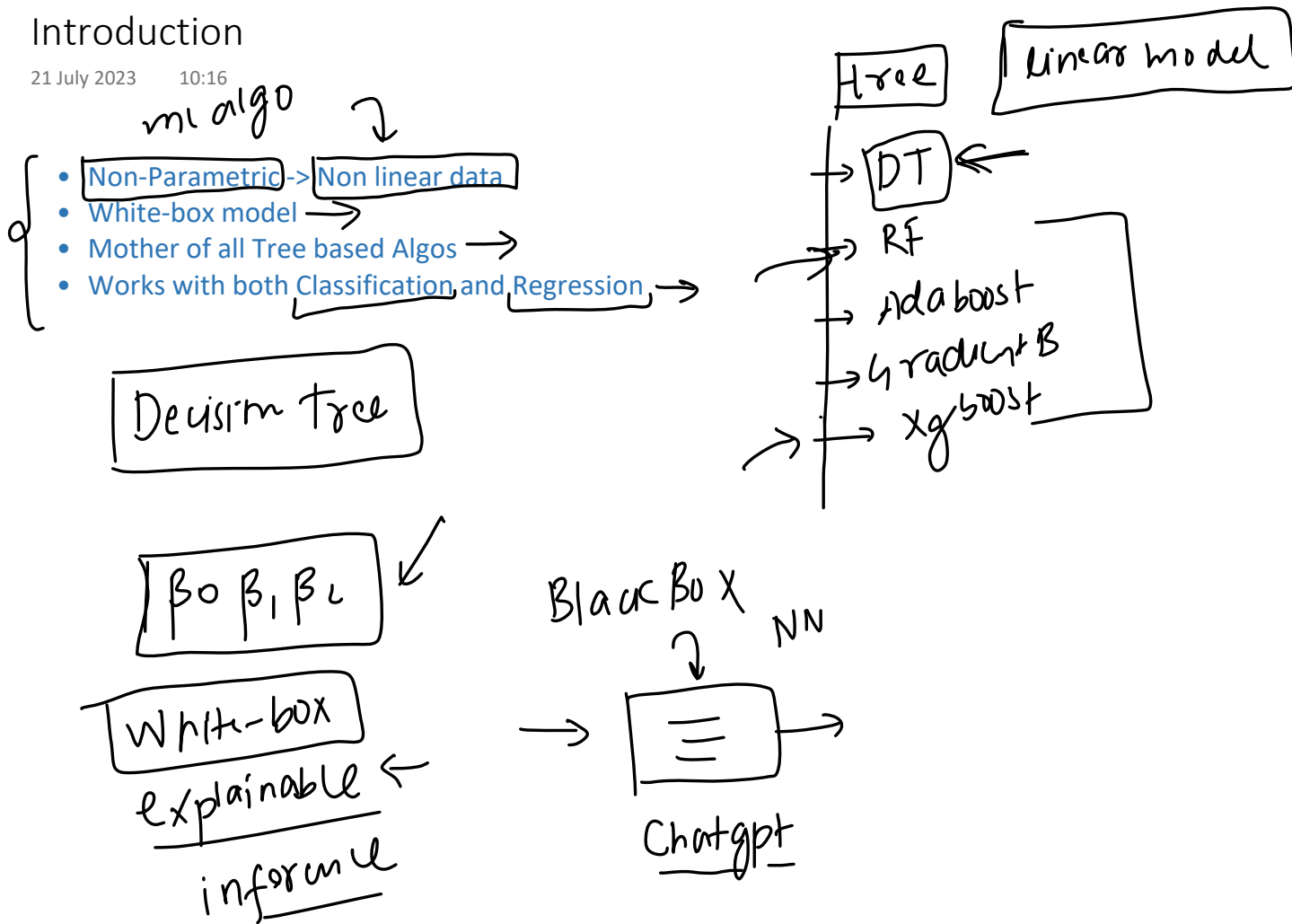


Introduction

21 July 2023 10:16



Intuition behind DT

20 July 2023 16:21

Gender	Occupation	Suggestion
F	<u>Student</u>	PUBG
<u>F</u>	Programmer	Github
<u>M</u>	Programmer	<u>Whatsapp</u>
<u>F</u>	Programmer	Github
M	<u>Student</u>	PUBG
M	<u>Student</u>	PUBG

$\{m, prog\} \rightarrow \{F, \underline{student}\} \leftarrow$

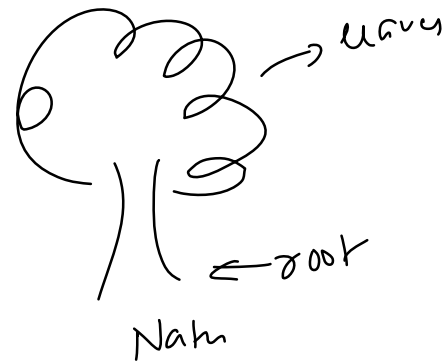
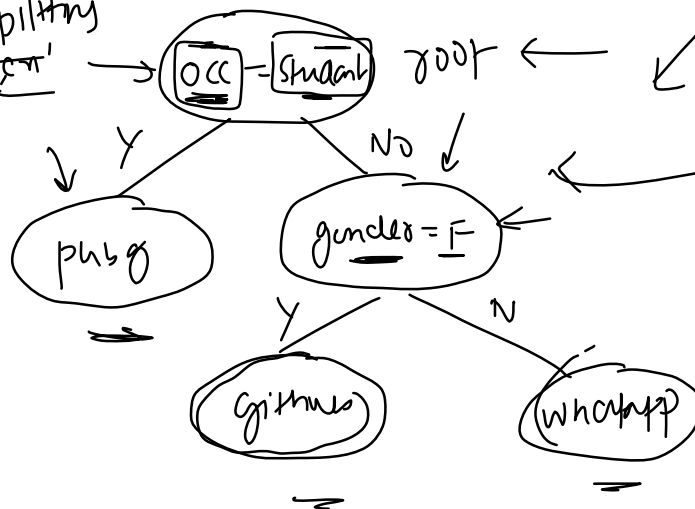
Nested
giant if-else structure

```

if occupation = student
  print pubg
else
  if gender = F
    print github
  else
    print whatsapp
  
```

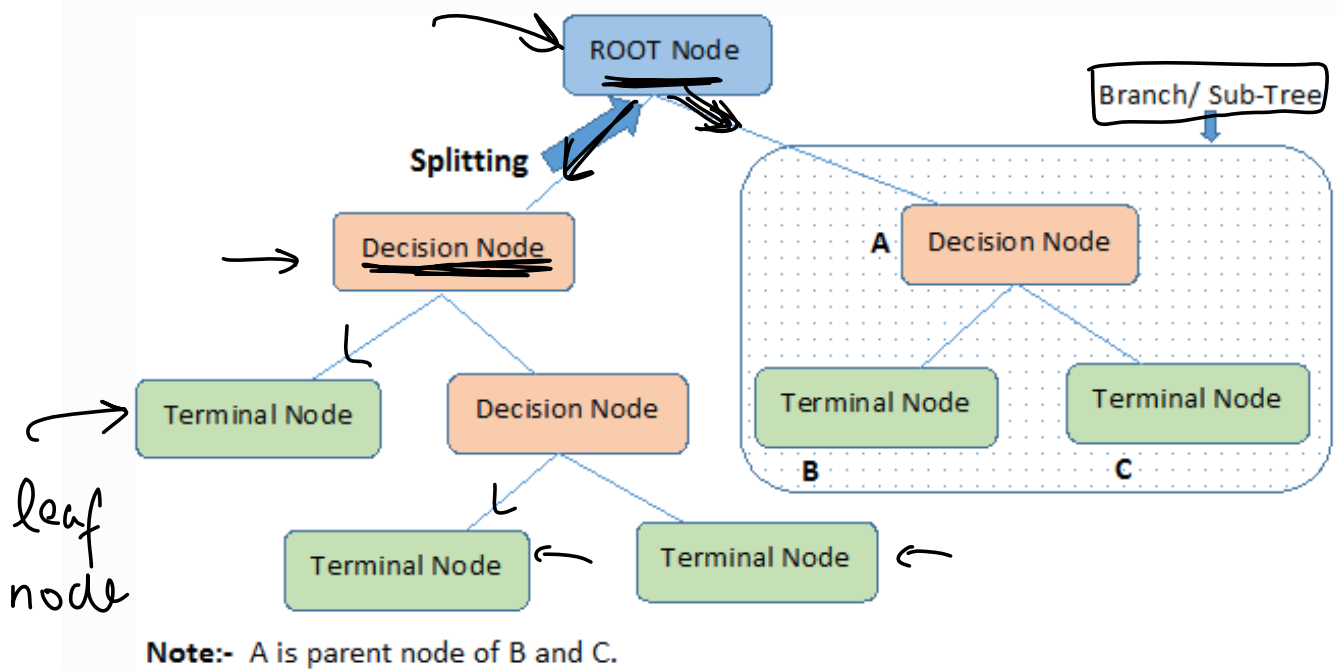
basic
dt

splitting
on 'occ'



Vocab

20 July 2023 16:22



Example 2

20 July 2023 16:25

	Degree_Type	Field	Average_Grade	Job_Outcome
0	Undergraduate	Science	89	Employed
1	Undergraduate	Arts	92	Unemployed
2	Postgraduate	Science	95	Employed
3	PhD	Science	85	Employed
4	Postgraduate	Arts	98	Unemployed
5	PhD	Arts	90	Employed
6	Undergraduate	Science	88	Unemployed
7	Postgraduate	Arts	93	Employed
8	Undergraduate	Arts	94	Unemployed
9	PhD	Science	86	Employed
10	Undergraduate	Arts	91	Unemployed
11	Postgraduate	Arts	96	Employed
12	PhD	Science	87	Employed
13	Undergraduate	Science	90	Unemployed
14	Postgraduate	Science	95	Employed

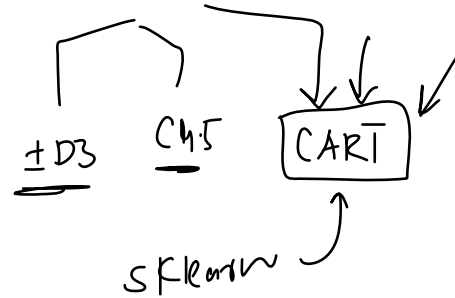
input

output

{UG, Science, 90}

15 rows

1.5 data



The CART Algorithm - Classification

21 July 2023 10:26

sklearn

Given training vectors $x_i \in R^n, i=1, \dots, I$ and a label vector $y \in R^I$, a decision tree recursively partitions the feature space such that the samples with the same labels or similar target values are grouped together.

Let the data at node m be represented by Q_m with n_m samples. For each candidate split $\theta = (j, t_m)$ consisting of a feature j and threshold t_m , partition the data into $Q_m^{left}(\theta)$ and $Q_m^{right}(\theta)$ subsets

$$Q_m^{left}(\theta) = \{(x, y) | x_j \leq t_m\}$$

$$Q_m^{right}(\theta) = Q_m \setminus Q_m^{left}(\theta)$$

The quality of a candidate split of node m is then computed using an impurity function or loss function $H()$, the choice of which depends on the task being solved (classification or regression)

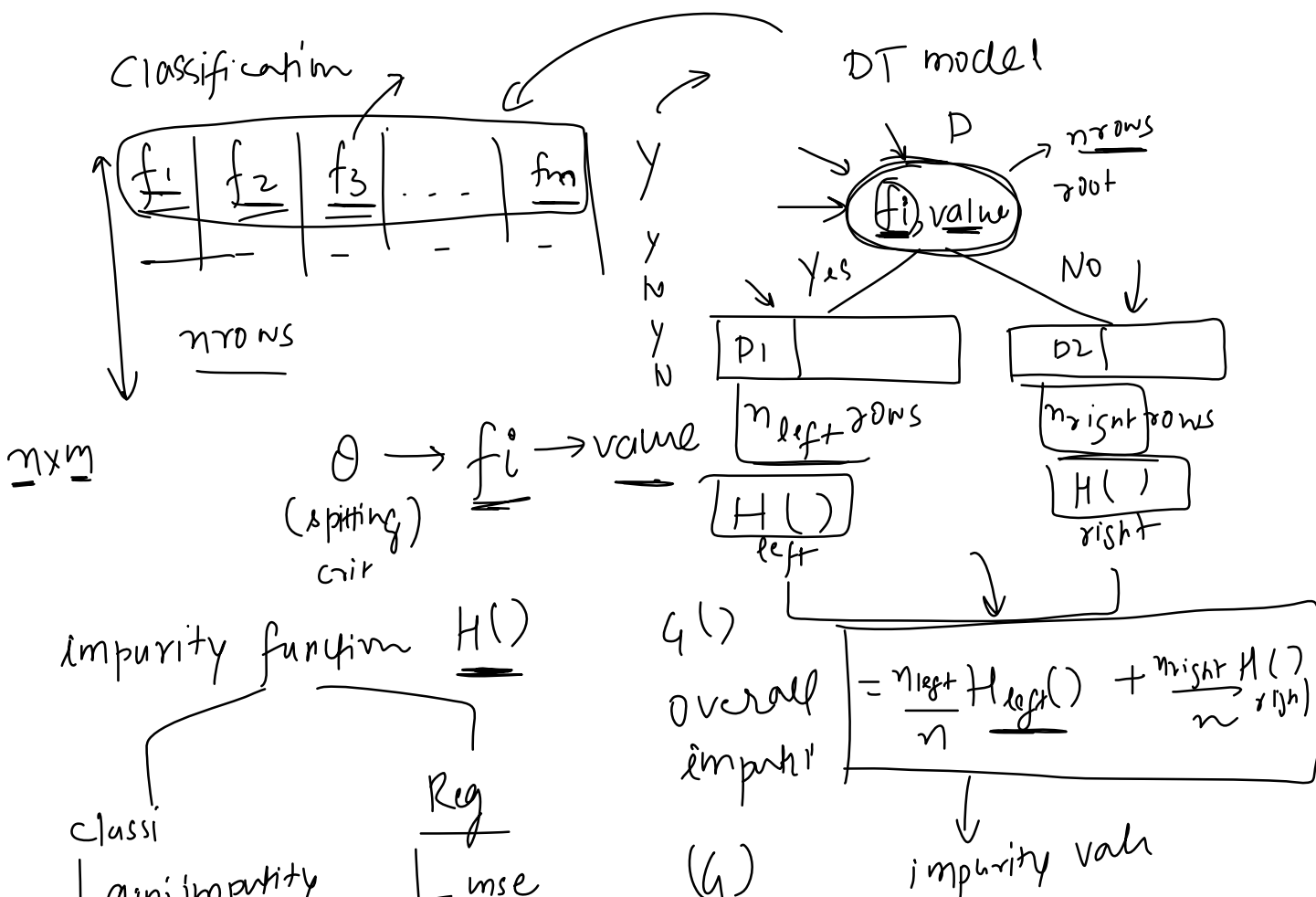
$$G(Q_m, \theta) = \frac{n_m^{left}}{n_m} H(Q_m^{left}(\theta)) + \frac{n_m^{right}}{n_m} H(Q_m^{right}(\theta))$$

Select the parameters that minimises the impurity

$$\theta^* = \underset{\theta}{\operatorname{argmin}} G(Q_m, \theta)$$

Recurse for subsets $Q_m^{left}(\theta^*)$ and $Q_m^{right}(\theta^*)$ until the maximum allowable depth is reached, $n_m < \min_{samples}$ or $n_m = 1$.

CART → Classification and Regression trees
 Algorithm to grow decision tree
sklearn
 Classi → Regression



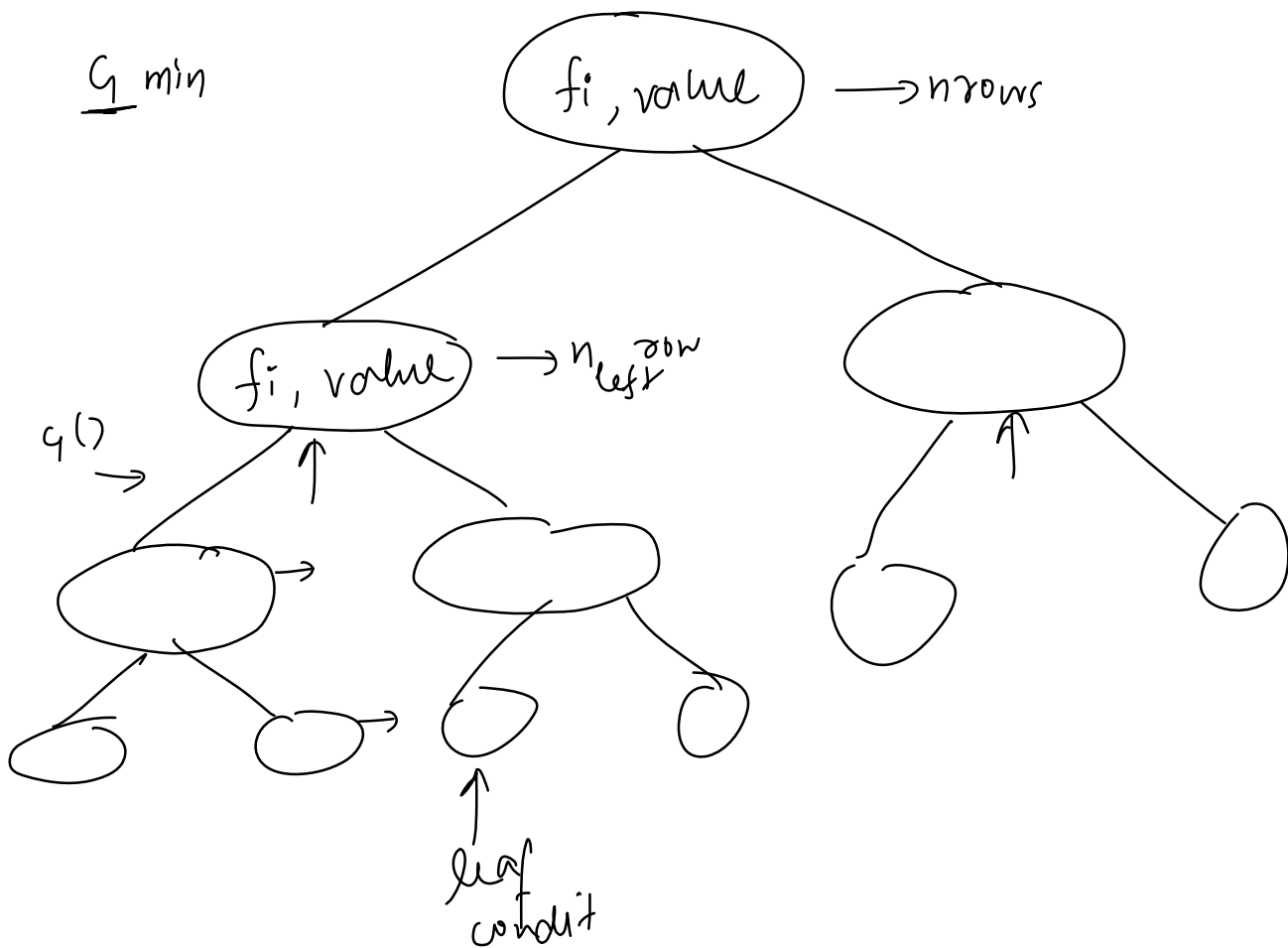
classi
 + gini impurity
 + entropy

~~+~~ mse
~~+~~ mac

(G)

↓
 impurity val

$f_1, value_1 \rightarrow G() \leftarrow (min)$
 $f_2, value_2 \rightarrow G()$
 $\dots \rightarrow G()$



Splitting Categorical Features

21 July 2023 17:45

	Degree_Type	Field	Average_Grade	Job_Outcome
0	Undergraduate	Science	89	Employed
1	Undergraduate	Arts	92	Unemployed
2	Postgraduate	Science	95	Employed
3	PhD	Science	85	Employed
4	Postgraduate	Arts	98	Unemployed
5	PhD	Arts	90	Employed
6	Undergraduate	Science	88	Unemployed
7	Postgraduate	Arts	93	Employed
8	Undergraduate	Arts	94	Unemployed
9	PhD	Science	86	Employed

multi class

binary class

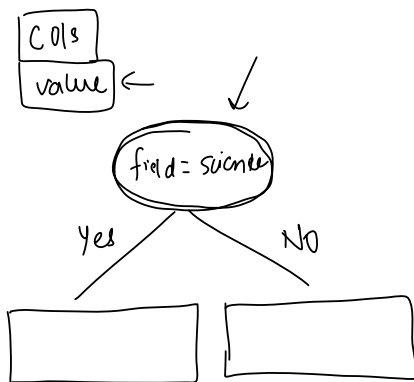
Numerical

$$gini\ impurity = 1 - \sum_{i=1}^k p_k^2$$

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

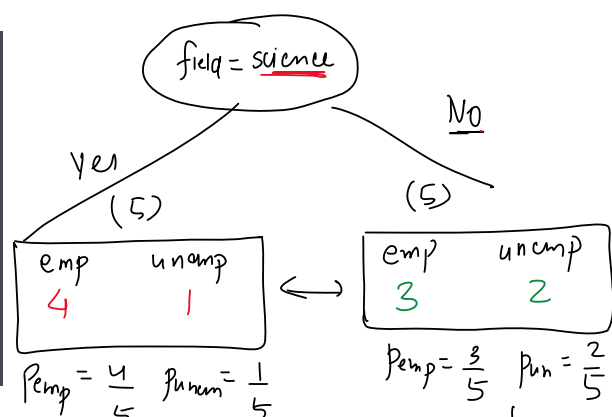
p_k = prob of each class

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



gini impurity

	Degree_Type	Field	Average_Grade	Job_Outcome
0	Undergraduate	Science	89	Employed
1	Undergraduate	Arts	92	Unemployed
2	Postgraduate	Science	95	Employed
3	PhD	Science	85	Employed
4	Postgraduate	Arts	98	Unemployed
5	PhD	Arts	90	Employed
6	Undergraduate	Science	88	Unemployed
7	Postgraduate	Arts	93	Employed
8	Undergraduate	Arts	94	Unemployed
9	PhD	Science	86	Employed



$p_{emp} = \frac{4}{5}$ $p_{unemp} = \frac{1}{5}$

$gini$

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

$$1 - \frac{16}{25} - \frac{1}{25} = \frac{8}{25}$$

$p_{emp} = \frac{3}{5}$ $p_{un} = \frac{2}{5}$

$gini$

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

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

$$\frac{5}{10} \times \frac{8}{25} + \frac{5}{10} \times \frac{12}{25}$$

$0.16 + 0.24 \rightarrow$

G() value

0.40

G → (field, science)

0.11

0.11

multi class

values → UG / PG / PhD

	Degree_Type	Field	Average_Grade	Job_Outcome
0	Undergraduate	Science	89	Employed
1	Undergraduate	Arts	92	Unemployed
2	Postgraduate	Science	95	Employed
3	PhD	Science	85	Employed
4	Postgraduate	Arts	98	Unemployed
5	PhD	Arts	90	Employed
6	Undergraduate	Science	88	Unemployed

degree, UG

$G() = 0.5$

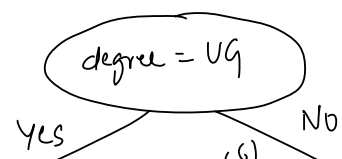
degree, PG

$G()$

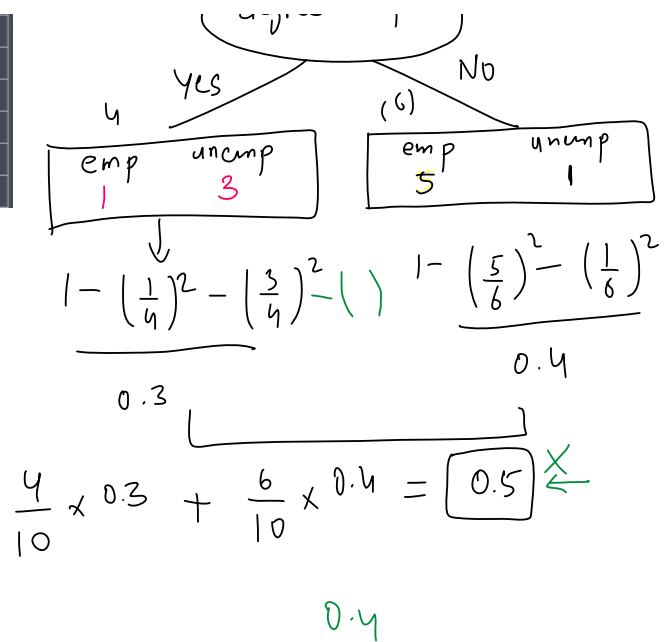
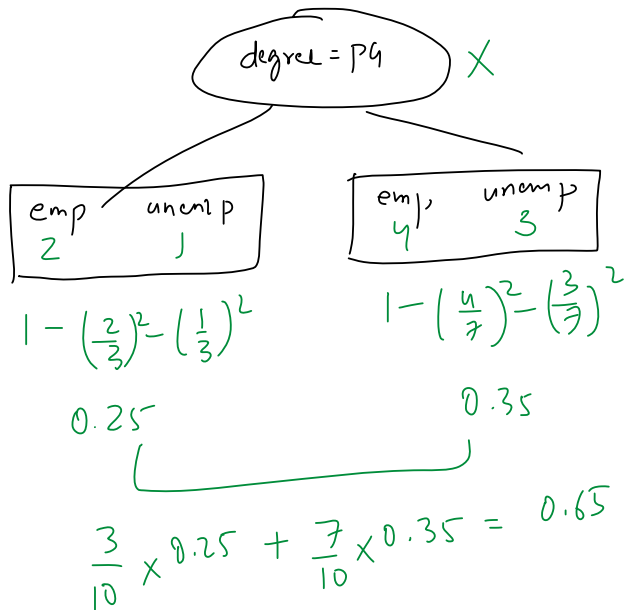
degree, PhD

$G()$

min



4	Postgraduate	Arts	98	Unemployed
5	PhD	Arts	90	Employed
6	Undergraduate	Science	88	Unemployed
7	Postgraduate	Arts	93	Employed
8	Undergraduate	Arts	94	Unemployed
9	PhD	Science	86	Employed

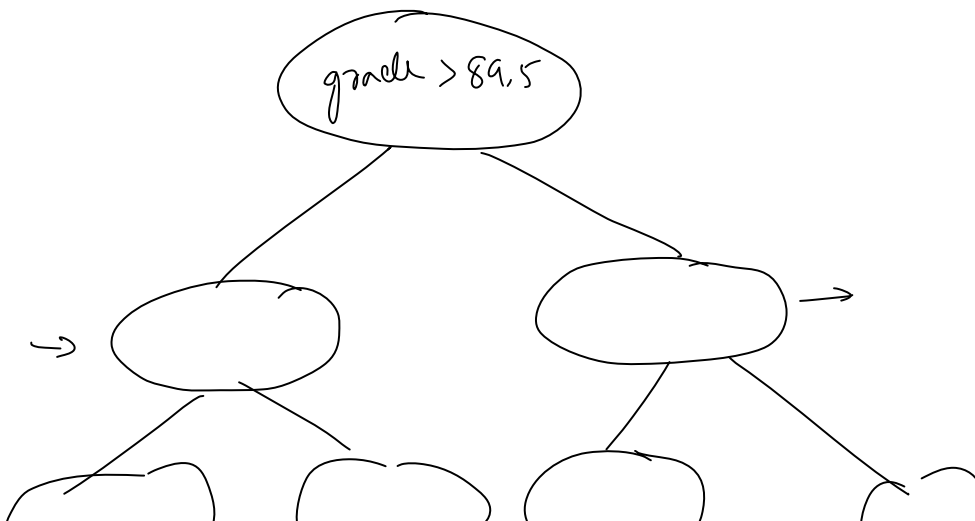
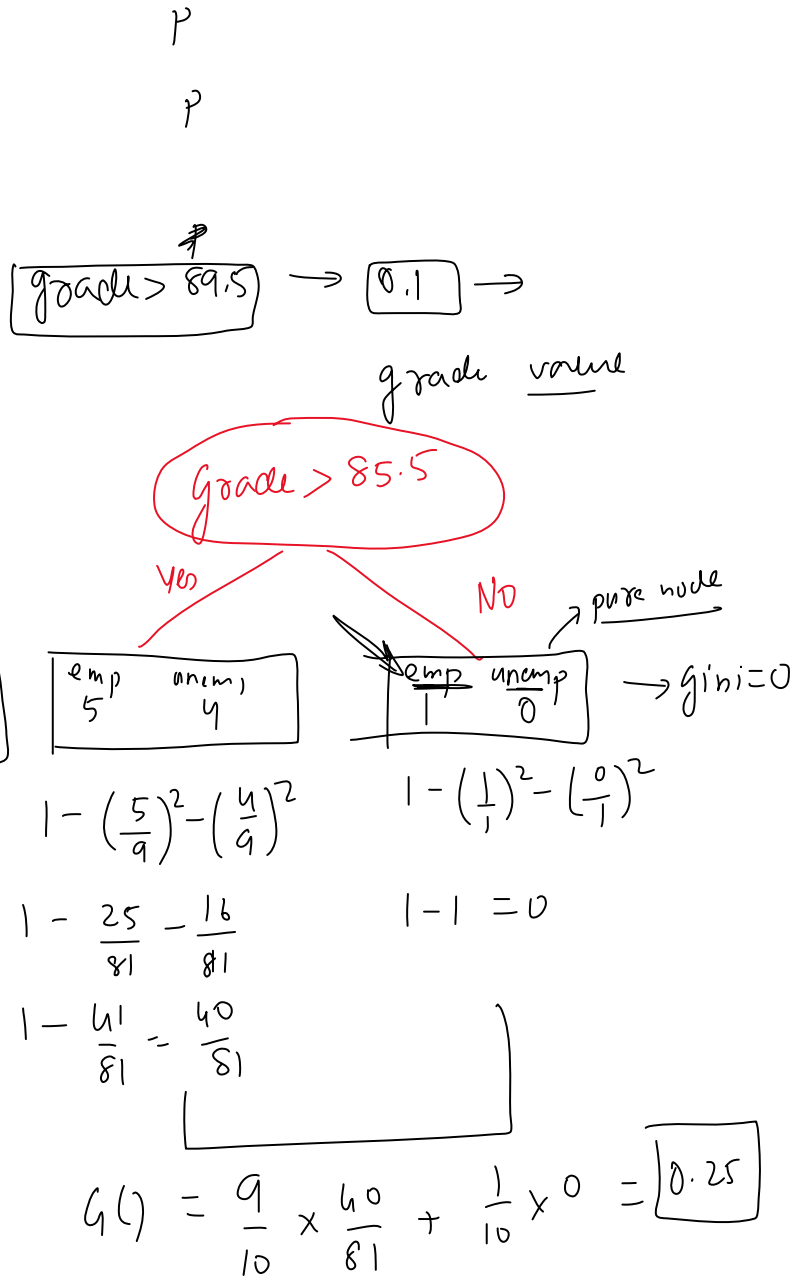


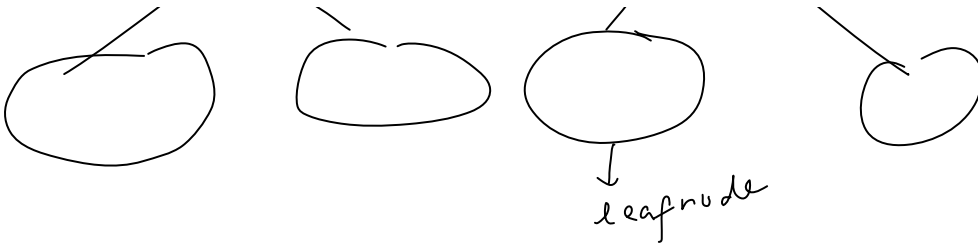
Splitting Numerical Features

21 July 2023 17:45

Average_Grade	Job_Outcome
89	Employed
92	Unemployed
95	Employed
85	Employed
98	Unemployed
90	Employed
88	Unemployed
93	Employed
94	Unemployed
86	Employed

0.4 →	85.5	85 - emp
0.3 →	87	86 - emp
0.2 →	88.5	88 - unemp
0.1	89.5	89 - emp
	91	90 - emp
0.8	92.5	92 - unemp
0.7	93.5	93 - emp
	94.5	94 - unemp
		95 - emp
		98 - unemp





Understanding Gini Impurity?

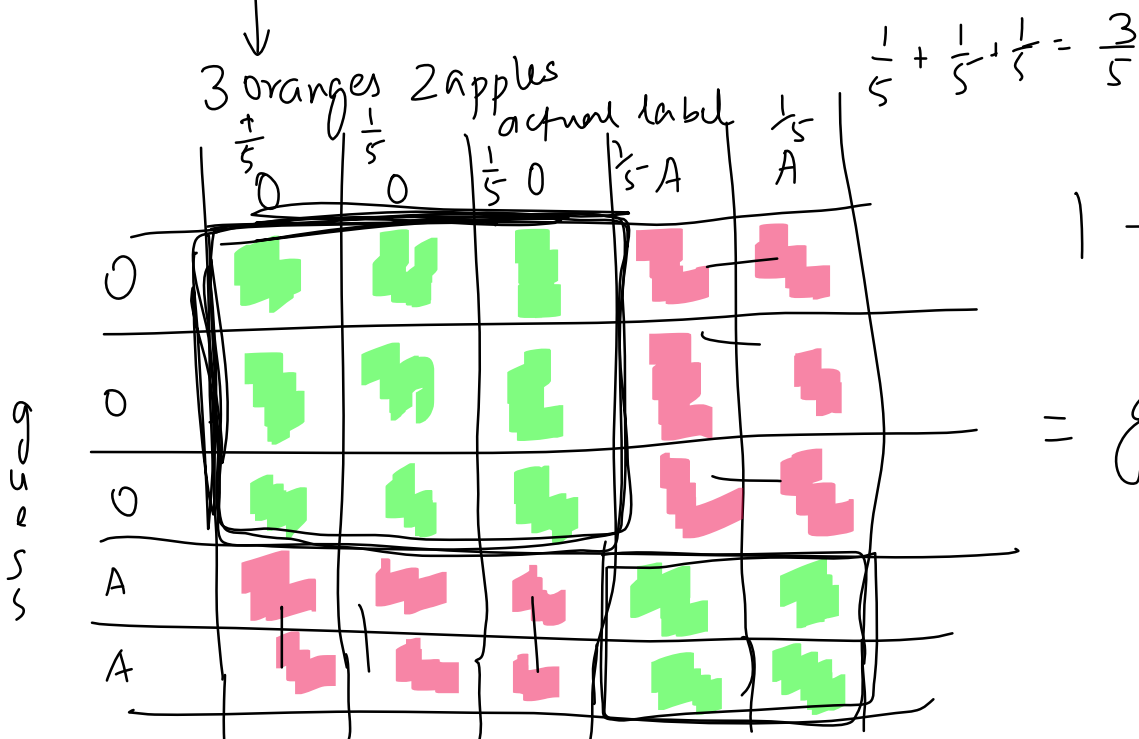
21 July 2023 16:40

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.

$$H() = 1 - \sum_{i=1}^K p_i^2 \quad \longleftrightarrow$$

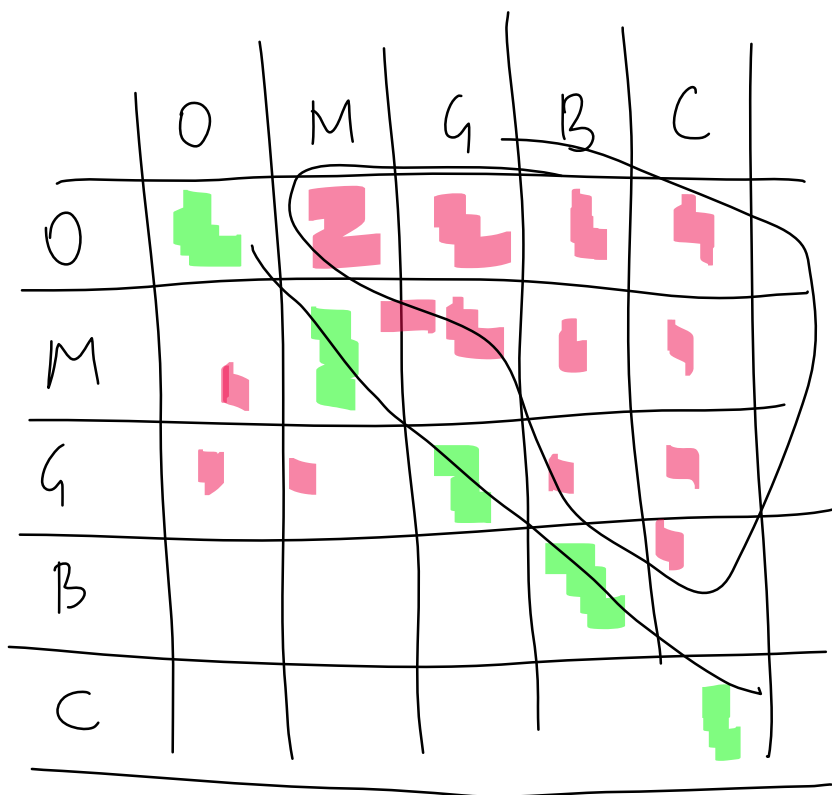
O_1, O_2, O_3
 A_1, A_2

O_1, O_2, O_3
 O_4, O_5

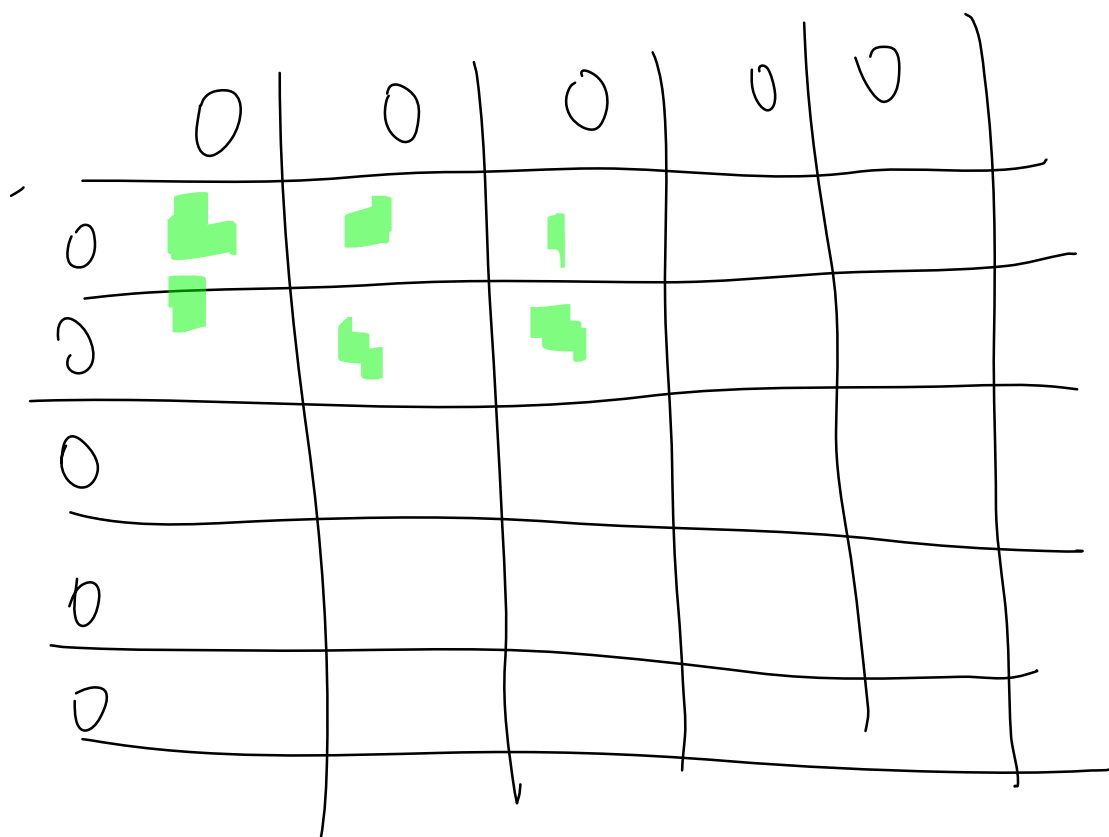


$$1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2$$

= gini impurity



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






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

$$1 - 1 = 0$$

$$1 > g_{in} \geq 0$$

0 → preferred

	o	m	g	b
o				
m				
g				
b				

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

$$\frac{9 - 3}{9} = \left(\frac{2}{3}\right)$$

$$1 - \frac{1}{9} - \frac{1}{9} - \frac{1}{9}$$

$$0.66$$

b

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

(1)

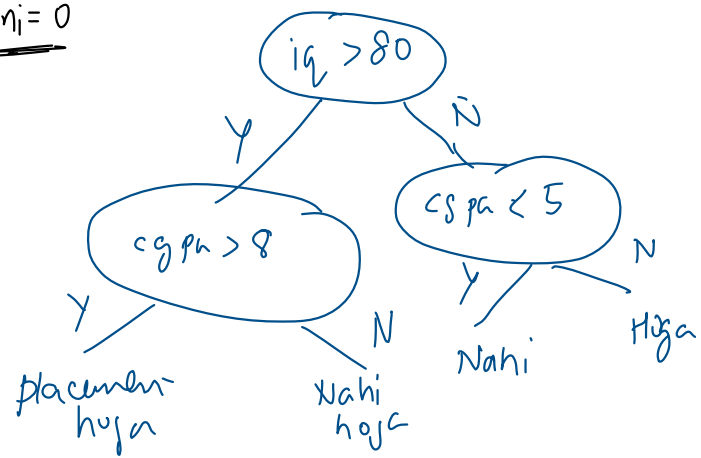
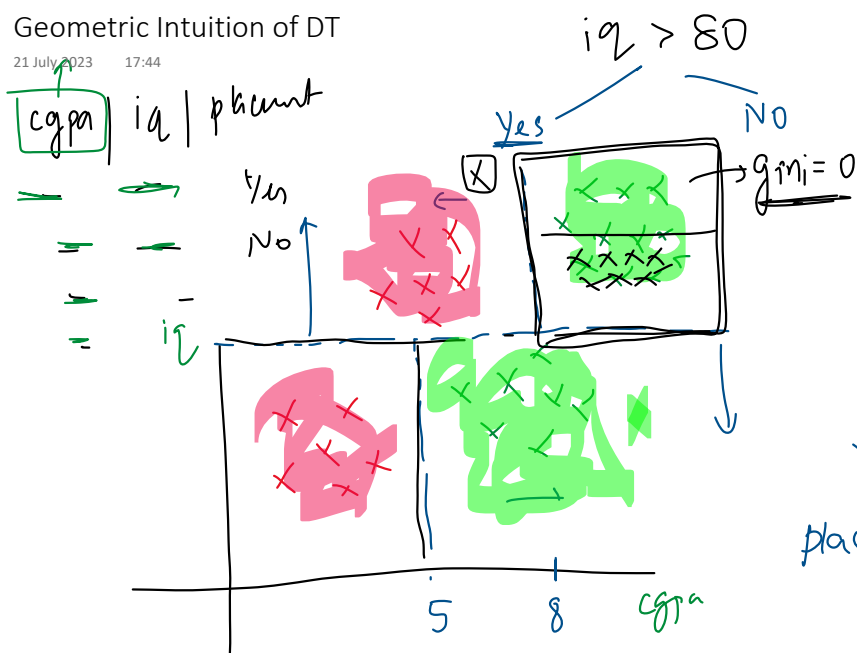
$$1 - \frac{4}{16} = \frac{12}{16}$$

$$0.75$$

Geometric Intuition of DT

21 July 2023

17:44



Code

21 July 2023 17:49

The CART Algorithm - Regression

	Subject	Grade_Level	Hours_Studied	Test_Score
0	Math	Freshman	4	59
1	Physics	Freshman	1	82
2	Physics	Freshman	4	81
3	Math	Junior	6	60
4	Physics	Sophomore	1	73
5	Physics	Junior	3	85
6	Physics	Junior	4	61
7	Physics	Freshman	9	78

Code

21 July 2023 17:49

Advantages and Disadvantages

21 July 2023 17:36

Advantages

- Simple to understand and to interpret. Trees can be visualized.
- Requires little data preparation. Other techniques often require data normalization, dummy variables need to be created and blank values to be removed. Note however that this module does not support missing values.
- The cost of using the tree (i.e., predicting data) is logarithmic in the number of data points used to train the tree.
- Able to handle both numerical and categorical data.
- Can work on non-linear datasets

Disadvantages

- Decision-tree learners can create over-complex trees that do not generalize the data well. This is called overfitting. Mechanisms such as pruning, setting the minimum number of samples required at a leaf node or setting the maximum depth of the tree are necessary to avoid this problem.
- Decision trees can be unstable because small variations in the data might result in a completely different tree being generated. This problem is mitigated by using decision trees within an ensemble.
- Predictions of decision trees are neither smooth nor continuous, but piecewise constant approximations as seen in the above figure. Therefore, they are not good at extrapolation.
- Decision tree learners create biased trees if some classes dominate. It is therefore recommended to balance the dataset prior to fitting with the decision tree.
- Predictions of decision trees are neither smooth nor continuous, but piecewise constant approximations as seen in the above figure. Therefore, they are not good at extrapolation.

