



POORNIMA FOUNDATION

LECTURE NOTES

Campus: PCE

Course: BTECH in CSE

Class/Section: III Yr. Section- A

Date: 24-02-21

Name of Faculty: Praveen Kumar Yadav

Name of Subject: Machine Learning

Code: 6CS4-02

Date (Prep.): 24-02-21

Date (Del.): 24-02-21

Unit No.: 6

Lect. No.: 07

OBJECTIVE: To be written before taking the lecture (Pl. write in bullet points the main topics/concepts etc., which will be taught in this lecture)

Gini Impurity for a Decision Tree

IMPORTANT & RELEVANT QUESTIONS:

what is Gini Impurity for a decision tree?

FEED BACK QUESTIONS (AFTER 20 MINUTES):

what are the factors to determine the decision tree?

OUTCOME OF THE DELIVERED LECTURE: To be written after taking the lecture (Pl. write in bullet points about students' feedback on this lecture, level of understanding of this lecture by students etc.)

good.

REFERENCES: Text/Ref. Book with Page No. and relevant Internet Websites:

scikit with ML.

Campus: PCE. Course: BTECH

Name of Faculty: Praveen Kumar Yadav

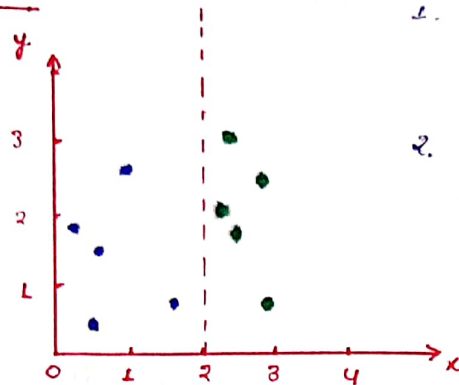
Class/Section: DI CE - A

Name of Subject: ML

Date: 24-02-21

Code: CEG-02

Gini Impurity:-



1. Randomly pick a datapoint in our dataset, then
2. Randomly classify it according to class distribution in the dataset.

(what is the probability we classify the data point

Incorrectly?)

The answer to that question is the Gini Impurity.

Probability

For eg-

Event	Probability
Pick blue classify Blue	25%
Pick blue classify Green	25%
Pick Green " Blue	25%
Pick Green " Green	25%

} Here we classify incorrectly in 2 events.

So our total probability is - $25\% + 25\% = 50\%$.

$$\text{Gini Impurity} = 0.5$$

So Gini Impurity is -

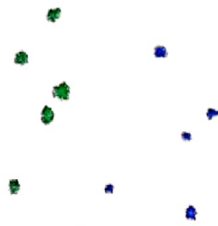
$$G = \sum_{i=1}^C P(i) * (1 - P(i))$$

So here we have $C=2$, $P(1) = 0.5$
 $P(2) = 0.5$

$$G = [P(1) * (1 - P(1))] + [P(2) * (1 - P(2))]$$

$$= [0.5 * (1 - 0.5)] + [0.5 * (1 - 0.5)]$$

$$G = 0.5$$



$P(1) = 0.5$
 $P(2) = 0.5$

Now calculate Gini Impurity After split in two branches.

$$P(1) = 1$$

$$P(2) = 0$$

$$G_{left} = 1 * (1 - 1) + 0 * (1 - 0) = 0$$

$$G_{right} = 0 * (1 - 0) + 1 * (1 - 1) = 0$$

$$P(1) = 0$$

$$P(2) = 1$$

so perfect split turned a dataset with 0.5 impurity into 2 branches with 0 impurity.

So a Gini Impurity of 0 is the lowest and best possible impurity. It can only be achieved when everything is from the same class.

Campus: PCE, Course: BTECH

Name of Faculty: Praveen Kumar Yadav

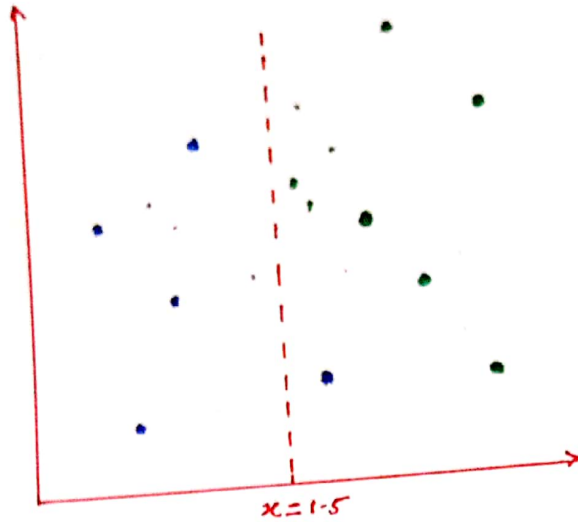
Class/Section: 20CE019

Name of Subject: ML

Date: 24.05.19

Code: 6556-01

Imperfect split:-



Here left branch has only blue so

$$G_{\text{left}} = 0$$

Right branch has 1 blue and 5 greens, so

$$G_{\text{right}} = \frac{1}{6} \times (1 - \frac{1}{6}) + \frac{5}{6} (1 - \frac{5}{6})$$

$$= \frac{5}{18} = 0.278$$

we have already calculated the Gini Impurity for -

Before split (the entire data set) = 0.5

Left branch: 0

Right " : 0.278

So determine the quality of split by weighing the Impurity of each branch by how many element it has -

$$\text{i.e. } (0.4 \times 0) + (0.6 \times 0.278) = 0.167$$

So Amount of Impurity we have "removed" with the

split is - $0.5 - 0.167 = 0.333$

↑
Gini Gain.

This is what's used to pick the best split in

decision tree.

Higher Gini Gain = Better split.

So **Gini Impurity** - is the probability of incorrectly classifying a randomly chosen element in the dataset. if it were randomly labeled according to the class distribution in the data set.

MINIMA FOUNDATION

DETAILED LECTURE NOTES

Apus: PCE. Course: BTECH
Name of Faculty: Praveen Kumar Yadav

Class/Section: VI CSE - A

Name of Subject: ML

Date: 24-02-20

Code: 654-02

when training a decision tree, the best split is chosen by "Maximizing the Gini gain" which is calculated by subtracting the weighted impurities of the branches, from the original impurity.