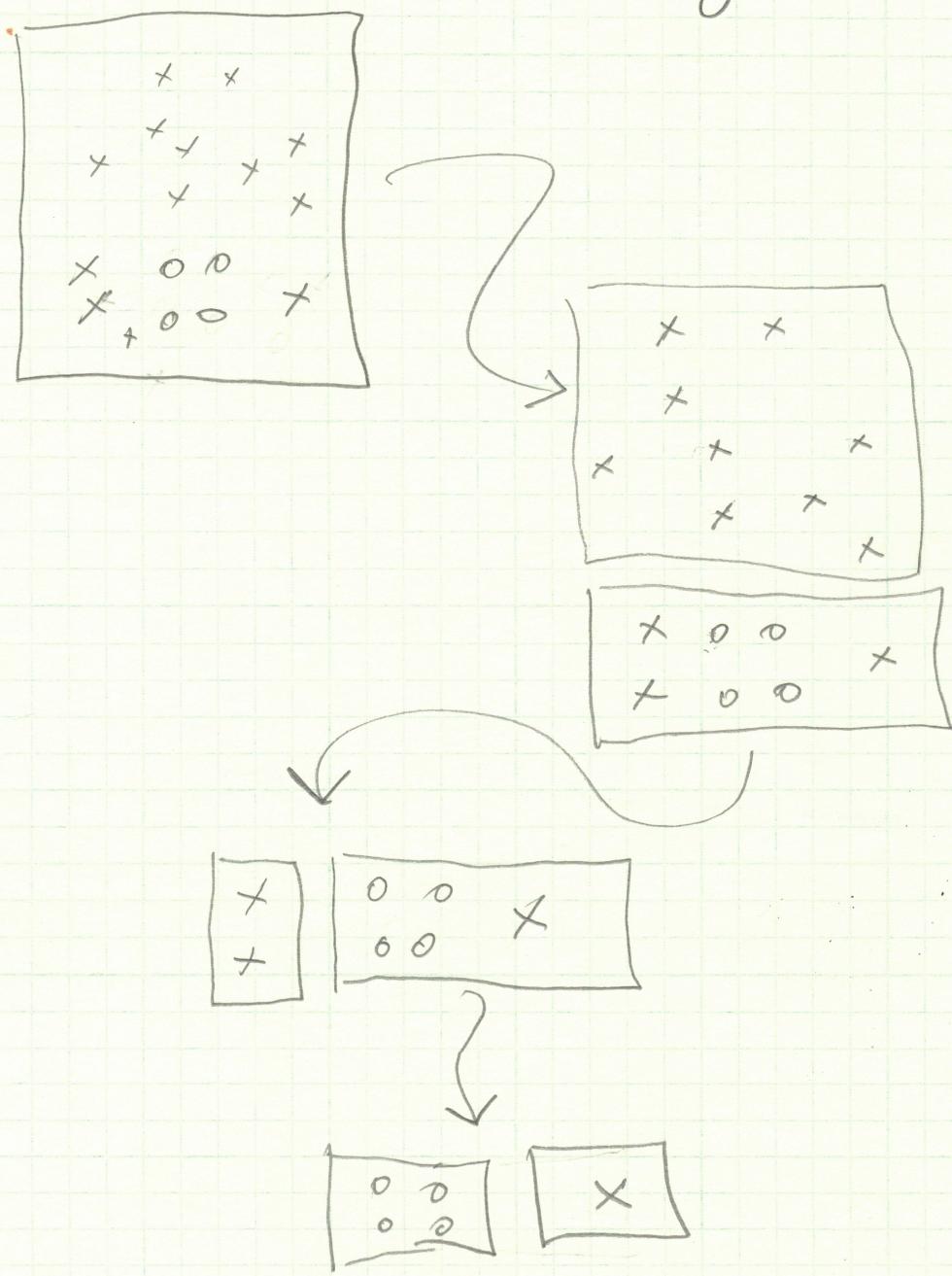


①

Decision Trees

Notes based on article
by Robert Kubler
towards Data Science

- Decision Trees separate / splits data into groups based on "information gain"
- One really good way to implement the algorithm is to do it recursively



①

We stop splitting when:

- the split result in a part being empty
- A certain recursion depth is reached
- not much data is in a split, making more splits unnecessary

> However, looking at the example, notice we split data into 2 parts:

- a homogeneous part
- a non-homogeneous part

The recursion, or splitting follows the non-homogeneous part of the data.

* * * * * +
* * * * + *
* * * 0 0 0
* * * 0 0 0

The idea idea
is too split
the data into
two parts -
these parts
should be
"cleaner" than
the complete data
before

$$\begin{array}{cccccc}
 & \times & \times & \times & \times & \times \\
 & \times & \times & \times & \times & \times \\
 \hline
 & \times & + & 0 & 0 & 0
 \end{array}$$

clean

messy

next
we split messy
set here

- We need a measure of cleanliness
here are 3 ways

$$x = [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0] \quad \text{clearest}$$

$$y = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0] \text{ mel!}$$

$$z = [10110010] \text{ messyest}$$

- A simple measure (for 2 class data)

$$\frac{| \text{number class}_0 - \text{number class}_1 |}{\text{total set size}} = \frac{| n_0 - n_1 |}{n}$$

(3)

$$\frac{|n_0 - n_1|}{n} = \frac{|(n-n_1) - n_1|}{n} = \frac{|n - 2n_1|}{n}$$

$$= \left| \frac{n}{n} - \frac{2n_1}{n} \right| = \left| 1 - 2\frac{n_1}{n} \right|$$

let $P_1 = \frac{n_1}{n}$; so $|1 - 2P_1|$

P_1 is the percentage of n_1 labels

let $n = n_0 + n_1$

- This is set up for 2 class data

$$x = [0 0 0 0 0 0 0 0]$$

$$|1 - 2P_1| = |1 - 2 \cdot 0| = \boxed{1}$$

$$P_0 = \frac{8}{8} = 1; P_1 = \frac{0}{8} = 0$$

$$y = [1 0 0 0 0 0 1 0]$$

$$P_0 = \frac{6}{8} = \frac{3}{4}; P_1 = \frac{2}{8} = \boxed{\frac{1}{4}}$$

$$|1 - 2 \cdot \frac{1}{4}| = .5$$

$$z = [1 0 1 1 0 0 1 0]$$

$$P_0 = \frac{4}{8} = \frac{1}{2}, P_1 = \frac{4}{8} = \frac{1}{2}$$

$$|1 - 2 \cdot \frac{1}{2}| = \boxed{0}$$

(4)

A more general formula is the Gini impurity measure:

$$g_{\text{ini}}(P_0, P_1) = P_0(1 - P_0) + P_1(1 - P_1)$$

The above is the 2 class case, below is the multi-class case

$$g_{\text{ini}}(P) = \sum_{i=1}^m P_i(1 - P_i)$$

here $P_0 = \frac{n_0}{n}$, $P_1 = \frac{n_1}{n}$

and

$$P_i = \frac{n_i}{n}$$

> looking at X, y, and z

X $g_{\text{ini}}(1, 0) = \frac{8}{8}(1 - 1) + \frac{0}{8}(1 - 0) = 0$

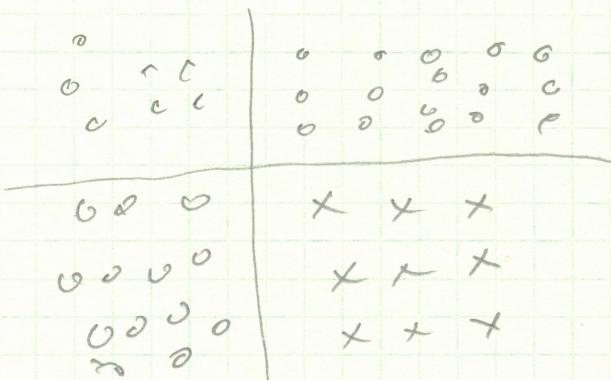
Y $g_{\text{ini}}\left(\frac{6}{8}, \frac{2}{8}\right) = \frac{6}{8}\left(1 - \frac{6}{8}\right) + \frac{2}{8}\left(1 - \frac{2}{8}\right)$
 $\frac{6}{8}\left(\frac{2}{8}\right) + \frac{2}{8}\left(\frac{6}{8}\right)$
 $\frac{12}{64} + \frac{12}{64} = \frac{24}{64} = \frac{3}{8}$

Z $g_{\text{ini}}\left(\frac{1}{2}, \frac{1}{2}\right) = \frac{1}{2}\left(1 - \frac{1}{2}\right) + \frac{1}{2}\left(1 - \frac{1}{2}\right)$
 $\frac{1}{4} + \frac{1}{4} = \frac{1}{2}$

⑤ So unlike our previous calculations, the more "disordered" or "impure" the set is, the higher the gini value

▷ Finding where to split the data.

We will split the data along either the X or y axis using a gini measure



Suppose 75% of the points are o
25% of the points are x

$$\text{gini}(.75, .25) = .25(1 - .25) + .75(1 - .75)$$
$$= .375$$



overall gini

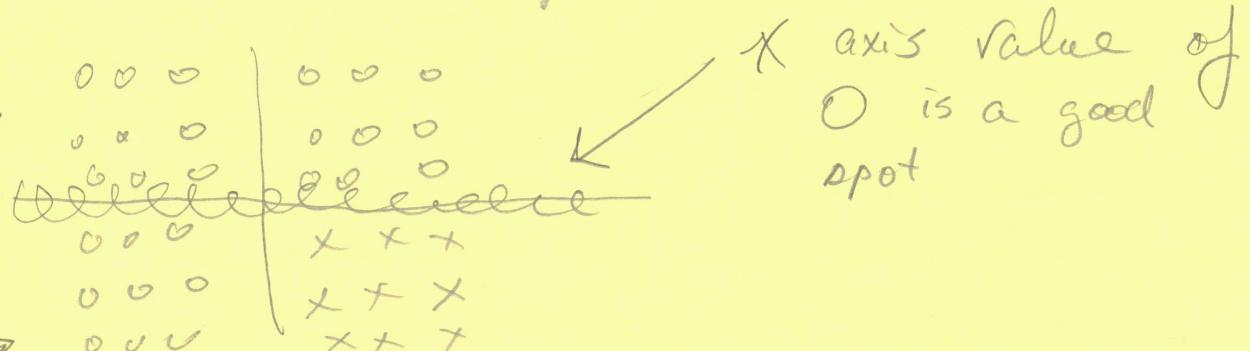
6

We want to find a split where the sum on 1 side of split shows

- o Very pure

and the other side

- not so pure



$$g_{\text{ini}}(1,0) = 1(1-1) + 0(1-0) \\ = 0$$

$$g_{\text{ini}}(0.5, 0.5) = 0.5(1 - 0.5) + 0.5(1 - 0.5) \\ = 0.25 + 0.25 \\ = 0.5$$

So using this split, we can average the gini values

$$\frac{(0 + .5)}{2} = .25$$

(7)

- Remember, the gini value without the split was .375

$.25 < .375$ so this is a better measure

- the lower the gini number, the better

→ we could have done a split on $y = 0$ and that would have been good too.

last thing, we should use a weighted average when calculating our average Gini #'s

Suppose part 1 has 50 points
part 2 has 450 points

instead of calculating

$$\text{gini Average} = \frac{(0 + .5)}{2} = .25$$

we should do:

$$\text{gini Average} = \frac{50}{50+450} (0) + \frac{450}{50+450} (.5) = \underline{\underline{.45}}$$

(8)

How do we find the split?

o o	o		o o	o
o o	o		o o	o
o o	o		o o	o

o o o	x x x
o o o	x x x
o o o	x x x

while (splitting) {

 find the y range of the data
 for every y value

 find the weighted
 gini average

 y value with lowest gini average
 is y split

 keep the data in the impure side of
 the split

 if both sides have gini values \leq pure
 enough

 Rsplitting = false

⑨ if (splitting) {

 find the x range of the data
 for every x value

 find the weighted
 gini average

 x value with the lowest gini
 average is the y split

 keep the data in the impure side
 of split

 if both sides have gini values \leq pure
 enough
 splitting = false

}

} It end while */