## Decision tree (for classification algorithms or regression (continuous))

* It has two version

         C5.0     CART

uses in

     classification            regression.
     (mostly used)
     (which output
     to which class
     its belong)

* In decision tree there two nodes

   Decision node        leaf node
   (have branches)     (have not branches.

* A decision will happen with the help of      yes
                                                  NO

simple example of decision tree

       A person want      → Decision node
       wants to eat
         Ice cream

   yes            No         → leaf node.

  go to the           sit at
  ice cream          home.
  store

* when decision tree classifiers things unto

   categories         numeric values

  classification       Regression
  tree             tree.

* we will focus on decision tree (Classification)

terminology of decision tree

Root

internal nodes or branches

internal nodes or branches

leaf node / node / leaf

leaf node / node

leaf node node

leaves node or leaf

* lets take an example for decision tree.

| loves ice cream | loves soda | Age | loves chocolate |
|---|---|---|---|
| Y | Y | 7 | N |
| N | N | 12 | N |
| No | Y | 18 | Y |
| ye | Y | 35 | Y |
| y | Y | 38 | Y |
| | N | 50 | N |
| N | N | 83 | N |

Y — yes
N — No

loves ice-cream POPCORN

T          F

loves chocolate        loves chocolates

Y/N         Y/N

N           Y
1           2

3:3        1:3

similarly

loves Soda

T          F

loves chocolate        loves chocolate

Y          Y

1:3        1:1        Y/N        N
                       0          3

* how to identify Root nodes

⇒ By looking independent variable.

⇒ In previous example you can see the split hasn't happened unto equal parts (yes or no) these splits are called <u>impurity</u>.

⟶ to solve this using methods

```
                ╱  ╲
               ╱    ╲
          Gini        Entropy ⟶ information
        impurity                    gain.
        (mostly
         used)
```

Gini impurity

how do calculate gini impurity

★ Gini impurity for a leaf = $1 - ($ the probability of yes $)^2$
                                   $- ($ the probability of No $)^2$

take previous example
of ice-cream
(left side)

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

$$= 0.375$$

and calculate similarly for right side.

$$= 0.444.$$

Now calculating the gini impurity

total Gini = weighted average of Gini impurities for the
impurity      the leaves

$$= \left(\frac{4}{4+3}\right)0.375 + \left(\frac{3}{4+3}\right)0.444$$

total impurity in leave on left

$\left(\text{total impurity in both leave}\right)$

$\left(\text{Gini impurity}\right)$ + ("similarly" for right) Gini impurity

Conclusion is

Gini impurity for clones ice cream = 0.405.

this is we have done just for categories what about numeric data

from the previous example we will take age column
step ① calculate average of adjacent date
② then calculate Gini impurity for each for value

Age
7  $\quad \frac{7+12}{2} = 9.5$  for example  → Age ≤ 9.5
12 → 15 → 0.343
18 → 25.5 → 0.476
35 → 36.5 → 0.476
38 → 44 → 0.343
50 → 66.5 → 0.429
83

Age 9.5

i=1 soda
soda → N=1 → 3 3
Y → 0

So total gini impurity = 0.429

So the conclusion is 0.343 is lowest so there are two 0.343 we will pick any one.

Similarly we took out gini impurity for age so we we take out ice-cream and soda and we have found out soda has the lowest (gini) impurity values is 0.214
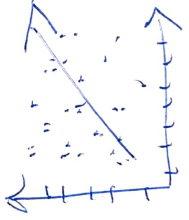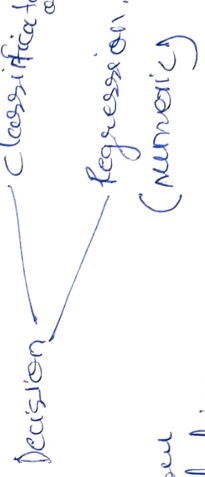
so likewise are we will put soda on the top and according to gini impurity we will distribute categories and whoever got the highest impurity will be the leaf node and those who got the highest majority values will win

example

| Soda | Soda |
|------|------|
| Y    | N    |
| 1    | 0    |
| love soda | dont love soda |

Problem Statement (regression tree)

what if data is not linear and we have to predict the values.
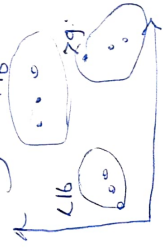
Then → solution is regression tree

So How regression work?

Classification
Decision
regression (numerics)

lets say we have one example.

so you can see there is three part in graph, each has threshold value
so we will plot tree on the basis of this

>16
<16    79
29

Postage ≤16 | Postage >16

average cost cluster
←4.24    2.6

how to build tree for regression.



Drugg effective 100%

0   10   20   30   Drug Dosage (mg).
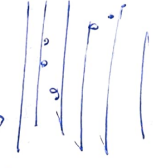
In above graph we have look only two dataPoint (circles). the line is average values on that basis we can build tree.
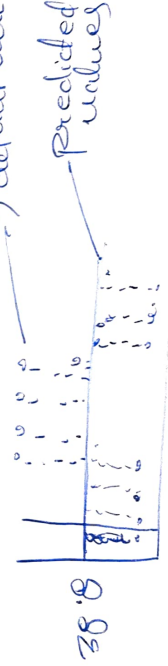
Dosage L average (circle)

average = 0%

average 38%



average 38%

10  29  30  40
all the cluster point

and if the values came (Dosage 20) it showing 100% but according to graph showing 38% don't meetly a poorPrediction. So according graph Decision tree at poorPrediction

* So we can we visualize the predicted values and actual values.
→ predicted values
→ actual values



38%

→ Predicted values

# what the Solution then ?

→ Sure can not be accurate best we can improved the decision tree.

To improve . threshold.

① take any datapoint and determines it to all datapoint by square like this

$$(datapoint - datapoint1)^2 + (datapoint - datapoint 2)^2 + (datapoint1 - datapoint 3)^2$$

= Ans.

whatever you gothe any plot this in graph.

SSR | Dosage threshold.

So whoever go the the lowest SSR that threshold we win, because it has lowest SSR.

So that threshold will be the root.

for example

Dosage ≤ 14.5

what about Next

↓
do the same process to calculate threshold values again.

So do split until the Doesy effectiveness #'s
same or you dont have data point to split.
but this

So divide the data according to groups or
cluster not by data point. If there
are not 1 there many independe
what we there are not 1 there many independe
nd measurable.

Process can lead to model variance

1) step calculate threshold for every column.
2) and pick that threshold (column) whoever
got the lowest values
(3) then make branch assending to descending
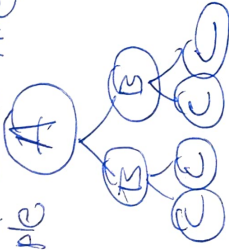(on the basis of threshold)

Decision tree
(work on

Divide and conquer technique (Y and N)

every Node is purely homogeneous region.

at the root node it represent the all the
data set and so on for example

you can check Accuracy of Decision tree by Holdout method.

Decision tree ( C5.0 ).
in C5.0 to select feature at branch it uses something called entrophy.

data with high entrophy ↓ low entroph

↓

diverse, little information

homogeneeu region

most important feature will determine by entrophy Score.

$$entrophy (S) = \sum_{i=1}^{c} -pi \log_2 (pi).$$

p = proportion.
i = features.

to calculate entrophy we have to create random Split for example. (ABCD)
(subset) ① (ABC)
② C(DB)
③ C(DB)
So on.

② then take 1st random and calculate entroph
   just say
     play      $\frac{3}{4}$  $- \log_2 (3/4)$
   $\begin{vmatrix} y \\ y & n \end{vmatrix}$        = 6.4

③ calculate entropy for another random also

for purely humogeneous region = Entropy will always be 0 .

Now we will calculate Entrophy before by calculate -ing complete dataset .

before contineous we should understand why we pecking up important feature by entrophy method or any methode

① if dont select the decision tree will have nodes to much;

② if dont seled the decision tree will dont predict values properly. (impurity will increase)

entrophy :- measures the purity of split .

wrost emtrophy $\quad \frac{N}{50\%} \Big| \frac{Y}{50\%}$ = 1 bits

best entrophy $\quad 100\% \Big| 0\%$ = 0 bits.

entrophy ranges to 0-1 and entrophy unit is bits

★ the split will go on till the leaf node is pearly homogeneous subset.

# information gain.
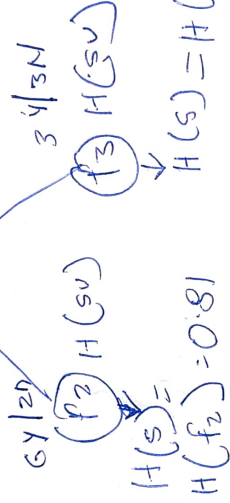
this will give the tree which split is best.

$$= Gain(S, A) = H(s) - \sum_{V \in VAL} \frac{|Sv|}{|S|} \; H(sv)$$

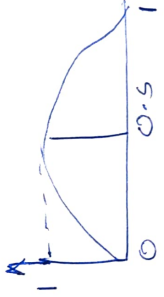$H(s)$ = entrophy for particular subset (after split what is the split.)

### example



$$H(s) = H(f_3) = 1$$

$H(s) = $

$H(f_2) = 0.81$

$$Gain(S, f_1) = H(s) - 8/14 \; H(f_2) - \frac{6}{14} H(f_3)$$

$$= 0.91 - 8/14 \times 0.91 - \frac{6}{14} \times 1$$

$$= 0.049$$

what this value indicate.

if this values values is high it means this structure is best to use the machine algorithm will use different Decision tree structure to identify which one is the best

to calculate impd

to calculate purity in set we use
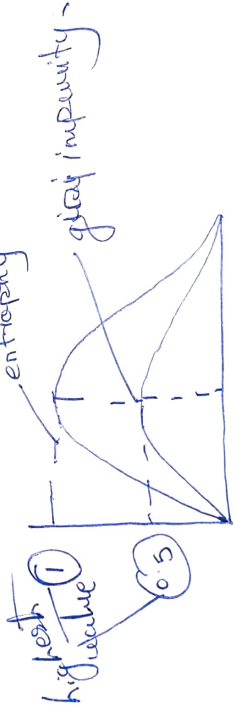
entropy

gini impurity

① entrophy =

② gini impurity.

$$= 1 - \sum_{i=1}^{n} (P)^2$$

$$= 1 - [(P_+)^2 + (P_-)^2]$$



but we are using both to calculate impurity then which one to use and what is the different.

graph for entroph and gini impurity



Gini impurity takes less time to do operation.

* entrop before (after creating random set ) original dataset
before

* information gain ( lies in 1 to 0 ) .

= (P)* entrop - entropy before - entrophy after

* information gain ( How useful is split ).
    0 - worst
    1 - best .

* to avoid overfitting issue we will use pruning

Deceision tree : C 5.0                    ( alpha Have
                                            zero )
                                               ↓
                                          research paper
                                          on decision
                                          tree .

* How long Decision tree
= it will go indefinitly untll algorithm have'nt
    got homogeneous region .

* CART ( classification and Regression trees)