

Decision Tree

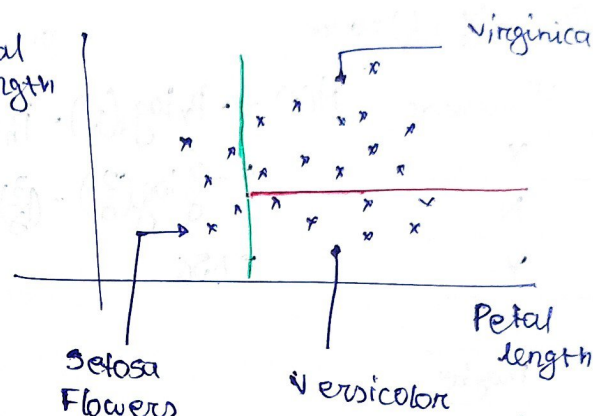
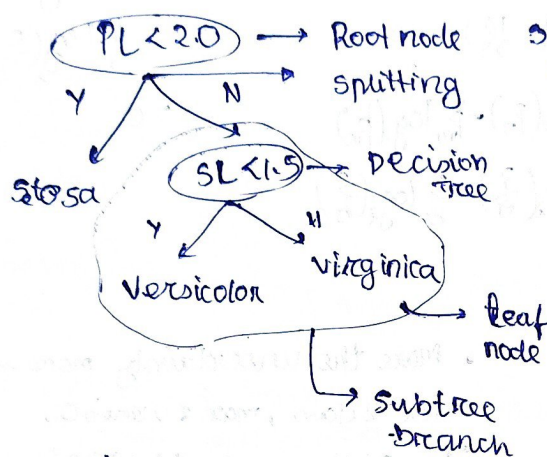
P-1

8:49pm
Thu May 1

Gender	Occupation	Suggestion
F	Student	Crane
F	programer	crithub
M	"	whatsapp
F	"	github
M	Student	Crane
M	Crane	"

```

if occupation == Crane
    print("Crane")
Else
    if gender == Female
        print("github")
    Else
        print("whatsapp")
    
```



Q. which is best features for splitting.
- splitting criteria ?

CART - Classification & Regression Tree

Entropy

• measure of purity/impurity

$$E(S) = - \sum_{i=1}^n P_i \log_2 P_i$$

Ex- if data have two classes

$$E(D) = -P_{yes} \log_2(P_{yes}) - P_{no} \log_2(P_{no})$$

where

P_i is simply the frequentist probability of an element/class 'i' in our data



Age	Purchase
21	Y
51	N
69	Y
86	N
76	N

Purchase
Y
N
N
N
N

Pur
N
N
N
N
N

$$H(d) = -p_Y \log_2(p_Y) - p_N \log_2(p_N)$$

$$= -2/5 \log_2(2/5) - 3/5 \log_2(3/5)$$

$$= 0.97$$

$$H(d) = -p_Y \log_2(p_Y) - p_N \log_2(p_N)$$

$$= -1/5 \log_2(1/5) - 4/5 \log_2(4/5)$$

$$= 0.72$$

$$H(d) = -p_Y \log_2(p_Y) - p_Y \log_2(p_Y)$$

$$= -0/5 \log_2(0/5) - 5/5 \log_2(5/5)$$

$$= 0$$

$$G = 1 - (1/25 + 9/25) = 0.48$$

$$G = 1 - (1/25 + 16/25) = 0.72$$

Purchase
Y
N
Y
N
maybe
N
maybe
maybe

$$H(d) = -p_Y \log_2(p_Y) - p_N \log_2(p_N) - p_m \log_2(p_m)$$

$$= -2/8 \log_2(2/8) - (1/8) \log_2(1/8) - 3/8 \log_2(3/8)$$

$$= 1.56$$

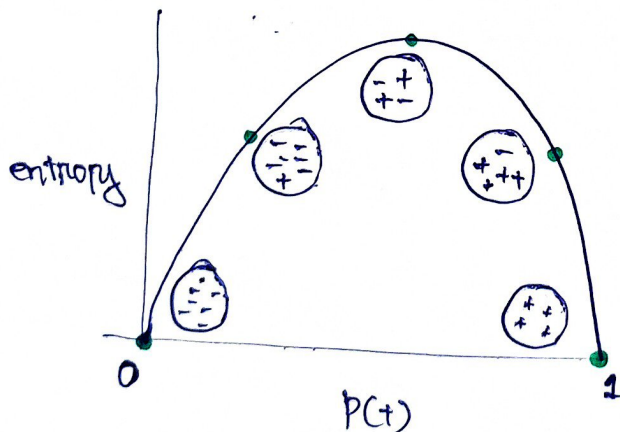
- More the uncertainty more is entropy
- For 2 class, max 1, min 0.
- For 2+ class, max 1+, min 0.

Y → 2

N → 3

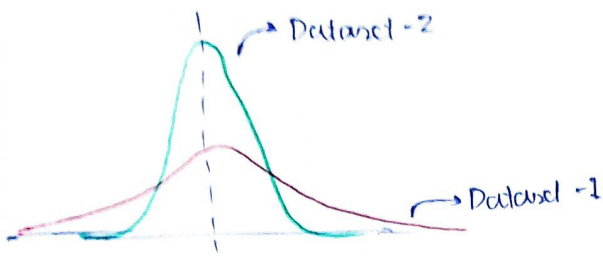
m → 3

Entropy vs Probability



Entropy for Continuous Variables

1. The more you know about the data then less entropy.



which datasets have higher entropy?
- whichever is less peaked.

Information Gain

- is a metric used to train Decision Trees.

$$\text{information gain} = E(\text{parent}) - \left\{ \begin{array}{c} \text{weighted} \\ \text{Average} \end{array} \right\} - E(\text{children})$$

g

Step 1 Entropy of Parent means full dataset. $E(P) = 0.97$

Step 2 Calculate for child

$$\begin{array}{ccc} & \text{outlook} & \\ \swarrow & \uparrow & \searrow \\ \text{sunny} & \text{overcast} & \text{Rain} \\ E(5) = 0.97 & E(0) = 0 & E(5) = 0.7 \end{array}$$

Step 3 weighted Entropy of children

$$\begin{aligned} W.E(\text{children}) &= \frac{5}{14} \times 0.97 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.7 \\ &= 0.69 \end{aligned}$$

Step 4 Info. Gain = $0.97 - 0.69 = 0.28$

So information Gain, when we split this data on the basis of outlook column is 0.28

Step 5 highest Entropy split on that column
Recursively
when Entropy = 0 then stop

Gini Impurity

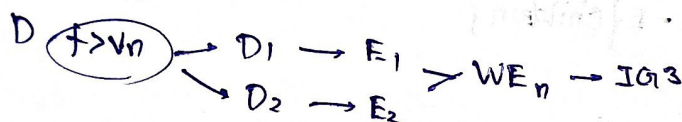
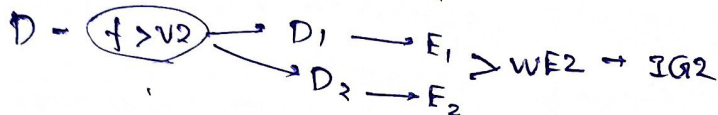
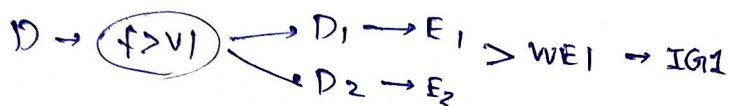
$$G_I = 1 - (P_Y^2 + P_N^2)$$

computationally fast

Numerical Data

S1. Sort data on the basis of one column

S2. Split entire data on the basis of every value of that column



S3. $\max \{ IG_1, IG_2, \dots, IG_n \}$



S4. Do recursively

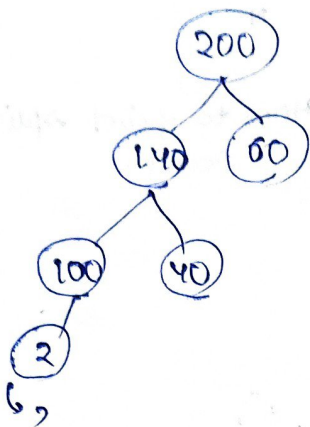
Expensive

- in Train time
- not in Test time

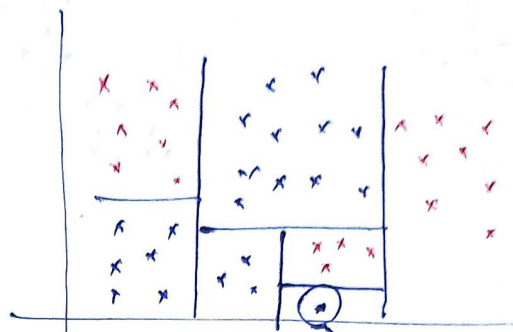
Overfitting

Performs well = Train
badly = Test

max_depth = None



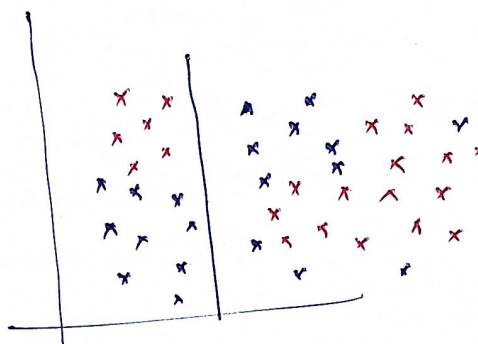
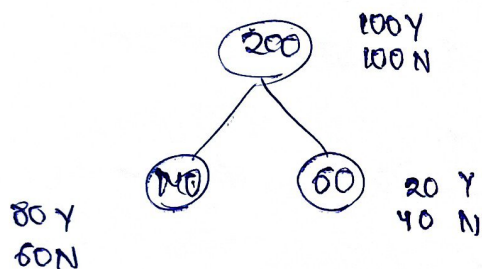
Noisy/outlier
Erroneous



Overfitting
we reach to
leaf 1

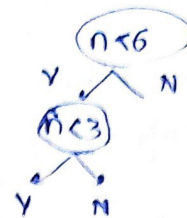
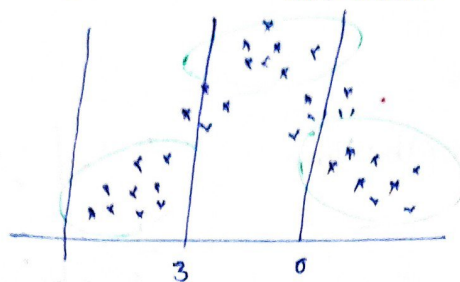
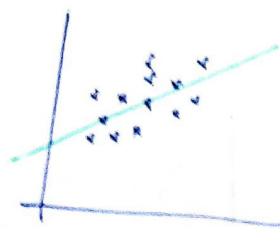
Underfitting

max_depth = 1

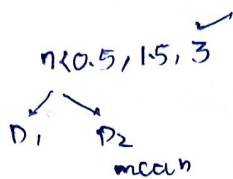
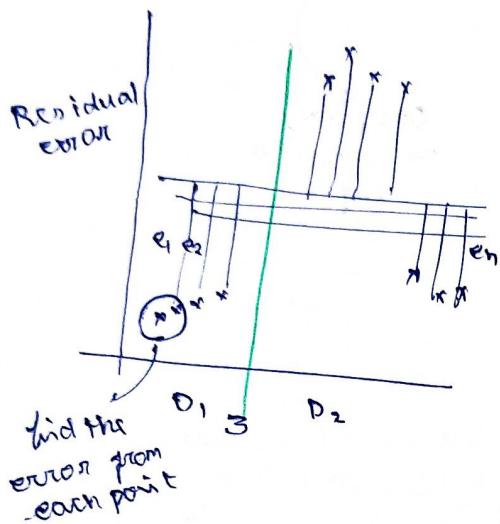


Regression Trees

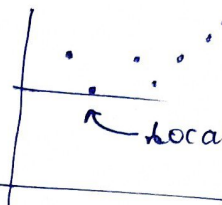
non-linear



How to find splitting points



$$SSE_1 = e_1^2 + e_2^2 + \dots + e_n^2$$



Here we do splitting