

Decision Trees:

KNN
↓
(Instance Based method)

Naive Bayes
↓
(probabilistic method)

Logistic Regress, linear regre, SVM
↓
(geometric, hyperplane)

✧ Decision Trees: similar to if... else condition

DT: nested if... else classifier

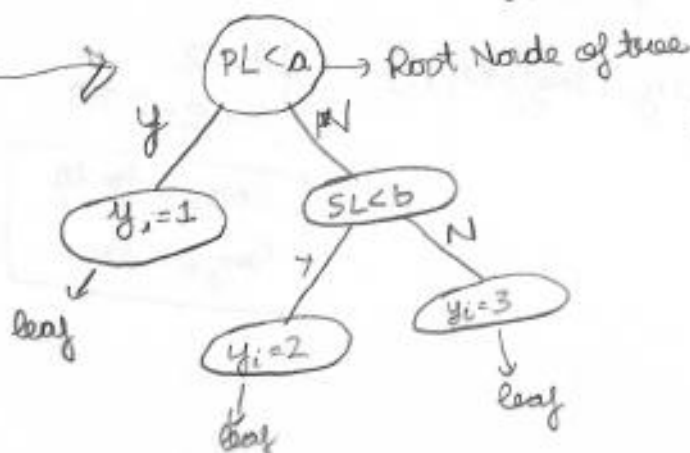
EDA: Iris dataset
↳ $y_i = \{1, 2, 3\}$
(SL, SW, PL, PW)

Sample { if PL < 5
then $y_i = 1$ }

model $x_i < (SL, PL, SW, PW)$
[if PL < a
 $y_i = 1$
else
 [if SL < b
 class = 2
 else
 class = 3]]

nested if... else conditions

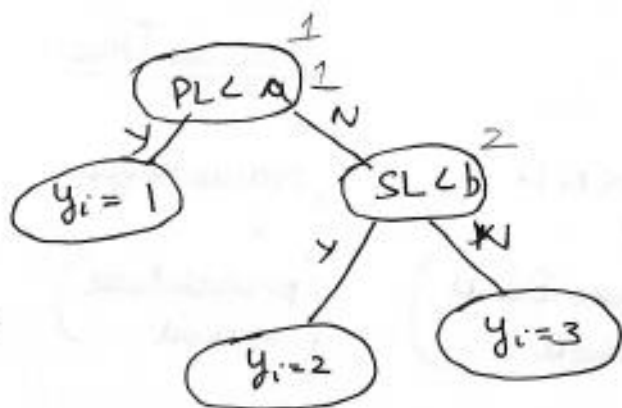
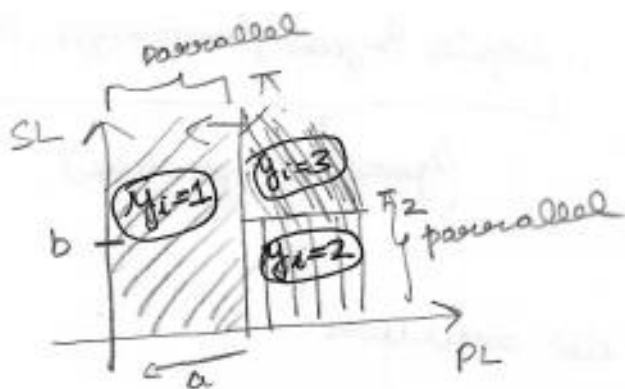
(Tree)



Terminating vertex are called leaf
starting → Root node
Vertex → node →
leaf nodes → decision node
Internal node → neither root nor leaf

non leaf node we make decisions.

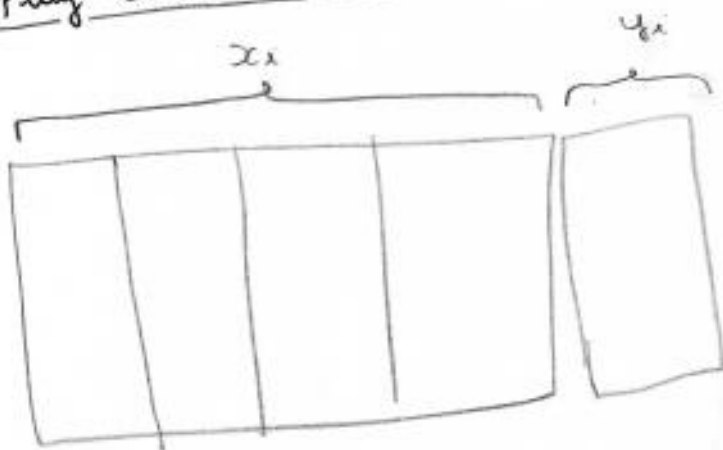
Geometric Intuition :



Corresponding to decision hyperplane is introduced

Note: All of your hyperplanes are axis-parallel in a decision tree.

→ Play Tennis Example



Entropy

Given Random variable $Y \rightarrow y_1, y_2, y_3, \dots, y_k$

$$H(Y) = - \sum_{i=1}^k p(y_i) \log_b(p(y_i))$$

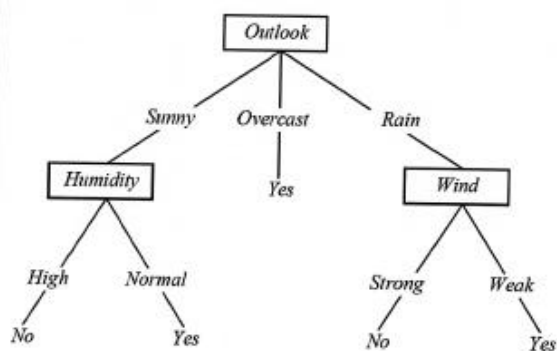
entropy

$$\begin{cases} b=2 \\ b=e=2.718 \end{cases}$$

$$\begin{aligned} \log_2 &= \frac{1}{\log} \\ \log e &= \ln \end{aligned}$$

Play Tennis Example

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No



$$H(Y) = - \sum_{i=1}^K P(y_i) \log(P(y_i))$$

$$P(y_i) = P(Y = y_i)$$

Y : play Tennis

Y_+, Y_- probability $\left\{ \begin{array}{l} P(Y_+) = \frac{9}{14} \\ P(Y_-) = \frac{5}{14} \end{array} \right.$

$$H(Y) = - \sum_{i=1}^K P(y_i) \log(P(y_i))$$

$$H(Y) = - \frac{9}{14} \log\left(\frac{9}{14}\right) - \frac{5}{14} \log\left(\frac{5}{14}\right) = 0.94$$

\uparrow $P(Y_+)$ \uparrow $P(Y_-)$ \rightarrow (% age of -ve pts in D)

$$\frac{\# \text{ +ve pts}}{\text{Total } \# \text{ pts}} = \% \text{ age of +ve pts in D}$$

Properties of Entropy

$Y \rightarrow Y_+, Y_-$ (2 class, 2 category)

Case 1: $Y_+ \rightarrow 99\%$

$D \rightarrow Y_- \rightarrow 1\%$

$$\left. \begin{array}{l} H(Y) = 0.99 \log 0.99 - 0.01 \log 0.01 \\ H(Y) = 0.0801 \end{array} \right\}$$

Case 2:

$D \rightarrow Y_+ \rightarrow 50\%$
 $D \rightarrow Y_- \rightarrow 50\%$

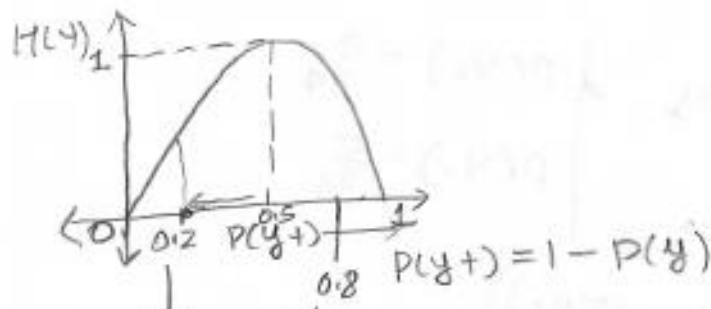
$$\left. \begin{array}{l} H(Y) = 0.5 \log 0.5 - 0.5 \log 0.5 \\ = 1 \end{array} \right\}$$

Case 3:

$D \rightarrow Y_+ \rightarrow 0\%$
 $D \rightarrow Y_- \rightarrow 100\%$

$$\left. \begin{array}{l} H(Y) = 0 \end{array} \right\}$$

NOTE: when both probability are equally probable that
 $H(Y) = 1$ case 2.
 One class fully dominates entropy become 0 case 3.



Same value

$$\begin{cases} 0.2 \leftarrow Y+ \\ 0.8 \leftarrow Y- \end{cases}$$

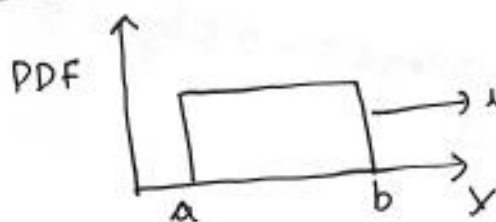
$$\begin{aligned} &= -0.2 \lg 0.2 - 0.8 \lg 0.8 \\ &= -0.8 \lg 0.8 - 0.2 \lg 0.2 \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{both are same.}$$

$$Y \rightarrow y_1, y_2, \dots, y_k$$

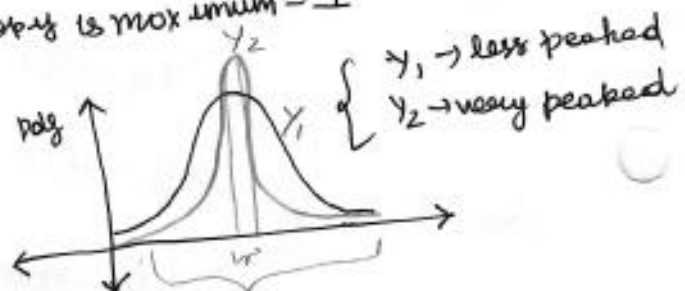
equi probable \rightarrow entropy is maximum

$y_1 \rightarrow$ most probable $\left. \begin{array}{l} \\ y_2, y_3 \dots \rightarrow 0 \end{array} \right\}$ entropy is minimum

PDF / Histogram based



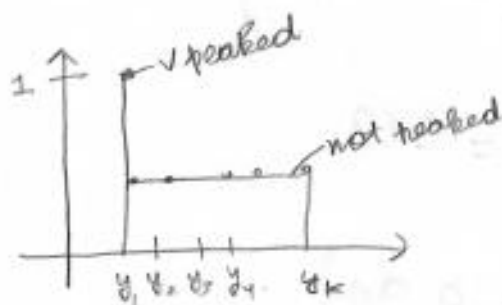
entropy is maximum = 1



$$H(Y_2) < H(Y_1)$$

more peaked the distribution, In case of Y_2 spread is less

On the other hand Y_1 has big spread.



NOTE → The more peaked a distribution is less is its entropy.

↓
If only one value dominates entropy is less
but if points ~~are~~ spread then entropy is maximum

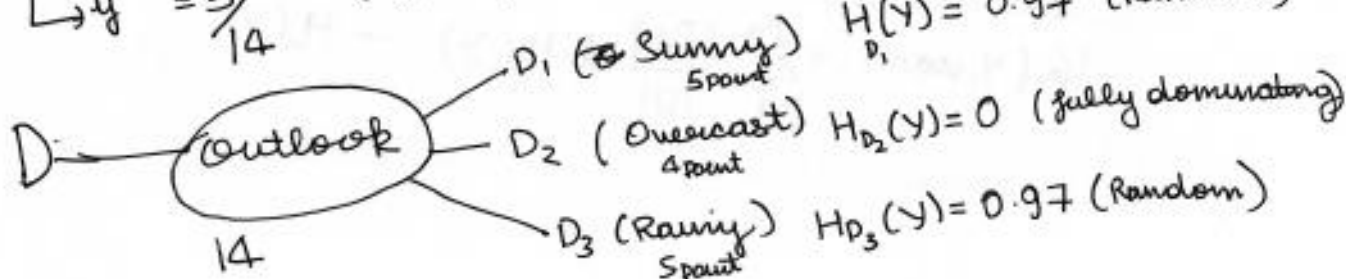
☆ Information gain

$$H(Y) = 0.94$$

$$y \rightarrow y^+ = \frac{9}{14}$$

$$\rightarrow y^- = \frac{5}{14}$$

(9,5) split



Information gain $H_D(Y) = 0.94 \rightarrow$ before breaking dataset
 $H_{D1}(Y) = 0.97, H_{D2}(Y) = 0, H_{D3}(Y) = 0.97$

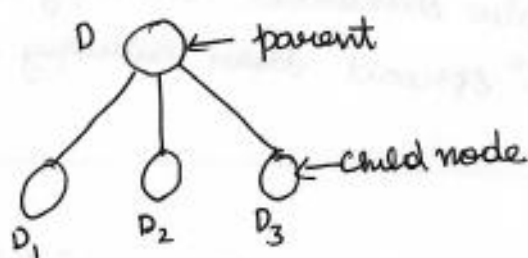
$$IG(Y, outlook) = \left[\left(\frac{5}{14} \times 0.97 \right) + \left(\frac{4}{14} \times 0 \right) + \left(\frac{5}{14} \times 0.97 \right) \right] - 0.94$$

Weighted entropy after D_1, D_2, D_3

$$\begin{aligned} &\left(\frac{5}{14} \times 0.97 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.97 \right) \\ &= \left(\frac{5}{14} \times 0.97 \times 2 \right) \end{aligned}$$

$$= \frac{5}{7} \times 0.97 = 0.69$$

$$IG(Y, outlook) = \left[\frac{5}{7} \times 0.97 - 0.94 \right]$$



$$H_D(Y) = [H_{D_1}(Y)] * \frac{|D_1|}{|D|} + [H_{D_2}(Y)] * \frac{|D_2|}{|D|} + [H_{D_3}(Y)] * \frac{|D_3|}{|D|}$$

$$Y \text{ var} \rightarrow Y_1, Y_2, \dots, Y_k$$

$$IG(Y, var) = \sum_{i=1}^k \frac{|D_i|}{|D|} * H_{D_i}(Y) - H_D(Y)$$

Gini Impurity ~ similar to Entropy

$$\begin{array}{ccc} I_G(Y) & \neq & I_G(Y) \\ \downarrow & & \downarrow \\ \text{Gini Impurity} & & \text{Information Gain} \end{array}$$

$$I_G(Y) = 1 - \sum_{i=1}^k (p(y_i))^2$$

$$Y \rightarrow \begin{cases} y_+ \\ y_- \end{cases}$$

$$Y \rightarrow y_1, y_2, y_3, \dots, y_k$$

Case I: $P(y_+) = 0.5$
 $P(y_-) = 0.5$

$$\begin{aligned} I_G(Y) &= 1 - [(0.5)^2 + (0.5)^2] \\ &= 1 - (0.25 + 0.25) \end{aligned}$$

Gini Entropy $\rightarrow I_G(Y) = 0.5$

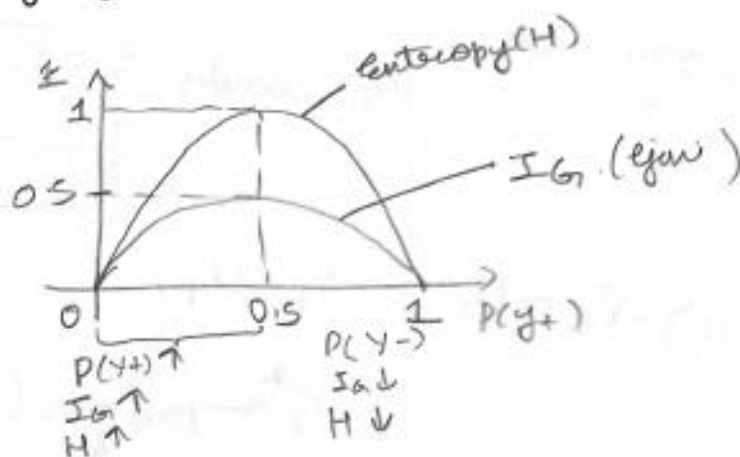
Entropy $\rightarrow H(Y) = 1$

Case II: $P(y_+) = 1$
 $P(y_-) = 0$

Gini Entropy $\rightarrow I_G(Y) = 1 - (1^2 + 0^2) = 0$

Entropy $\rightarrow H(Y) = 0$

2-Category case: y_+, y_- $P(y_+) = 1 - P(y_-)$



$$I_G(y) \text{ (Gini Entropy)}$$

$$1 - (P(y+)^2 + (P(y-))^2)$$

↓
No log

more computationally
efficient to compute
Gini Impurity

$$H(y) \text{ (Entropy)}$$

$$-P(y+) \log P(y+) - P(y-) \log P(y-)$$

↓
log

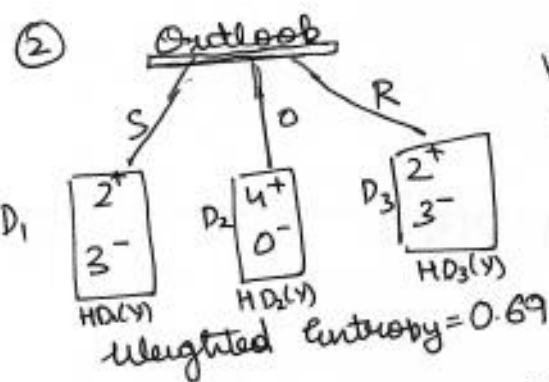
has log ~~less~~ more time complexity

NOTE: Gini Entropy is used over Entropy due to ~~its~~ computational efficiency

BUILDING A DECISION TREE

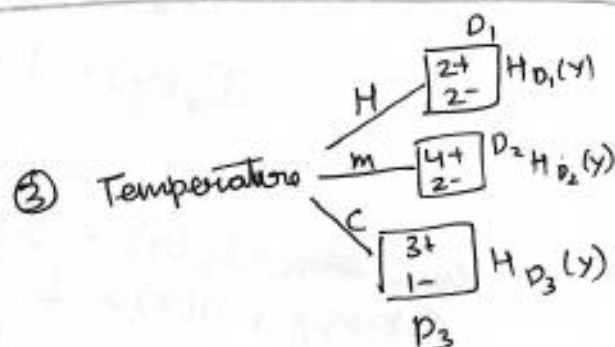
We start with Top to down.

① $D \rightarrow 9+, 5-$
 $H_D(y) = 0.94$



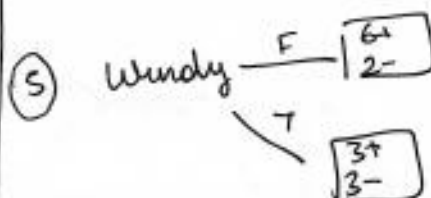
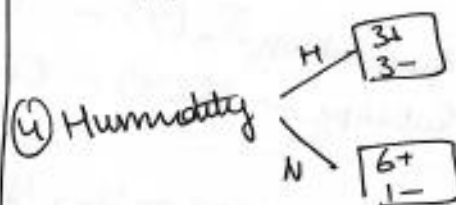
Information Gain = $0.94 - 0.69$
 $I_G(y, \text{Outlook}) = 0.25$

$H_{D1}(y) = 0.47$
 $H_{D2}(y) = 0$
 $H_{D3}(y) = 0.97$



Weighted entropy =

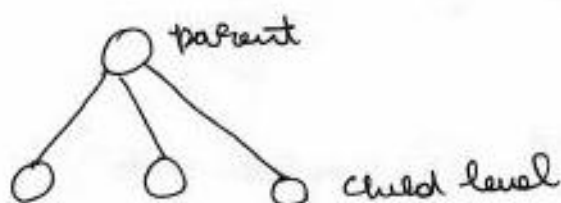
$$\frac{4}{14}$$



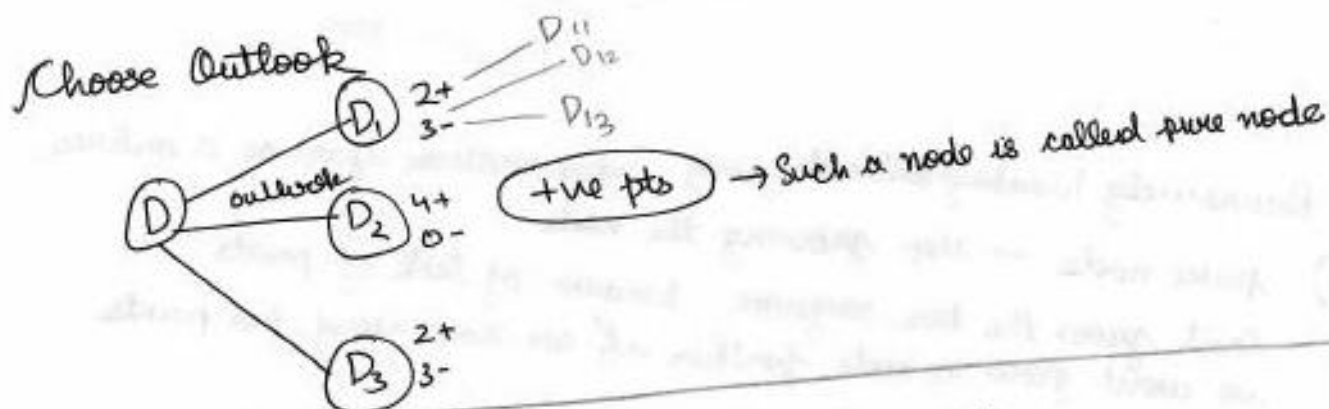
$$D_g(s, t) = g(t) - P_L g(t_b) - P_R g(t_r)$$

$$S^* \leftarrow \arg \min (D_g(s, t))$$

$IG(Y, f) = \text{entropy @ parent-level}$
 $- \text{weighted entropy @ child level}$



$$IG(Y, f) = H_D(Y) - \sum_{i=1}^k \frac{|D_i|}{|D|} * H_{D_i}(Y)$$



$D \rightarrow \text{Tree}$ $IG(Y, f) = H_D(Y) - \sum_{i=1}^k \frac{|D_i|}{|D|} * H_{D_i}(Y)$

$\underbrace{\frac{|D_i|}{|D|}}_{\text{Weight}}$
 $\underbrace{H_{D_i}(Y)}_{\text{Entropy}}$

$IG(Y, \text{outlook}) = 0.25$

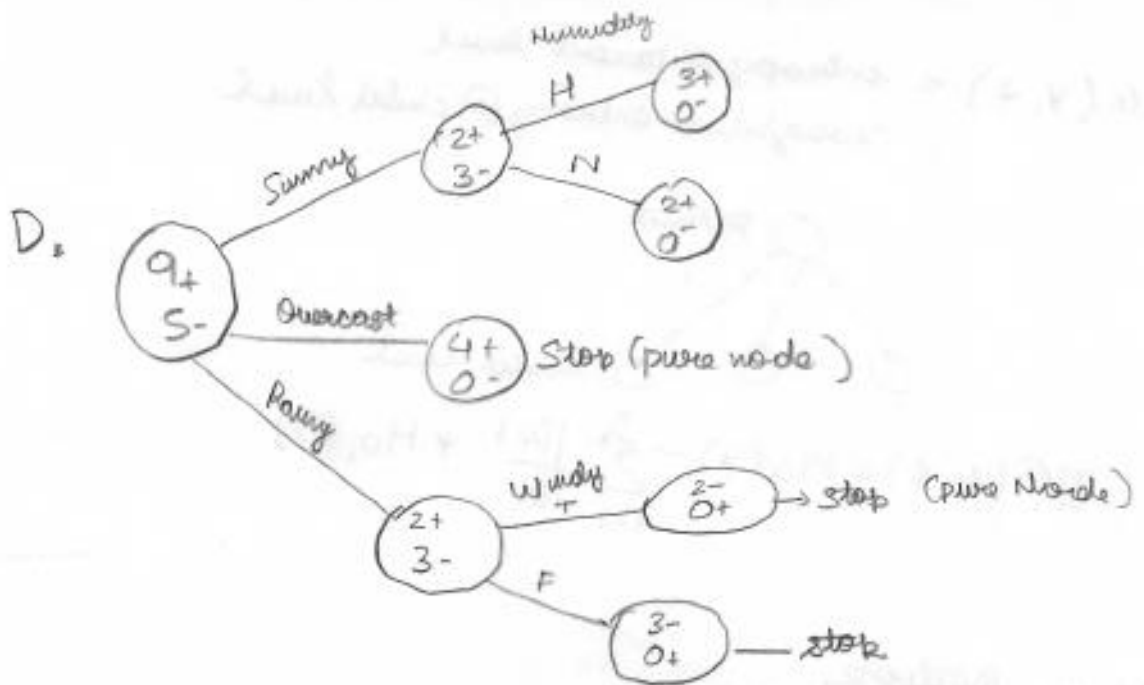
$IG(Y, \text{Temp}) = -$

$IG(Y, \text{Humidity}) =$

$IG(Y, \text{Windy}) =$

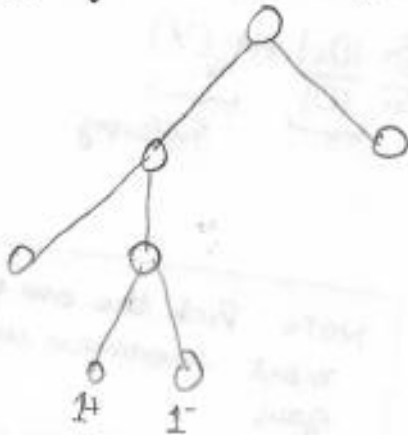
NOTE: Pick the one which has most / maximum information gain

$IG(Y, f) = (\text{entropy @ parent-level}) - (\text{weighted entropy @ child level})$

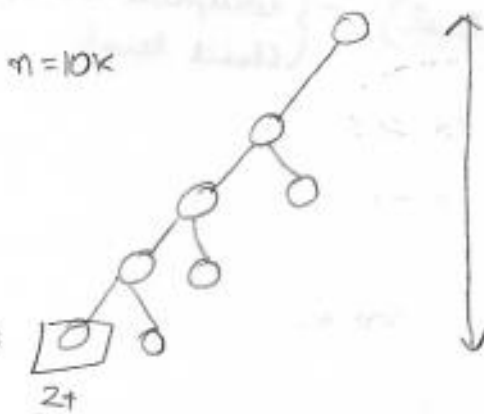


Recursively breaking each node using Information Gain as a criteria

- ① pure node \rightarrow stop growing the node
- ② Can't grow the tree anymore because of lack of points
we won't grow a node further if we have very few points



- ③ if we are too deep we stop growing tree.



Depth of the tree \uparrow ; overfitting \uparrow (few pts)

-6-

Depth of tree is small; underfitting

Too prevent overfitting and underfitting

Decision Tree \rightarrow hyperparameter \rightarrow depth - CV (Cross Validation)

\rightarrow logistic reg.
 \rightarrow SVM.

Splitting numerical features

Construct a DT: Splitting a node \rightarrow Information Gain

$|G| \rightarrow$ entropy

\rightarrow Gini impurity \rightarrow computationally efficient

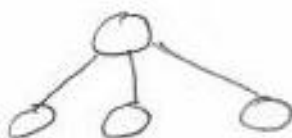
f_1	y
2.2	1
2.6	1
3.5	0
3.8	0
4.6	1
5.3	0

f_1 : numerical

\rightarrow Integer
 \rightarrow real valued

Split based on categorical variables

f_2 : 3-categories

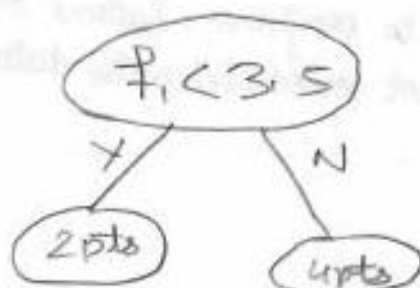


① Sort the numerical feature (in ascending order)

$f_1 < 2.2$ (one possible condition)

at most
n possible
conditions

$\left\{ \begin{array}{l} f_1 < 2.2 \\ f_1 < 2.6 \\ f_1 < 3.5 \\ f_1 < 3.8 \\ f_1 < 4.6 \\ f_1 < 5.3 \end{array} \right.$



We calculate information gain for each n and compute maximum information gain

num f_1	caten f_2

Feature Standardization

	f_j
x_i	x_{ij}

$$x'_{ij} = \frac{x_{ij} - \mu_j}{\sigma_j}$$

μ_j } mean and
 σ_j } standard deviation

In case of decision tree \rightarrow Not a distance base method

f_j

Threshold



\rightarrow sort this column

(In decision tree it is dependent on order)

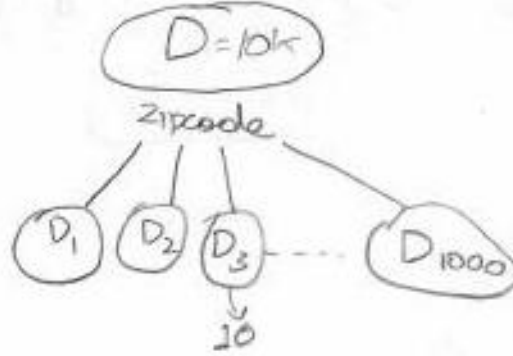
Note

No need to perform feature standardization in decision tree as it is not concerned with distance but it with order.

Categorical features with many categories

pincode / zipcode \rightarrow 1000's

Zipcode	y
pincode	
—	
—	
—	
—	



hack (feature engineering)

Pincode $y_i = \{0, 1\}$

Pincode/Zipcode \rightarrow Categorical

convert it

numerical feature

$$P(y_i=1 | P_j)$$

For example $P(y_i=1 | P_j)$

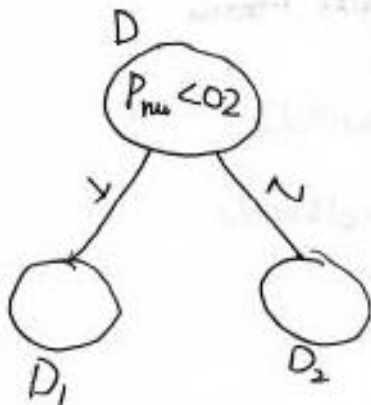
20 times P_j occurs

19 times $y_i=1$

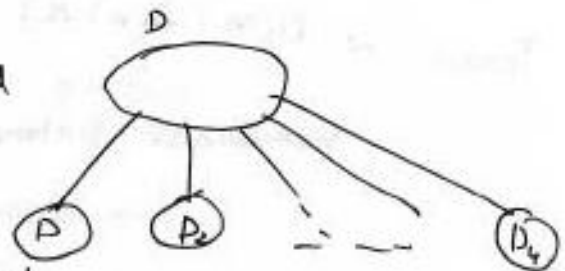
$$P(y_i=1 | P_j) = \frac{19}{20}$$

After it

Conversion from Categorical to numerical feature using conditional probabilities



if not used
It is presented if used probability



Overfitting und underfitting

depth $\uparrow \Rightarrow$ possibility of having very few pts @ a leaf node \uparrow

→ Interpretability of model also decreases if $C > C'$ & $C > C''$

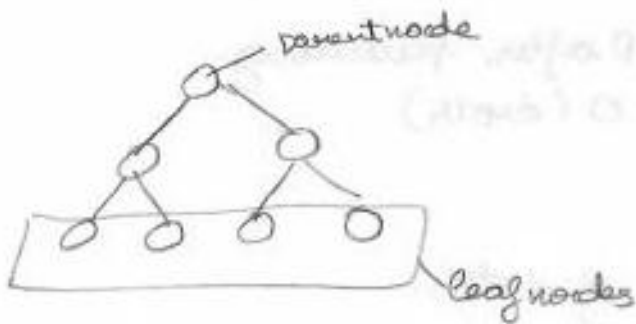
$$x_a \rightarrow y_a$$

Also ~~T~~
After Training:-

a) Runtime Space :- store my D-Tree

$x_1 \rightarrow y_2$

if else } nested if else



\uparrow # internal-nodes
+
leaf nodes.

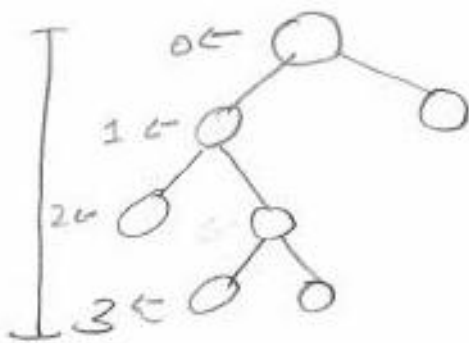
Space complexity $O(\#nodes)$

depth = 5 or 10

depth $\uparrow \Rightarrow$ interpretability \downarrow

Run time space :- reasonable

Run time complexity $x_2 \rightarrow y_2$



$K = \max$ - depth of any leaf nodes
at most 3 comparisons.
 $O(3)$.

DT :- Super good when you have large data.
dimensionality is small
low latency $\rightarrow O(\text{depth})$

Regression using DT

Pros of Decision Trees

- Interpretable : humans can understand decisions
- easily handles irrelevant attributes (gain = 0)
- can handle missing data
- very compact : # nodes $\ll D$ after pruning
- very fast at testing time : $O(\text{depth})$

Cons :

- only axis aligned splits of data
- exponentially many possible trees

Example Decision tree
Book Data mining pag 338

-9-

Examp 8.1

$$\text{Info}(D) = -\frac{9}{14} \log_2 \left(\frac{9}{14} \right) - \frac{5}{14} \log_2 \left(\frac{5}{14} \right) = 0.940 \text{ bits}$$

$$\begin{aligned} \text{Info}_{\text{age}}(D) &= \left[\frac{5}{14} \times \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) \right. \\ &\quad \left. + \frac{4}{14} \left(-\frac{4}{4} \log_2 \frac{4}{4} \right) + \frac{5}{14} \times \left(-\frac{3}{5} \log_2 \frac{3}{5} - 2 \log_2 \frac{2}{5} \right) \right] \\ &= 0.694 \text{ bits} \end{aligned}$$

Hence the Information Gain

$$\begin{aligned} \text{Gain}(\text{age}) &= \text{Info}(D) - \text{Info}_{\text{age}}(D) \\ &= 0.940 - 0.694 \\ &= 0.246 \text{ bits} \end{aligned}$$

Similarly,

$$\text{Info Gain}(\text{Income}) = 0.029 \text{ bits}$$

$$\text{Info Gain}(\text{Students}) = 0.151 \text{ bits}$$

$$\text{Info Gain}(\text{credit Rating}) = 0.048$$

Info. gain for age is maximum so it is selected

