

Decision Trees

30.1 Geometric intuition of Decision trees

known : is called Instance based method

NB : Probabilistic methods

Log : geometric in nature

Linear

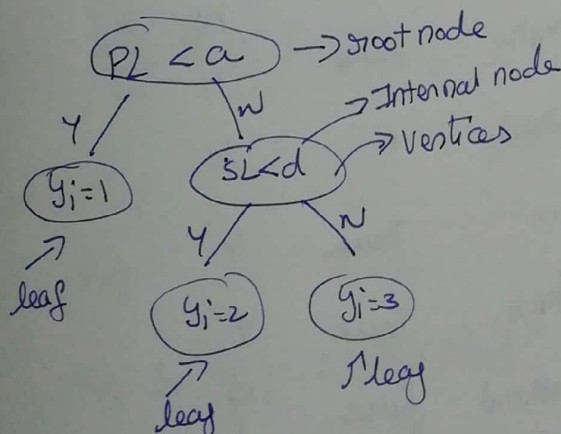
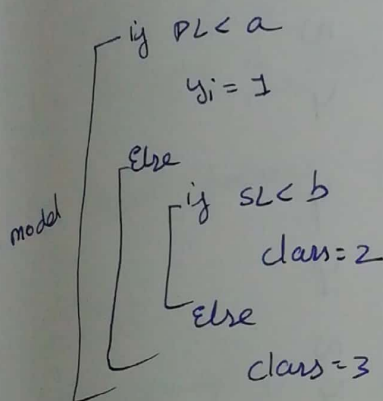
Sum

↳ Kernel trick

→ DT is ~~not~~ Programmers think like (if... else condition

→ DT is nested if-else classification

$$x_i = \langle SL, PL, SW, PL \rangle$$



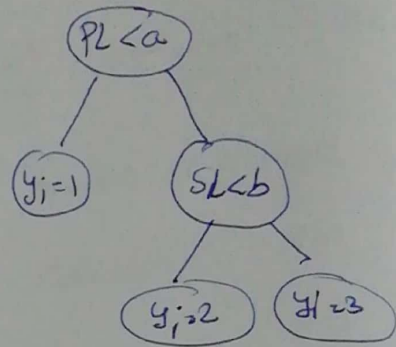
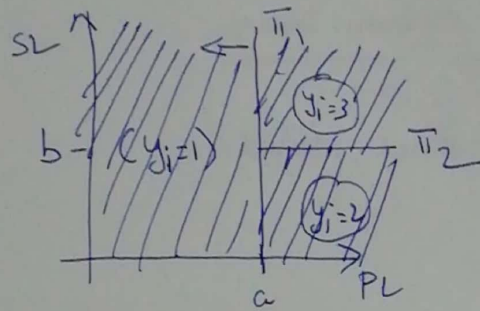
root node :

leaf node

internal nodes

At all non leaf nodes we will take decisions

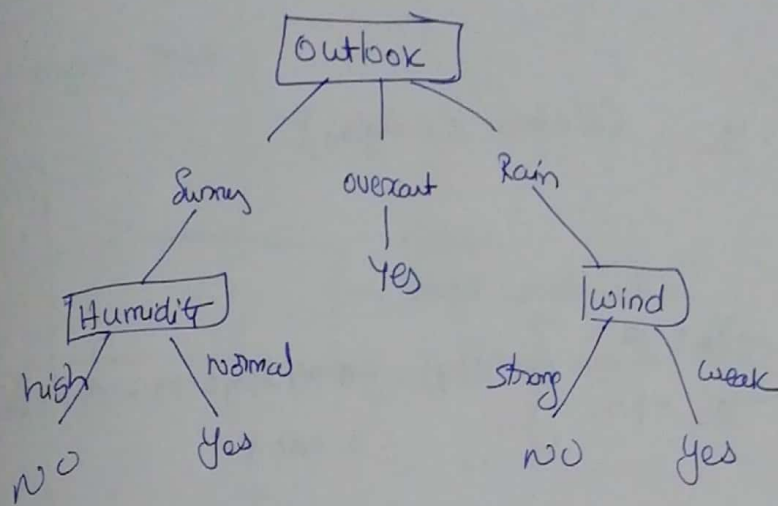
Geometric Intuition



- All of your hyperplanes are axis parallel
- DT: a set of axis parallel hyperplanes

30.2 Sample decision tree

outlook	temperature	humidity	windy	playtennis
Sunny	hot	high	false	N
Sunny	"	"	True	N
overcast	"	"	false	Y
Rainy	mild	"	"	Y
"	cool	normal	False	Y
"	cool	"	True	N
Overcast	"	"	True	Y
"	"	"	"	N
Sunny	mild	high	false	Y
Rainy	"	normal	"	Y
Sunny	"	"	True	Y
overcast	mild	high	True	Y
"	hot	normal	false	Y
Rainy	mild	high	True	N



Ex: If outlook = sunny
 If humidity = high
 $y_i = \text{No}$

$x_q = [\text{sunny}, \text{hot}, \text{high}, \text{r}] ; y_q = \text{No}$

30.3 Building a decision Tree : Entropy

Biggest challenge is building a decision tree from $\mathcal{D}_{\text{train}}$.

Entropy

i.e. $Y \Rightarrow y_1, y_2, y_3, \dots, y_K$ [Ex: Playtennis \rightarrow Yes, No]

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

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

$$P(y_{\text{Yes}}) = 9/14$$

$$P(y_{\text{No}}) = 5/14$$

$$\begin{cases} \log_2 = \log \\ \log_e = \ln \end{cases}$$

$$= -9/14 \log(9/14) - 5/14 \log(5/14)$$

$$= 0.94$$

% age of splits in \mathcal{D}

$$\frac{\text{\# the pts}}{\text{total \#pts}} = \% \text{ age of the pts in } \mathcal{D}$$

Properties

$Y \rightarrow y_+, y_-$ (2 class, 2 categories)

Case 1:

$$\left. \begin{array}{l} y_+ = 99\% \\ y_- = 1\% \end{array} \right\} H(Y) = -0.99 \log 0.99 - 0.01 \log(0.01) \\ = 0.0801$$

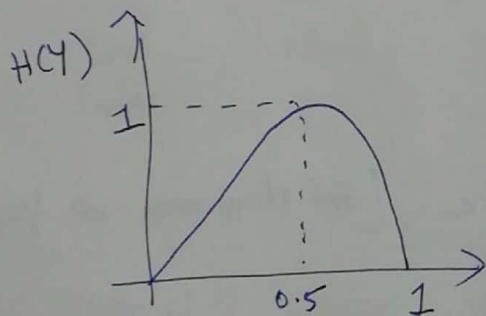
Case 2:

$$\left. \begin{array}{l} y_+ = 50\% \\ y_- = 50\% \end{array} \right\} H(Y) = -0.5 \log 0.5 - 0.5 \log 0.5 \\ = 1$$

Case 3:

$$\left. \begin{array}{l} y_+ = 0\% \\ y_- = 100\% \end{array} \right\} H(Y) = 0 \quad \begin{array}{l} 1 \log(1) = 0 \\ 0 \log(0) = 0 \end{array}$$

$$\boxed{\therefore P(y_+) = 1 - P(y_-)}$$



Ex

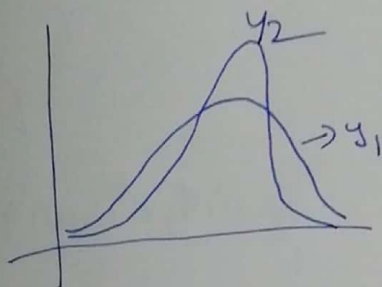
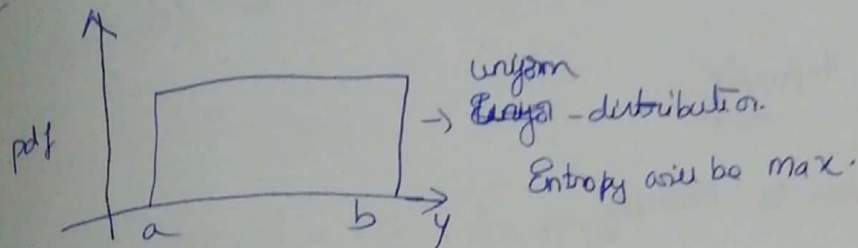
Given a r.v Y

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

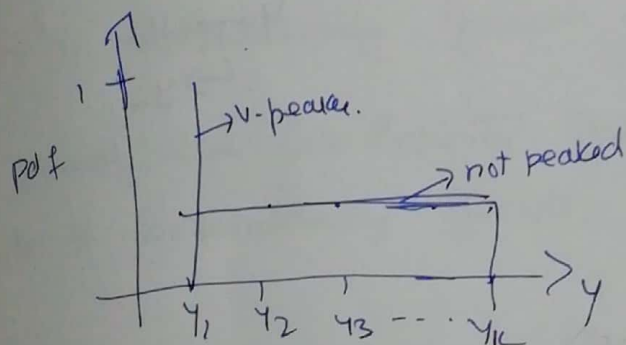
Equi probable \rightarrow Entropy is maximum

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

PDF for decision Tree



$$H(y_2) < H(y_1)$$



(discrete case)

- \rightarrow more peaked a distribution, less is its Entropy
- \rightarrow If one class dominates Entropy is minimal

30.4 Building a decision Tree \div Information gain

32 - decision-tree.pdf

homepage.cs.wri.edu/faculty/hamel/Courses/2016/spring2016/cs581/lecture-notes/

$$\text{Information Gain}(Y, \text{outlook}) = \left(\frac{5}{14} * 0.97 \right) + \left(\frac{4}{14} * 0 \right) + \left(\frac{5}{14} * 0.97 \right) -$$

(0.97)

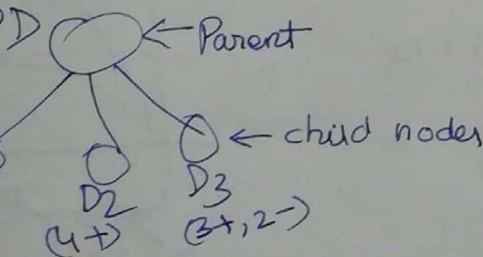
$$H(y) = -\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} \log_2 \left(\frac{2}{5} \right)$$

$$= -0.7369 - 0.463974$$

$$= 0.603025$$

$$\approx 0.97$$

$H_D(y)$



$$\frac{D_1}{D} H_{D_1}(Y) + \frac{D_2}{D} H_{D_2}(Y) + \frac{D_3}{D} H_{D_3}(Y)$$

total # of points in D

total # of points in D_3
Entropy of D_3

$$Y \xrightarrow{\text{Var}} \begin{matrix} x_1 & x_2 & & x_k \\ y_1 & y_2 & \dots & y_k \end{matrix}$$

$$I_G(Y, \text{Var}) = \sum_{i=1}^K \frac{|x_i|}{|X|} \cdot H_{D_i}(Y) - H_D(Y)$$

30.5 Gini Impurity

→ It is similar to Entropy

$I_G(Y) \rightarrow$ Gini Impurity of Y

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

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

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

Case 1:

$$P(y_+) = 0.5$$

$$P(y_-) = 0.5$$

$$I_G(Y) = 1 - (0.25 + 0.25)$$

$$= 0.5$$

$$H(Y) = 1 \rightarrow \text{Entropy}$$

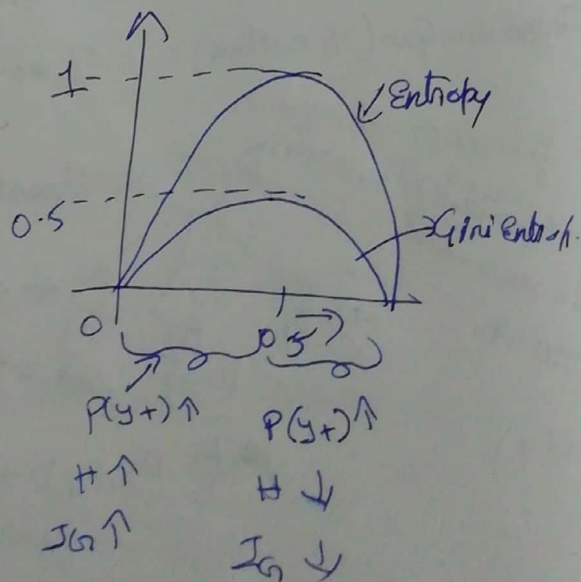
Case 2:

$$P(y_+) = 1$$

$$P(y_-) = 0$$

$$I_G(Y) = 1 - (1 + 0) = 0$$

$$H(Y) = 0$$



As both Entropy & Gini are reducing & increasing w.r.t to $p(y+)$. Then why we need two.

$$H(y)$$

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

$$I_G(y)$$

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

→ Gini is Computationally Efficient to Calculate than Entropy

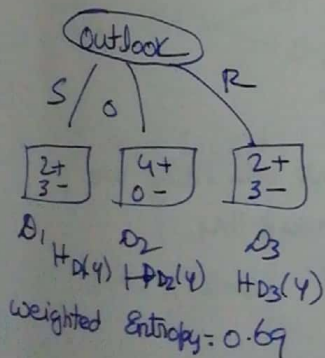
→ ppl generally use Gini Impurity than Entropy because it's faster & Computationally Efficient.

30.6 Building a decision Tree : Construction a DT

$$\textcircled{1} \textcircled{0} \rightarrow 9+, 5-$$

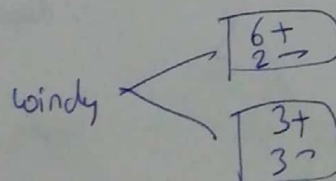
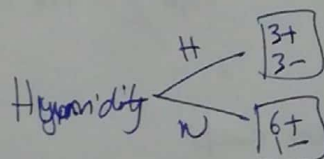
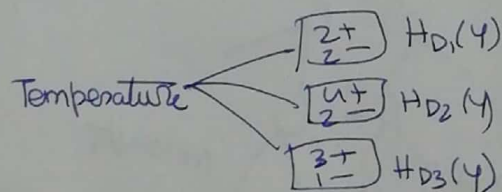
$$H_D(y) = 0.94$$

Step 1

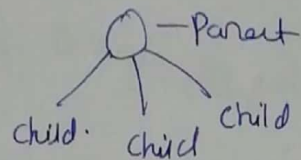


$$I.G. = 0.94 - 0.69$$

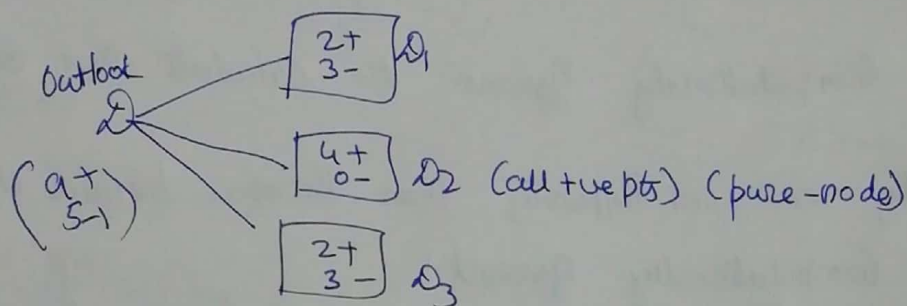
$$I_G = 0.25$$



$$IG(y, f) = H_D(y) - \sum_{i=1}^K \frac{|D_i|}{|D|} \cdot H_{D_i}(y)$$

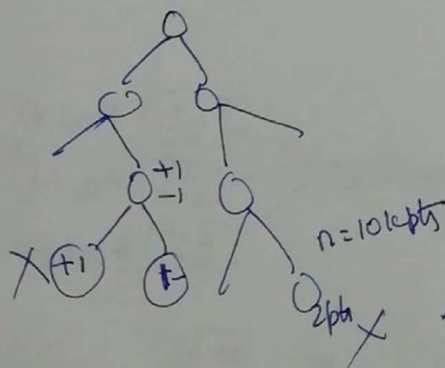
$$IG(y, f) = \text{Entropy @ parent-level} - \text{weighted entropy @ child level}$$


-) we have to pick no variable which has maximum information gain.



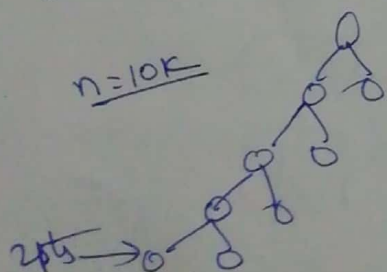
Eventually we will get 3

- ① Pure node \rightarrow stop growing the node
- ② Can't grow the tree any more because of lack of PT



As there too points are noisy we are overfitting

- ③ If we are too deep



\therefore depth of the tree \uparrow ; overfitting \uparrow
(few pts)

depth is small \rightarrow underfit

DT: hyperparameters \div depth

30.7 Splitting numerical features

Most important step in Splitting a node, For this we need IG

we typically use Entropy, Instead of we can use Gini impurity

Gini is 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

↳ Integer

↳ real valued.

Step 1 : Sort the numerical feature in ascending order

Step 2 :

$f_1 < 2.2$	$f_1 < 4.6$
$f_1 < 2.6$	$f_1 < 5.3$
$f_1 < 3.5$	
$f_1 < 3.8$	

$f_1 < \tau_1 \begin{cases} D_1 \\ D_2 \end{cases}$

$f_1 < \tau_2 \begin{cases} D_1 \\ D_2 \end{cases}$

⋮

$f_1 < \tau_n \begin{cases} D_1 \\ D_2 \end{cases}$

For Every value we will calculate IG,
we will pick the max IG

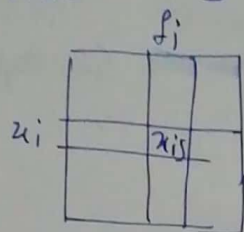
f_1 (num) f_2 (cat)

IG	IG
----	----

we will choose the feature with max IG.

30.8 Feature Standardization

Logistic Regression }
SVM } feature standardization
KNN }



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

Decision Trees: not a distance based method

↳ It depends on the order of the columns.

→ do not need to do feature standardization/normalization.

30.9 Building a decision Tree: Categorical features with many values

→ If ^{levels} no. of classes in a feature is more

Ex: pincode

↳ 1000's

~~Zipcode~~

Feature Engineering hack

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

P1

P2

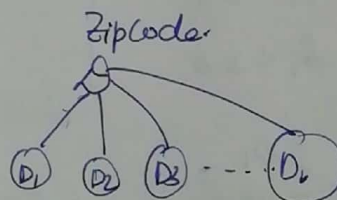
P3

P4

P1

P2

⋮

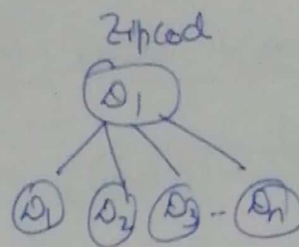
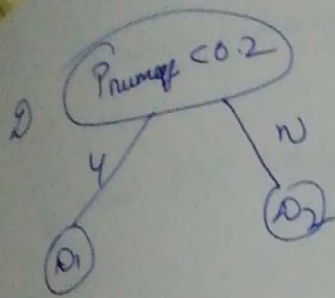


The splits/datasets will become extremely small.

$$P(y_i = 1 | P_j) = 1/20$$

$$= \frac{\# \text{ of times } y_i = 1 \text{ \& } P_j}{\# P_j}$$

We will remove the Zipcode from the data and replace with numeric feature



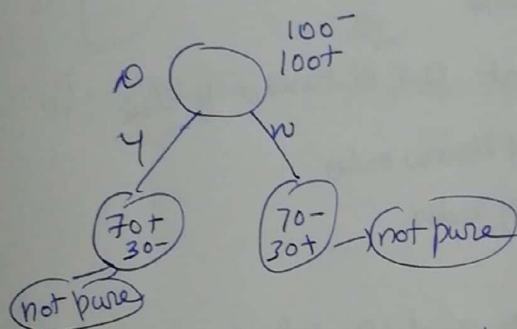
→ we will get rid of the problem called data sparsity

30.10 overfitting & underfitting

→ As the depth \uparrow , The possibility of having very few pts @ leaf nodes increases

→ Interpretability of the model \downarrow when depth \uparrow

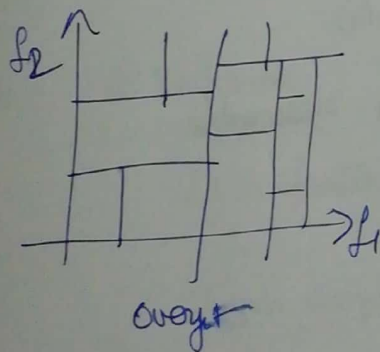
→ when depth is low, underfit the model.



→ A tree at depth 1 is called decision stump.

→ depth should be discovered using Cross Validation.

→ we depth is large, the plane is divided into multiple sections



30.11 Train & Runtime Complexity

Train $\div O(n(\log n)d)$ $n = \# \text{ pts to train}$
 $d = \text{dim.}$

numeric features \div threshold

\hookrightarrow algorithmic method.

$n \log n \rightarrow$ sorting

If we have large dimensions, decision tree is not the best model.

After Training

Runtime Space \div storing DT is easy
 $\hookrightarrow O(\text{nodes})$ \hookrightarrow If else

people convert D.T in nested if else

$\hookrightarrow \#$ of internal nodes

$\hookrightarrow \#$ leaf nodes

\rightarrow Typically we won't DTs with a depth more than 5 or 10

depth $\uparrow \rightarrow$ Interpretability \downarrow

\rightarrow runtime Space Complexity is reasonable

\rightarrow runtime Time Complexity $O(\text{depth})$

\hookrightarrow is very very reasonable

\rightarrow DT is good

\hookrightarrow large data

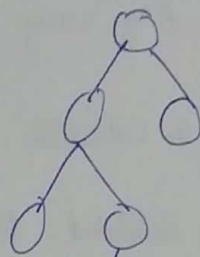
\hookrightarrow dim is small / reasonable

\hookrightarrow low latency $\rightarrow O(\text{depth})$

30.12 Regression using Decision Trees

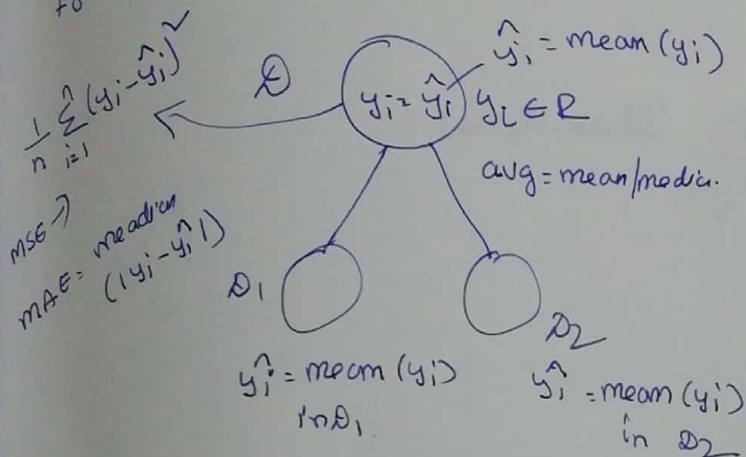
IG is used for classification algorithm

for classification



$n_1: +ve$
 $n_2: -ve$
 If $n_1 > n_2$
 $y_q = +ve$

For regression we will use MSE & MAE



30.13 Cases

Imbalanced data \div we have to balance the data

- \rightarrow up-sampling, down-sampling
- \rightarrow class weight

\rightarrow Impacts Entropy & MSE calculations.

large data \div @ each node, split

\downarrow
Each feature IG

\downarrow
Train complexity to train DT increases.

we should avoid \rightarrow one hot encoding

categorical feature with lots of levels \rightarrow It is unequal to convert to numeric features using $P(y_i = 1 | f = c_i)$

Similarity metric : DT need no features explicitly

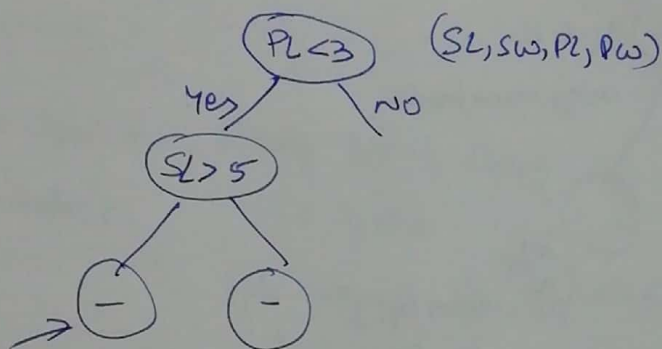
multiclass classification : Any way we have "one versus rest"

I_G & Entropy can be calculated even if we have y_1, y_2, \dots, y_K

DT naturally can be extended to multiclass classification.

Decision regions : non-linear
axis-parallel hypercuboids

Feature interaction :



$(PL < 3) \cap (SL > 5)$

↳ This is called Feature Interaction

Outlier

depth \uparrow : Outlier will impact

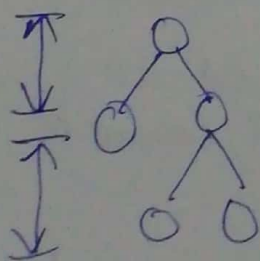
↳ Tree will become unstable

Interpretable

↳ Every time if else, it is super interpretable

Feature Importance

For every feature f_i we will get reduction in H or I_G



$\left\{ \begin{array}{l} \rightarrow \text{normalized} \\ \text{Sum up reduction in } H \text{ due to } f_i \end{array} \right.$