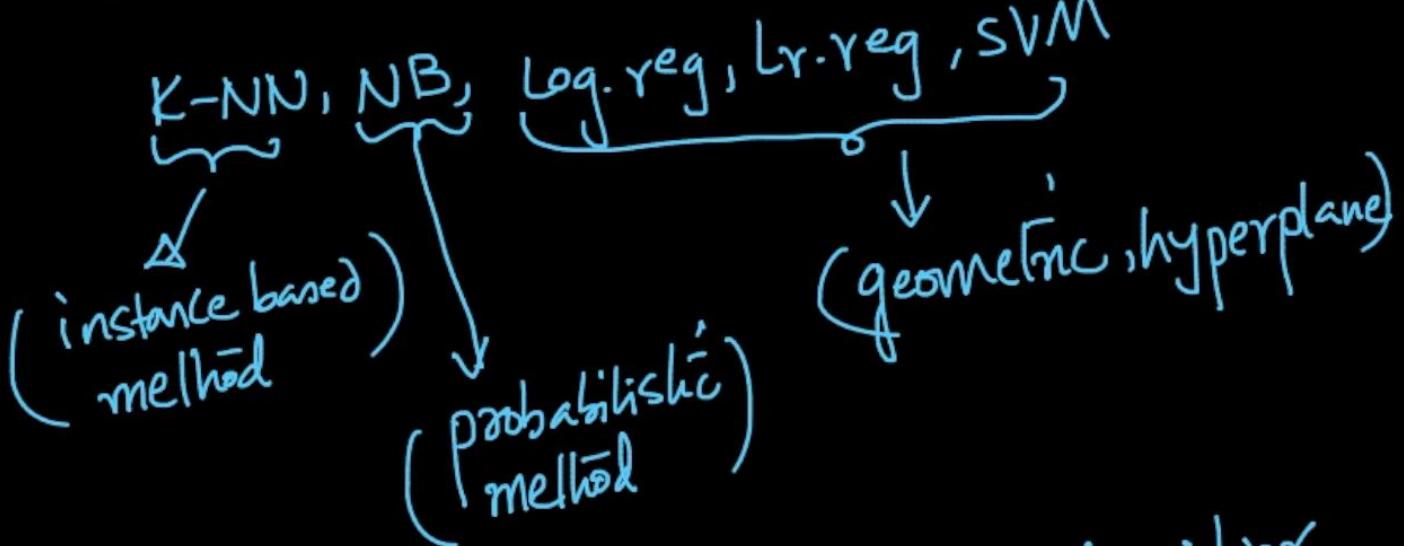


✓ Decision Trees: (DT) → (if ... else)



(DT): nested if...else classifier



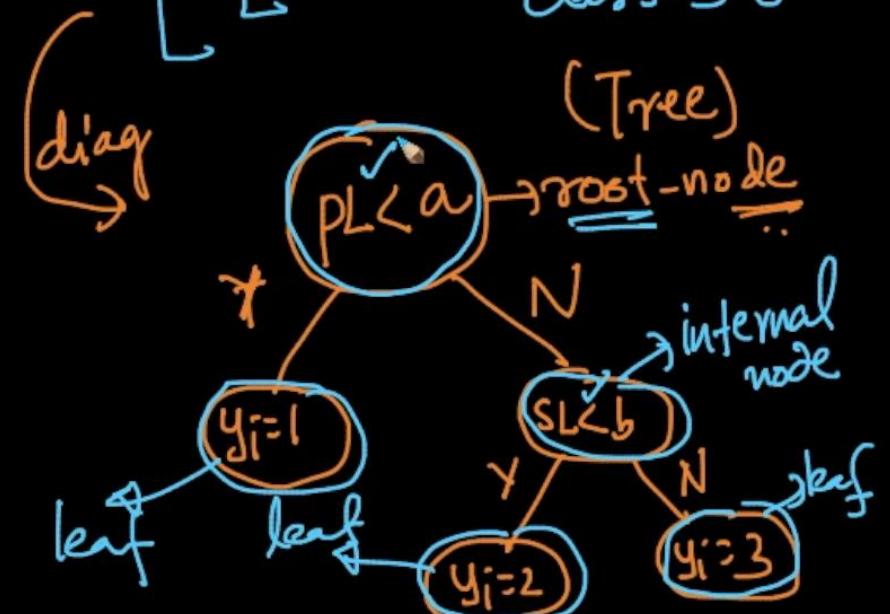
EDA: Iris dataset  
 $\hookrightarrow y_i = \{1, 2, 3\}$   
 $(SL, SW, PL, PW)$

simple  $\left\{ \begin{array}{l} \text{if } PL < 5 \\ \text{then } y_i = 1 \end{array} \right\}$

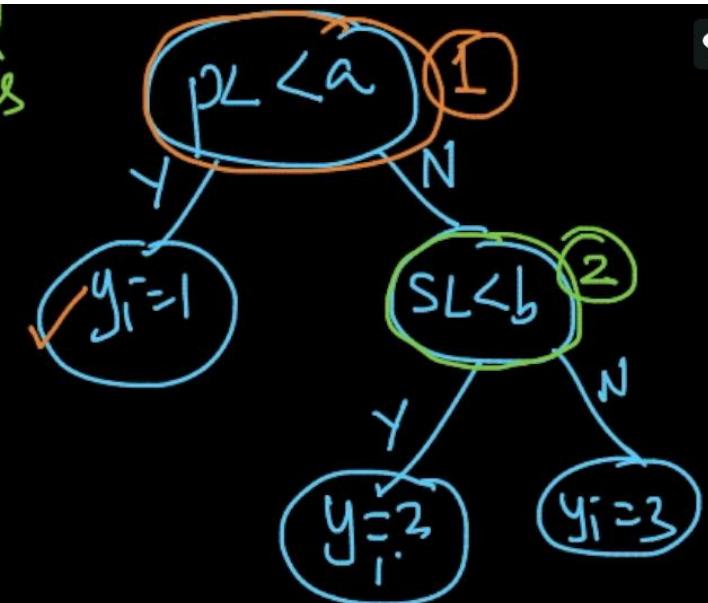
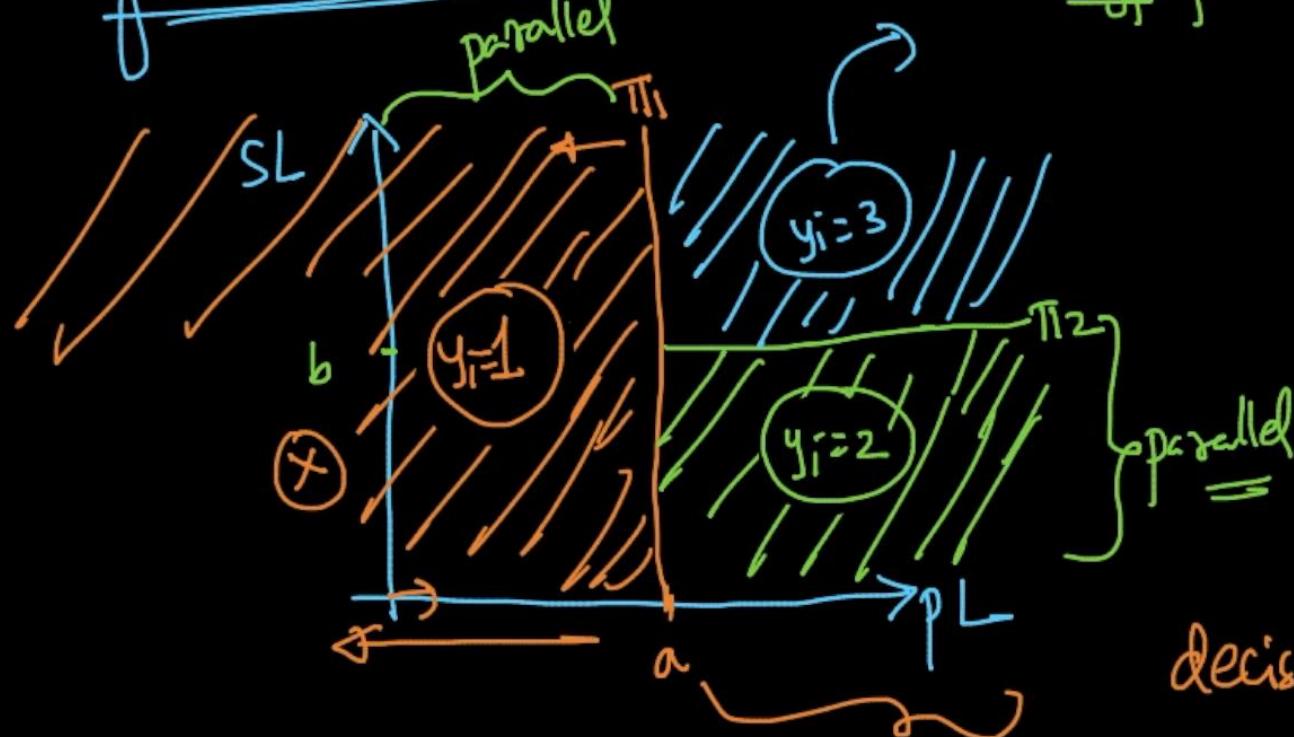
✓ nested if .. else conditions

↓  
diagram (Tree)

$x_i = \langle SL, PL, SW, PW \rangle$   
 $= \left\{ \begin{array}{l} \text{if } PL < a \\ \quad y_i = 1 \\ \text{else} \quad \left\{ \begin{array}{l} \text{if } SL < b \\ \quad \text{class} = 2 \\ \text{else} \quad \text{class} = 3 \end{array} \right. \end{array} \right. \right\}$

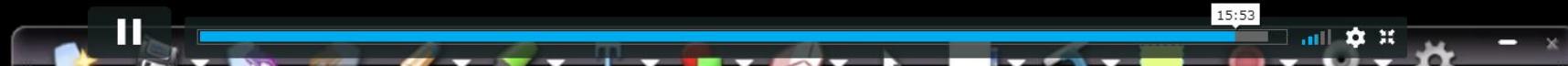


geom-intuition:  $DT$  :- set of axis parallel hyperplanes



decision :-  $\pi$

✓ all of your hyperplanes are axis-parallel



node - vertex

root-node

leaf-node

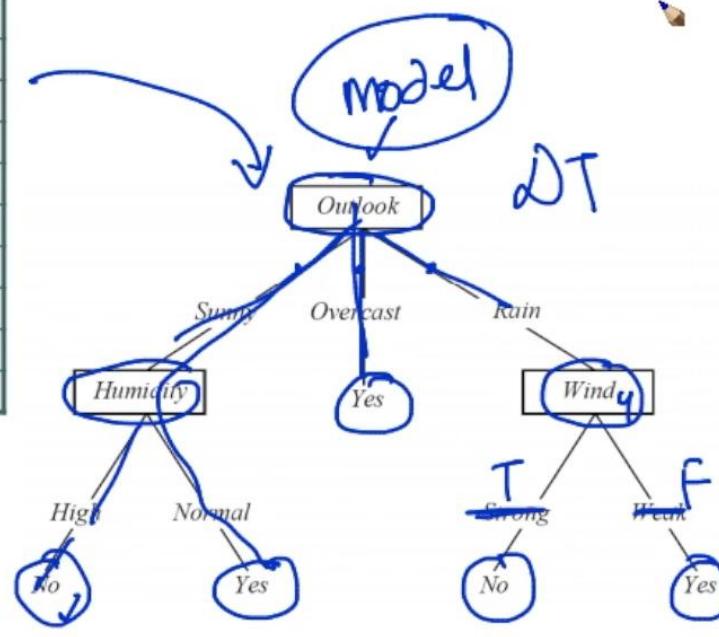
internal-node

non-leaf nodes:- decisions  
in a DT

# 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

- if outlook = Sunny  
if humidity = High  
 $y_i = \text{No}$

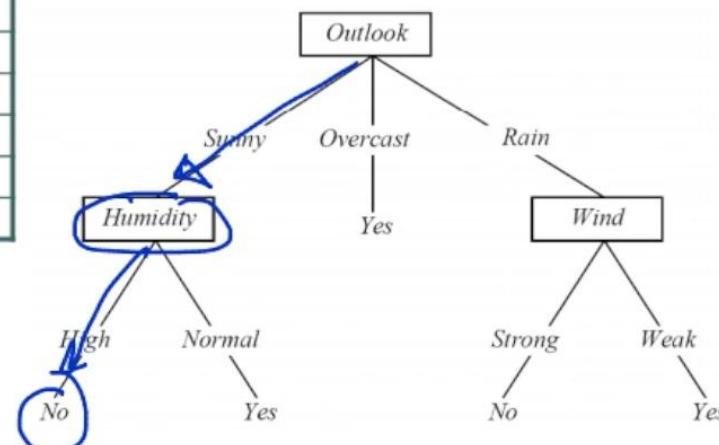


# Play Tennis Example

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

$x_4 = [\cancel{\text{Sunny}}, \cancel{\text{hot}}, \cancel{\text{high}}, \text{T}, \cancel{\text{X}}]$

$y_4 = \text{NO}$



✓ Building a DT: Entropy → information theory  
→ physics, info thy

✓  $D_{\text{Train}} \rightarrow DT$



# Play Tennis Example

$\text{Y}$

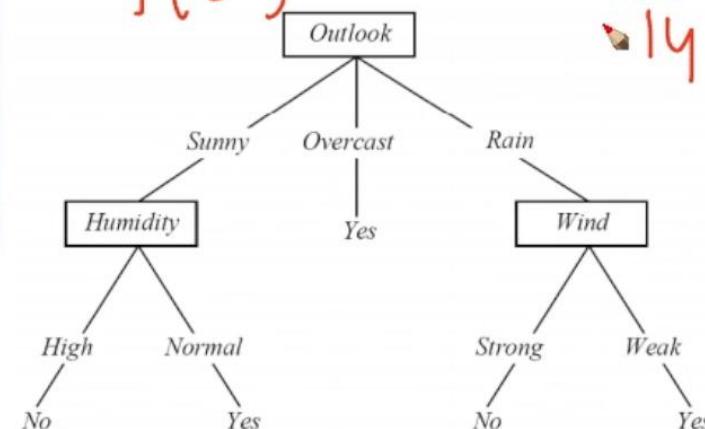
Outlook	Temperature	Humidity	Windy	Play Tennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal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

2 ↗

$\text{Y} \rightarrow \text{play Tennis}$

$$\text{Y}_+, \text{Y}_- \left\{ P(\text{Y}_+) = \frac{9}{14} \right\}$$

$$P(\text{Y}_-) = 1 - P(\text{Y}_+) = \frac{5}{14}$$



$y, v \quad Y \rightarrow y_1, y_2, y_3, \dots, y_k$

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

$$p(y_i) =$$

$$\underline{p(y_i) = P(Y=y_i)}$$

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

$$\boxed{\log_2 = \lg}$$
$$\boxed{\log_e = \ln}$$



$$H(Y) = - \sum_{i=1}^k p(y_i) \log_2(p(y_i))$$

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

✓

$\frac{\# \text{ +ve pls}}{\text{Total } \# \text{ pls}}$

$p(y_+)$

$p(y_-)$

$\% \text{ age of +ve pls in D}$



$$H(Y) = - \sum_{i=1}^k p(y_i) \log_2(p(y_i))$$

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

$\downarrow$

$p(y_+)$

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

$p(y_-)$

% age of -ve pls in D



Properties:  $\text{Y} \rightarrow Y_+, Y_-$  (2 class, 2 category)

Case 1:  $\begin{cases} Y_+ \rightarrow 99\% \\ Y_- \rightarrow 1\% \end{cases}$   $H(Y) = -0.99 \lg 0.99 - 0.01 \lg 0.01$   
 $= 0.0801$

Case 2:  $\begin{cases} Y_+ \rightarrow 50\% \\ Y_- \rightarrow 50\% \end{cases}$   $H(Y) = -0.5 \lg 0.5 - 0.5 \lg 0.5$   
 $= 1$

Case 3:  $\begin{cases} Y_+ \rightarrow 0\% \\ Y_- \rightarrow 100\% \end{cases}$   $H(Y) = 0$



Properties:  $\text{Y} \rightarrow Y_+, Y_-$  (2 class, 2 category)

Case 1:  $\begin{cases} Y_+ \rightarrow 99\% \\ Y_- \rightarrow 1\% \end{cases}$

$$H(Y) = -0.99 \lg 0.99 - 0.01 \lg 0.01 = 0.0801$$

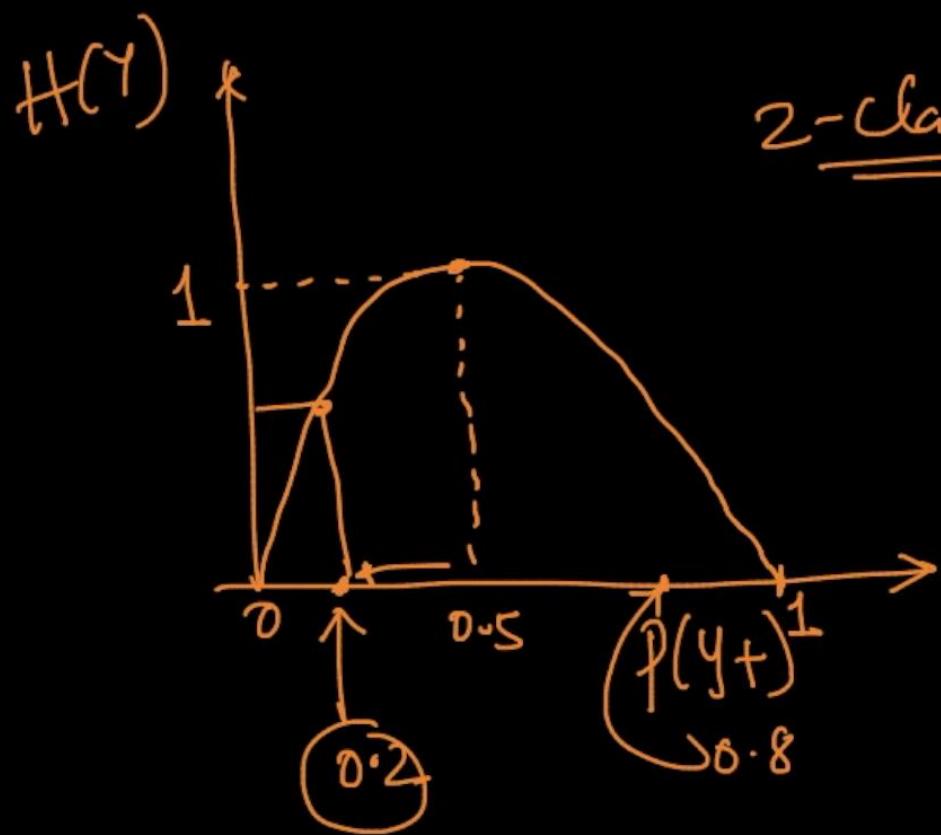
Case 2:  $\begin{cases} Y_+ \rightarrow 50\% \\ Y_- \rightarrow 50\% \end{cases}$

$$H(Y) = -0.5 \lg 0.5 - 0.5 \lg 0.5 = 1$$

Case 3:  $\begin{cases} Y_+ \rightarrow 0\% \\ Y_- \rightarrow 100\% \end{cases}$

$$H(Y) = 0$$





2-class case

$$\begin{aligned}
 & -\sqrt{-0.2 \lg 0.2} - \sqrt{-0.8 \lg 0.8} \\
 & = -\sqrt{-0.8 \lg 0.8} - \sqrt{-0.2 \lg 0.2}
 \end{aligned}$$

$$\left[ \underline{\sqrt{p(y_+)} = 1 - \sqrt{p(y_-)}} \right]$$

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



$$Y \rightarrow \underbrace{y_1, y_2, \dots, y_k}$$

equi-probable  $\rightarrow$  entropy is maximum

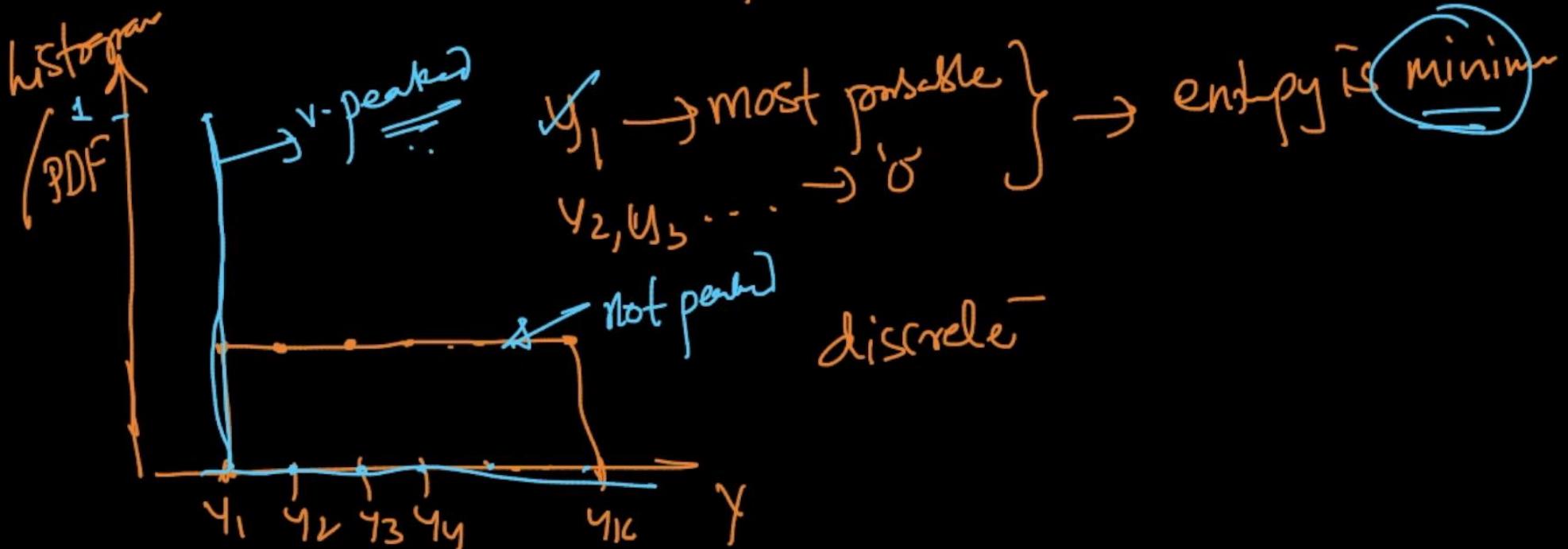
$y_1 \rightarrow$  most probable }  
 $y_2, y_3, \dots \rightarrow$  }  $\rightarrow$  entropy is minimum



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

equi-probable

✓ entropy is maximum



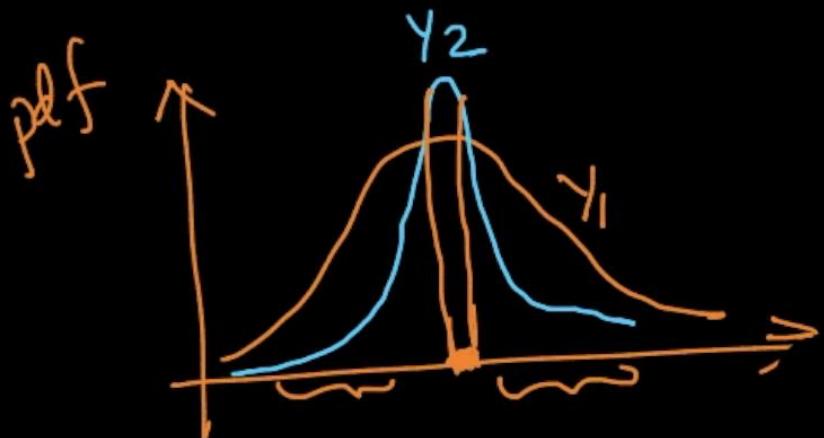
PDF

PDF

more peaked a dist is, less is its entropy

→ uniform dist

entropy is max



{  
     $y_1 \rightarrow$  less peaked  
     $y_2 \rightarrow$  very peaked

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



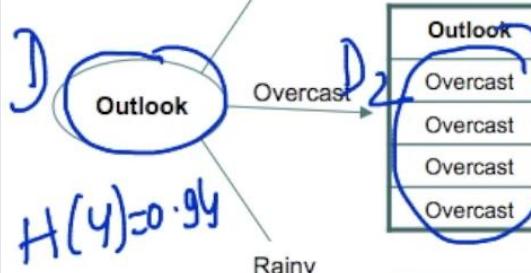
homepage.cs.uri.edu/faculty/hamel/courses/2016/spring2016/csc581/lecture-notes/32-decision-trees.pdf

# Partitioning the Data Set

$$y \rightarrow y_+ \\ y_0 \\ y_-$$

$D_1$

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



$$H(Y) = 0.94$$

$$-\frac{2}{14} \log_2 \frac{2}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

9,5

$D_3$

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

$E = .97$

$E = 0$

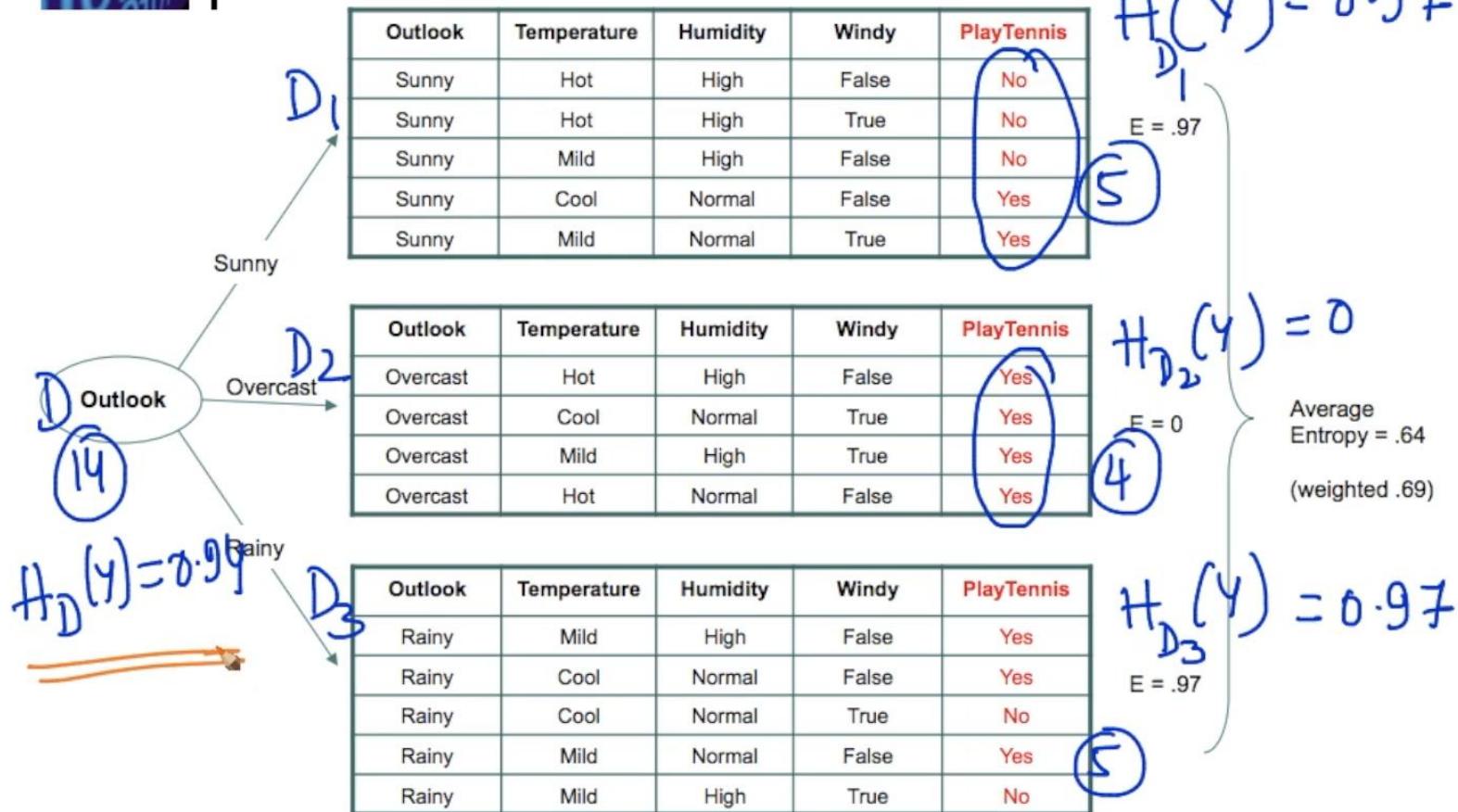
$E = .97$



Average  
Entropy = .64  
(weighted .69)

10101010  
10101010  
10101010

# Partitioning the Data Set



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

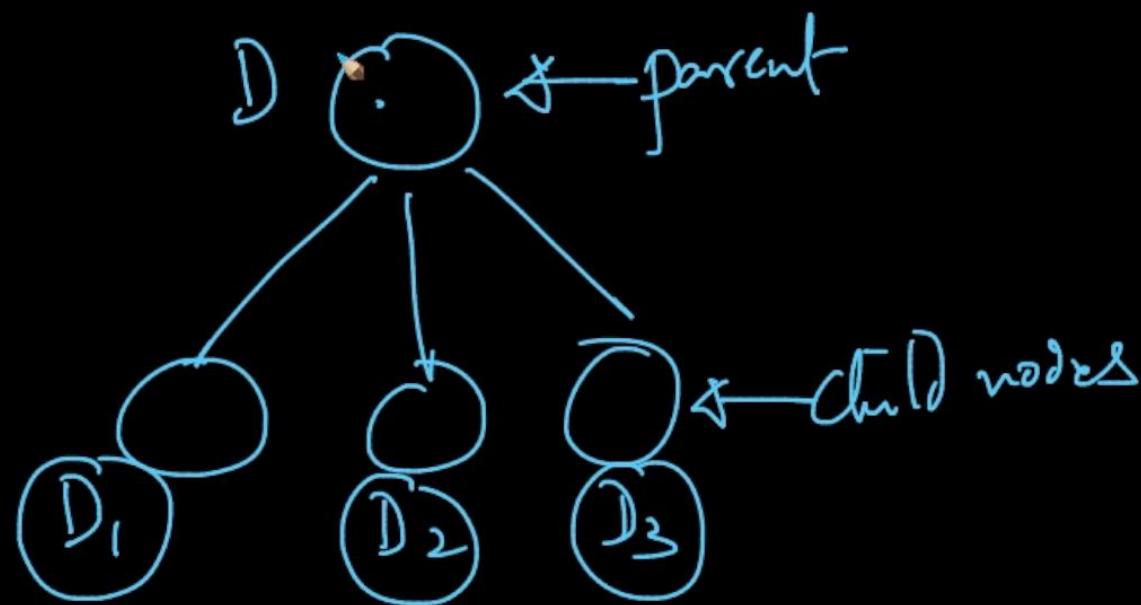
Weighted entropy  
after  $D_1, D_2, D_3$

$$\left( \frac{5}{14} \times 0.97 \right) + \cancel{\left( \frac{4}{14} \times 0 \right)} + \left( \frac{5}{14} \times 0.97 \right)$$

$$\frac{5}{14} \times 0.97 \times 2$$

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





$$-H_D(Y) + \left( \frac{|D_1|}{|D|} H_{D_1}(Y) + \frac{|D_2|}{|D|} H_{D_2}(Y) + \frac{|D_3|}{|D|} H_{D_3}(Y) \right)$$



Gini Impurity ~ similar to Entropy

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

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

Case 2:  $P(y_+) = 1$

$$P(y_0) = 0$$

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

$\uparrow$   
 $H(Y)$

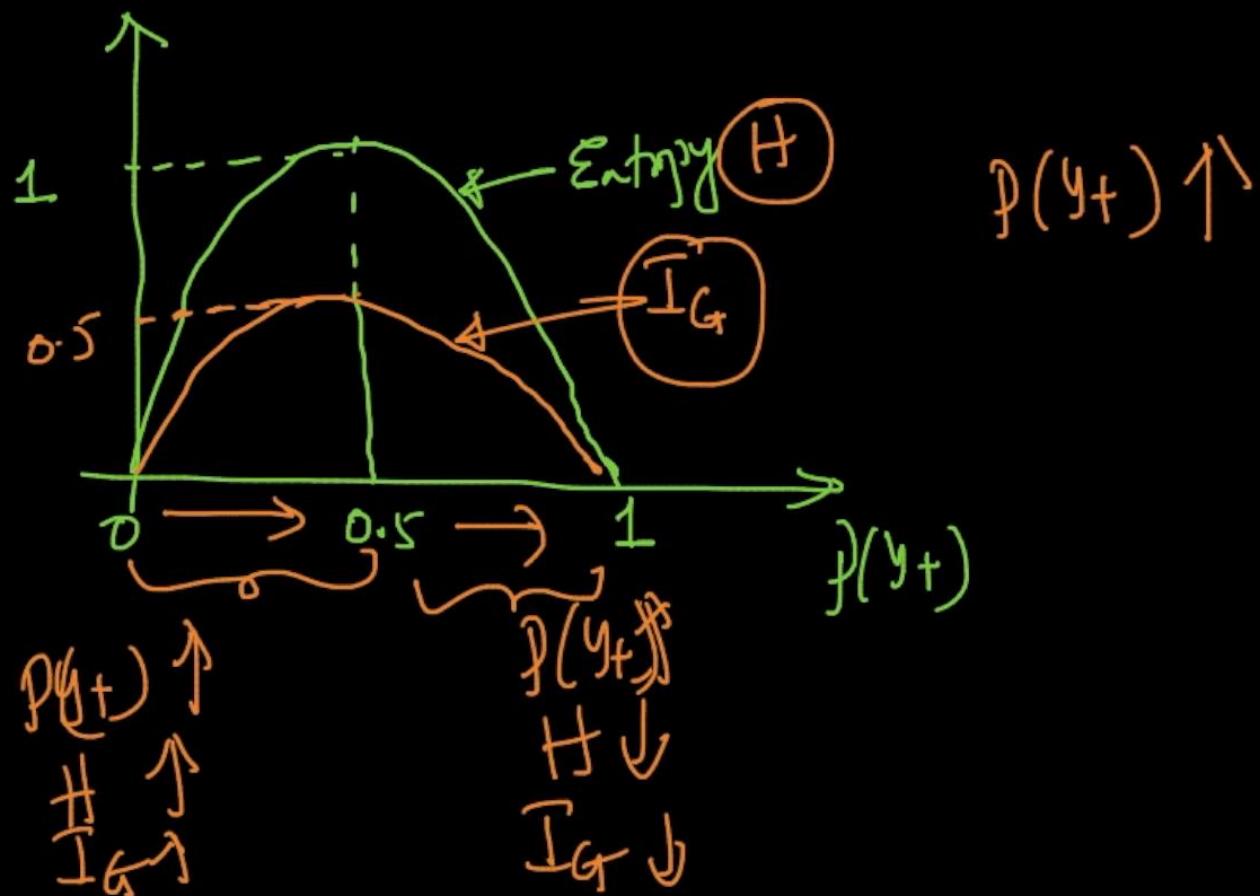
$$\left\{ \begin{array}{l} \text{Case 1: } P(y_+) = 0.5 \\ P(y_-) = 0.5 \end{array} \right.$$

$$I_G(Y) = 1 - (0.25 + 0.25) \\ = 0.5$$

$$H(Y) = 1$$



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



$$\left\{ \begin{array}{l} \text{More} \\ \text{computationally} \\ \text{efficient} \end{array} \right.$$

$$H(Y) = - \left[ P(y_+) \log_2 P(y_+) + P(y_-) \log_2 P(y_-) \right]$$

no-log

~~$H(Y)$~~

$$\left\{ \begin{array}{l} -P(y_+) \log_2 P(y_+) \\ -P(y_-) \log_2 P(y_-) \end{array} \right.$$

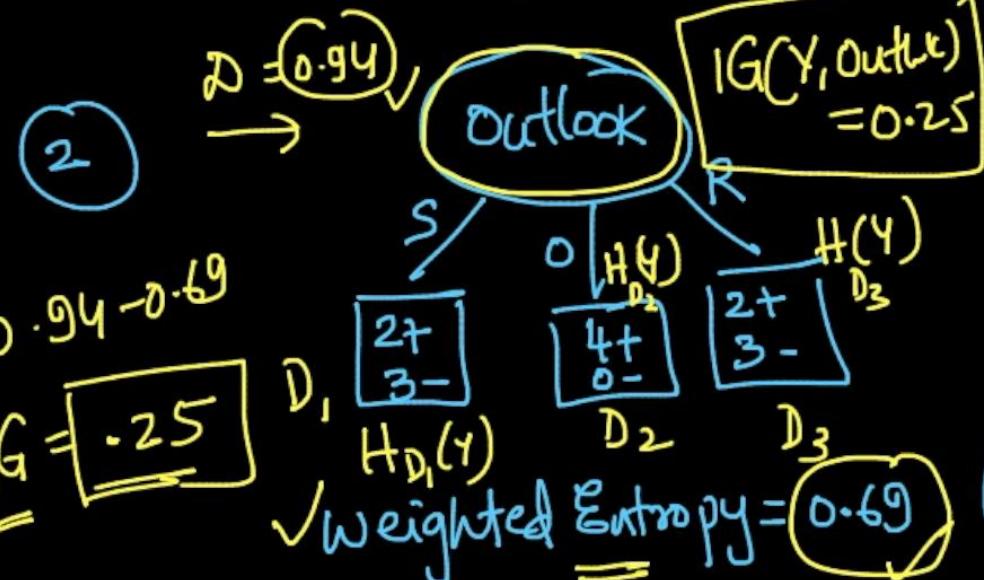
log

Construct a DT:  $H, \underline{IG} \underline{IG}$

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

① ✓  $D \rightarrow 9+, 5-$

$$H(Y) = \boxed{0.94}$$



$$D (0.94)$$

③ Temperature

$$\begin{aligned} & D_1 \left( \frac{2+}{2-} \right) H_{D_1}(Y) \\ & D_2 \left( \frac{4+}{2-} \right) H_{D_2}(Y) \\ & D_3 \left( \frac{3+}{1-} \right) H_{D_3}(Y) \end{aligned}$$

④ Humidity

$$\begin{array}{c} H \\ \diagup \quad \diagdown \\ \boxed{3+ 3-} \\ \diagup \quad \diagdown \\ N \end{array}$$

⑤ Windy

$$\begin{array}{c} F \\ \diagup \quad \diagdown \\ \boxed{6+ 2-} \\ \diagup \quad \diagdown \\ T \end{array}$$



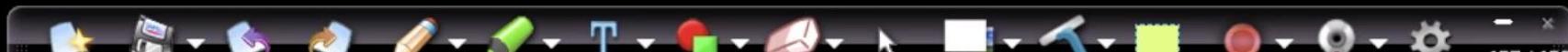
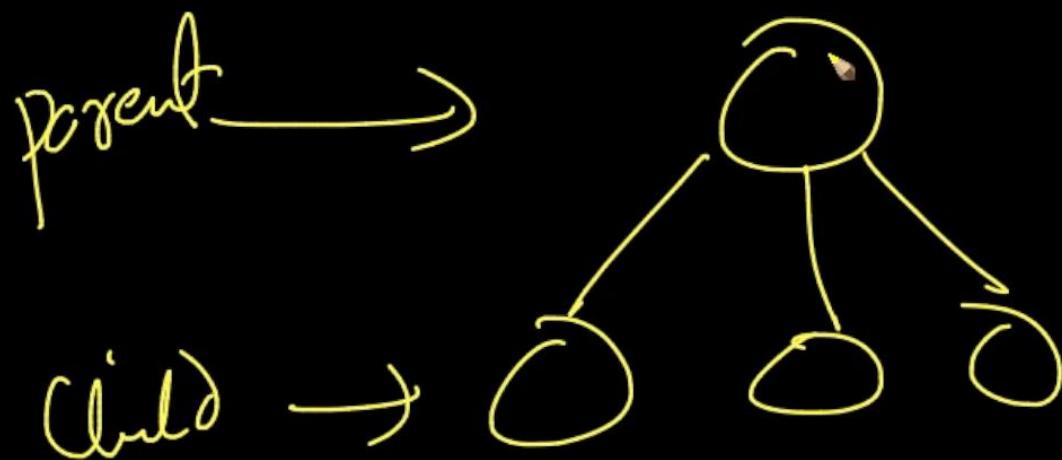
$$IG(Y, f) = H_D(Y) - \left( \sum_{i=1}^k \frac{|D_i|}{|D|} \cdot H_{D_i}(Y) \right)$$

$D \rightarrow$

→ choosing the root node

$$\left\{ \begin{array}{l} IG(Y, \text{outlook}) = 0.25 \\ IG(Y, \text{Temp}) = - \\ IG(Y, \text{Humidity}) = - \\ IG(Y, \text{Windy}) = - \end{array} \right.$$

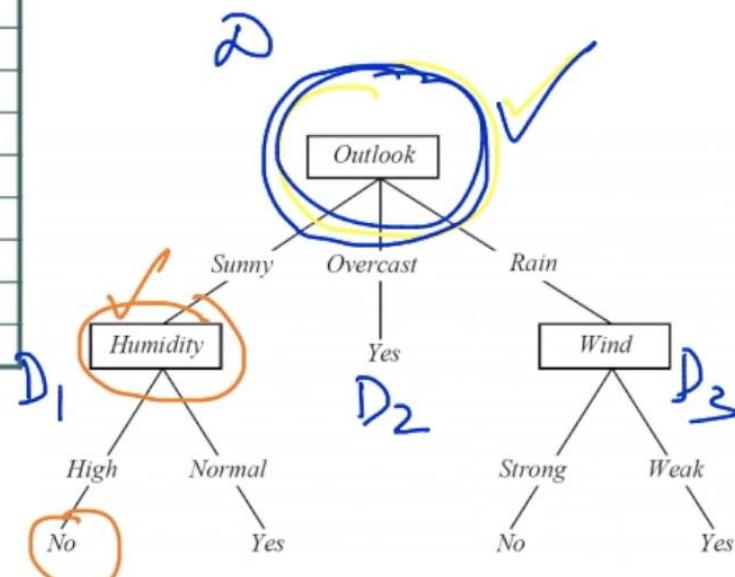


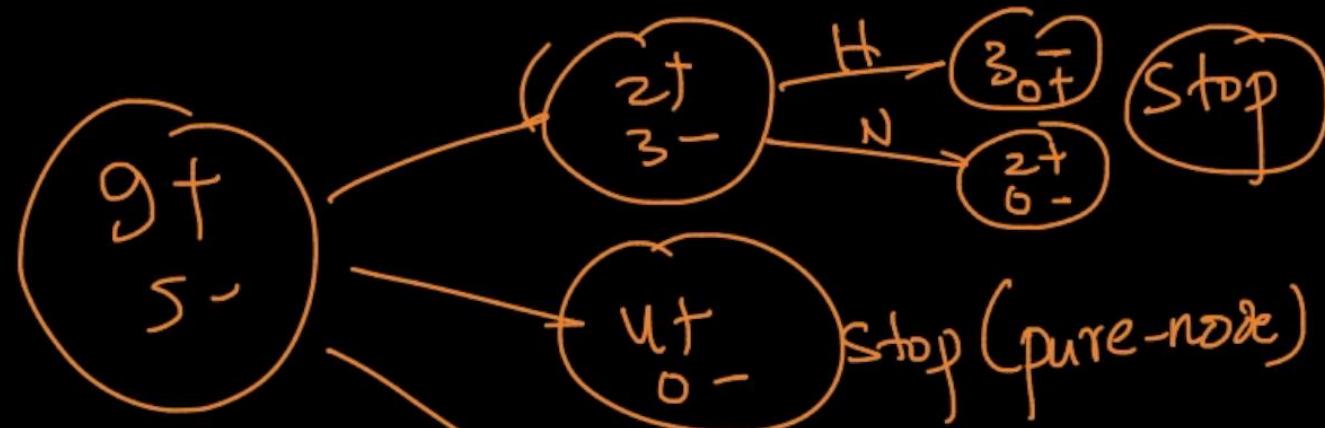
$$IG(Y, f) = \text{entropy}@\text{parent-level} - \text{weighted entropy}@(\text{child level})$$


# Play Tennis Example

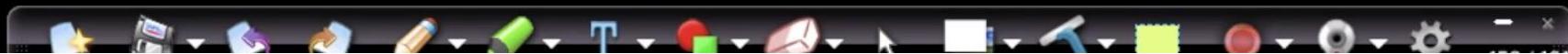
Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No ✓
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal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

$f_1 \ f_2 \ f_3 \ f_4$



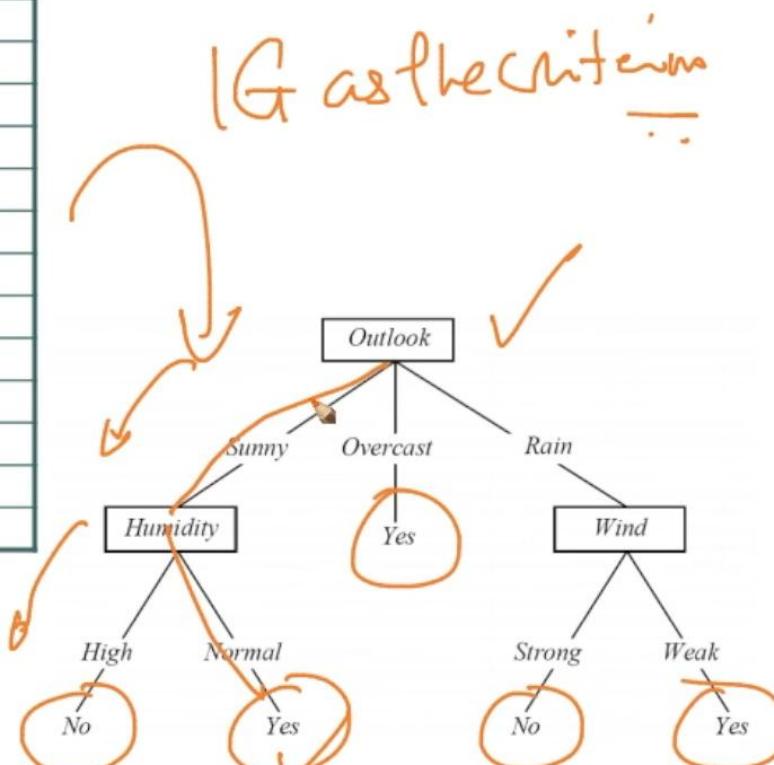


{ recursively breaking each node  
 using IG as the criterion



# 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

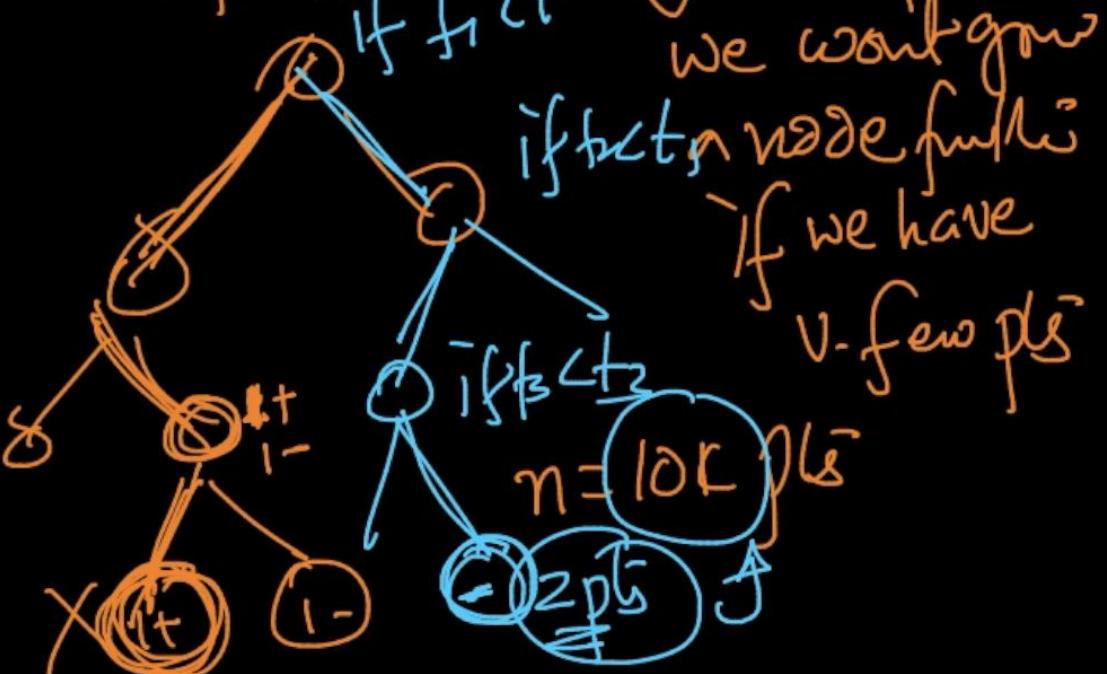
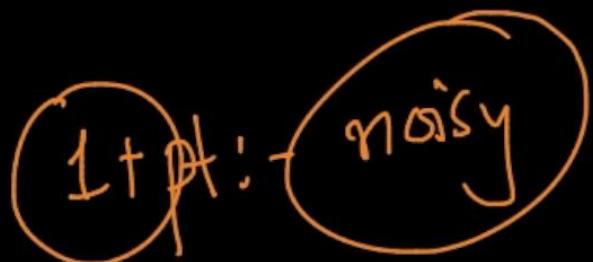


①

pure-node  $\rightarrow$  Stop growing the node

②

Can't grow the tree anymore because of  
lack of pls

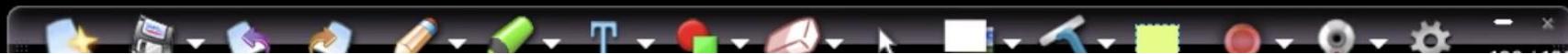
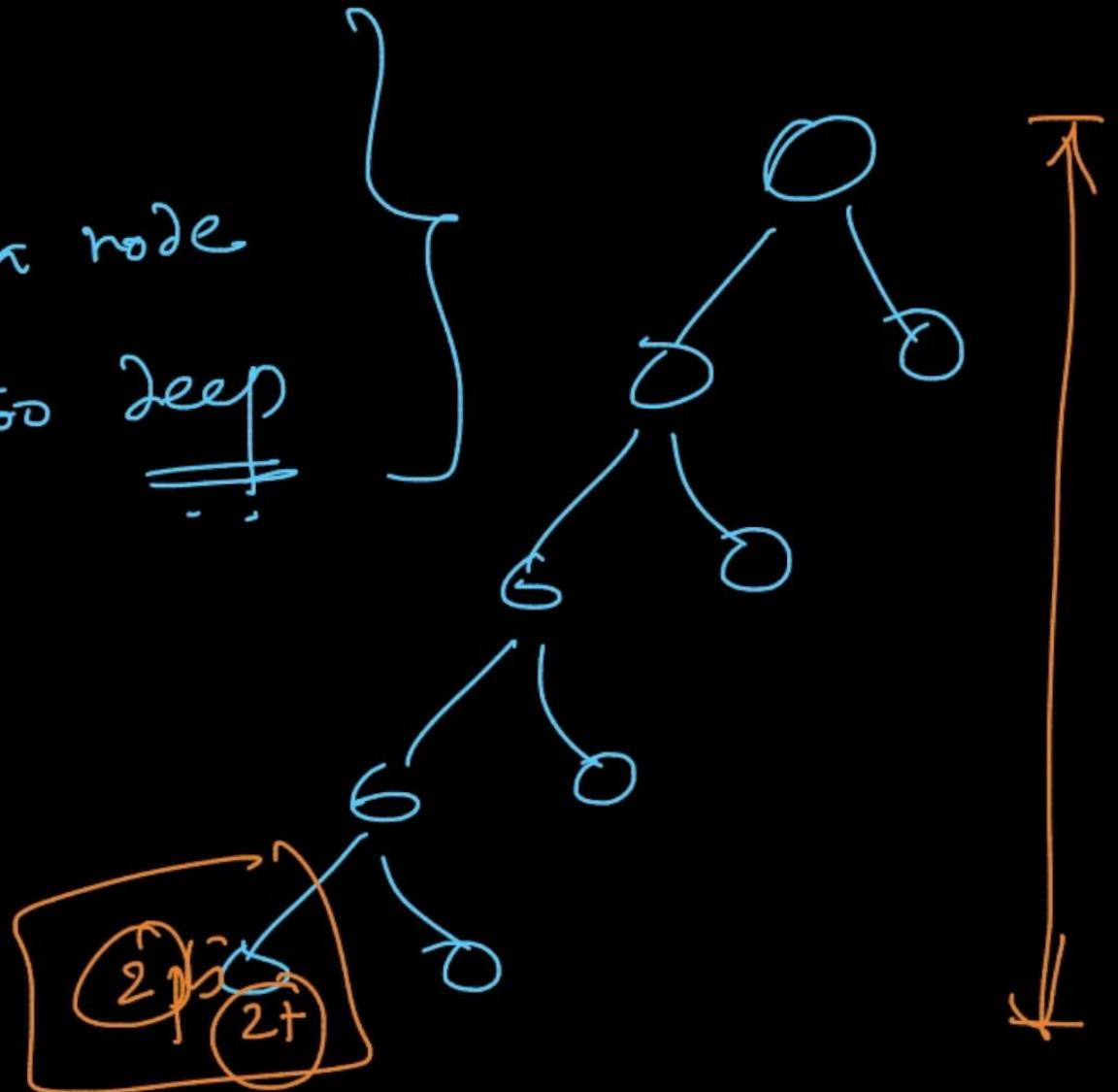


① pure node

② few pls @ a node

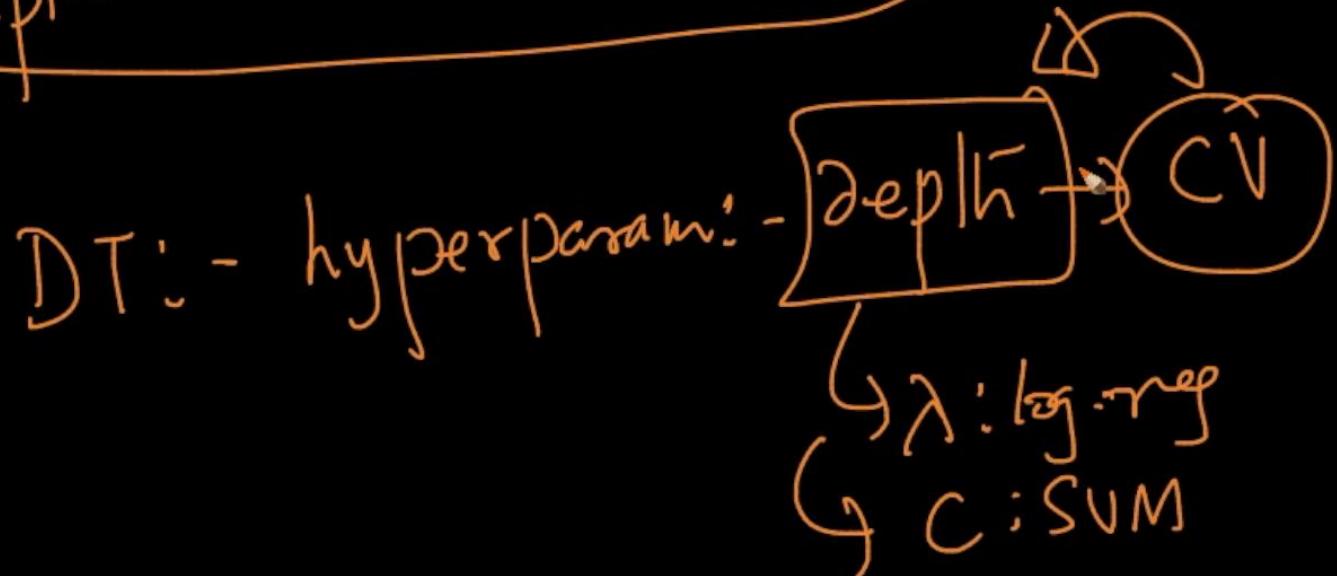
③ if we are too deep

n  $\in$  lot



{ depth of the tree ↑ ; overfitting ↑  
(few pts)

depth is small  $\Rightarrow$  underfit



## Splitting numerical features

Construct a DT :- Splitting a node → IG

IG :-  
entropy  
Gini · impurity → computationally efficient

discrete v.v (n) Categorical features



$f_i$	$y$
2.2	1
2.6	1
3.5	0
3.8	0
4.6	1
5.3	0

$f_i$ : numerical  
 ↗ integer  
 ↗ real-value

→ Split based on Cat. Var  
 $f_2$  ↓ 3 Categories

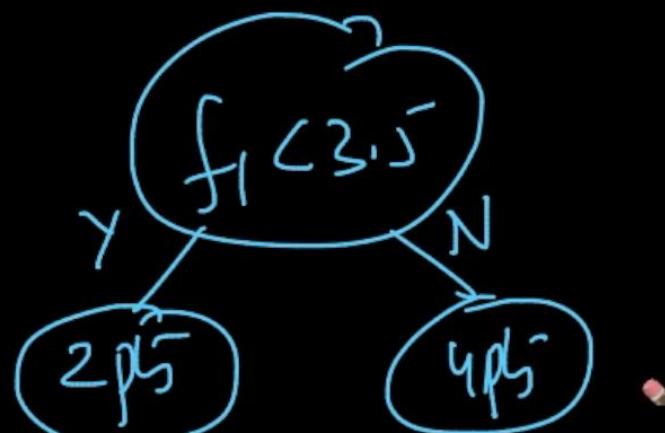


$f_1$	$y$
2.2	1
2.6	1
3.5	0
3.8	0
4.6	1
5.3	0

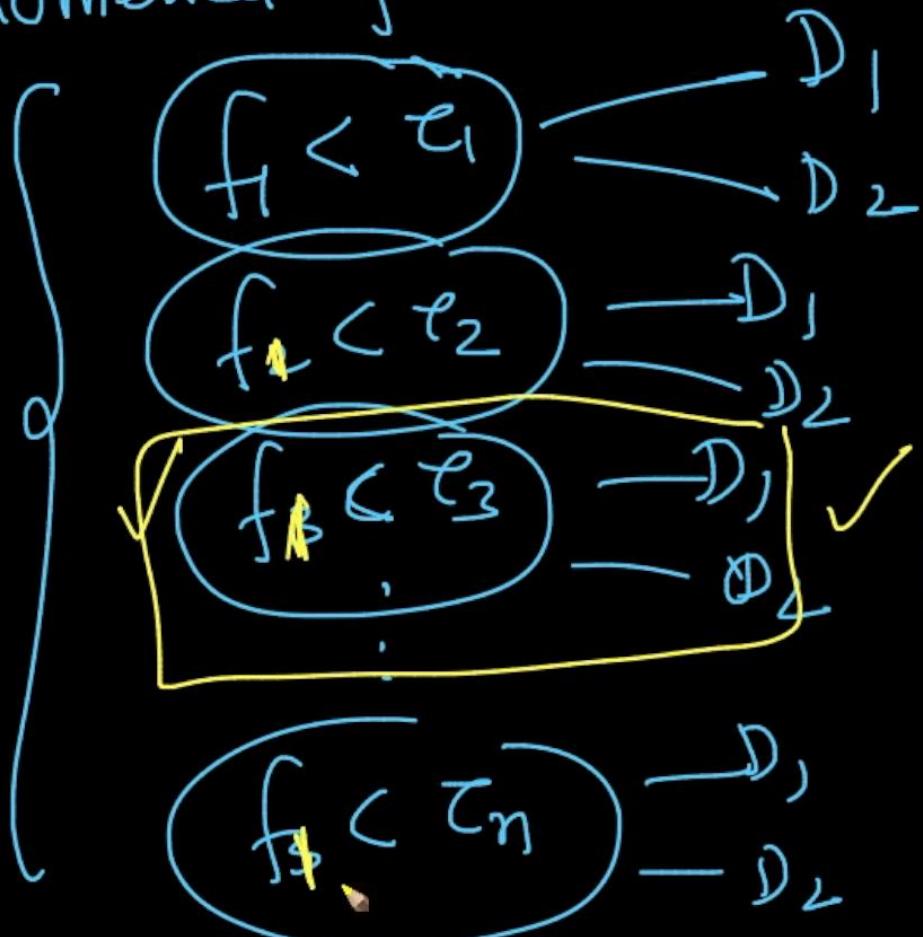
$f_i$ : numerical  
 ↘ integer  
 ↘ real-value

- ① Sort the numerical feature in desc-order
- ②

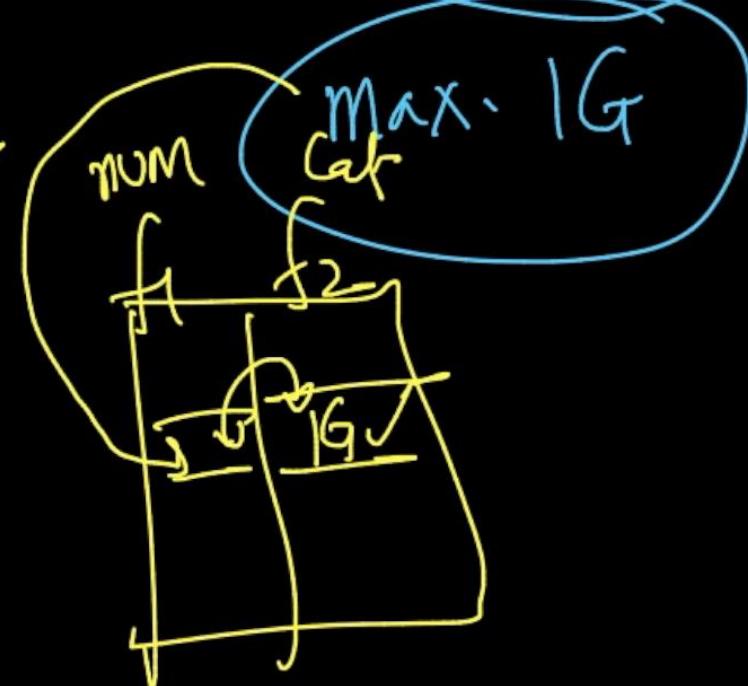
$$\begin{cases} f_1 < 2.2 \\ f_1 < 2.6 \\ f_1 < 3.5 \\ f_1 < 3.8 \end{cases}$$



✓ numerical - feature

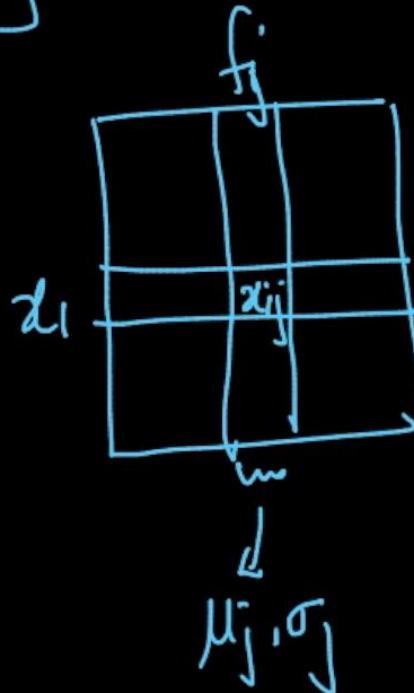
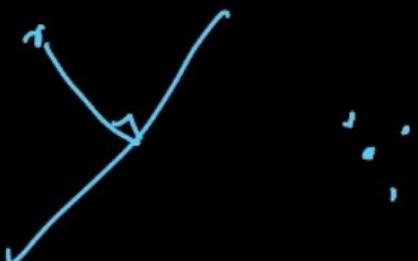


$n$ -possible  
splitting  
criterion



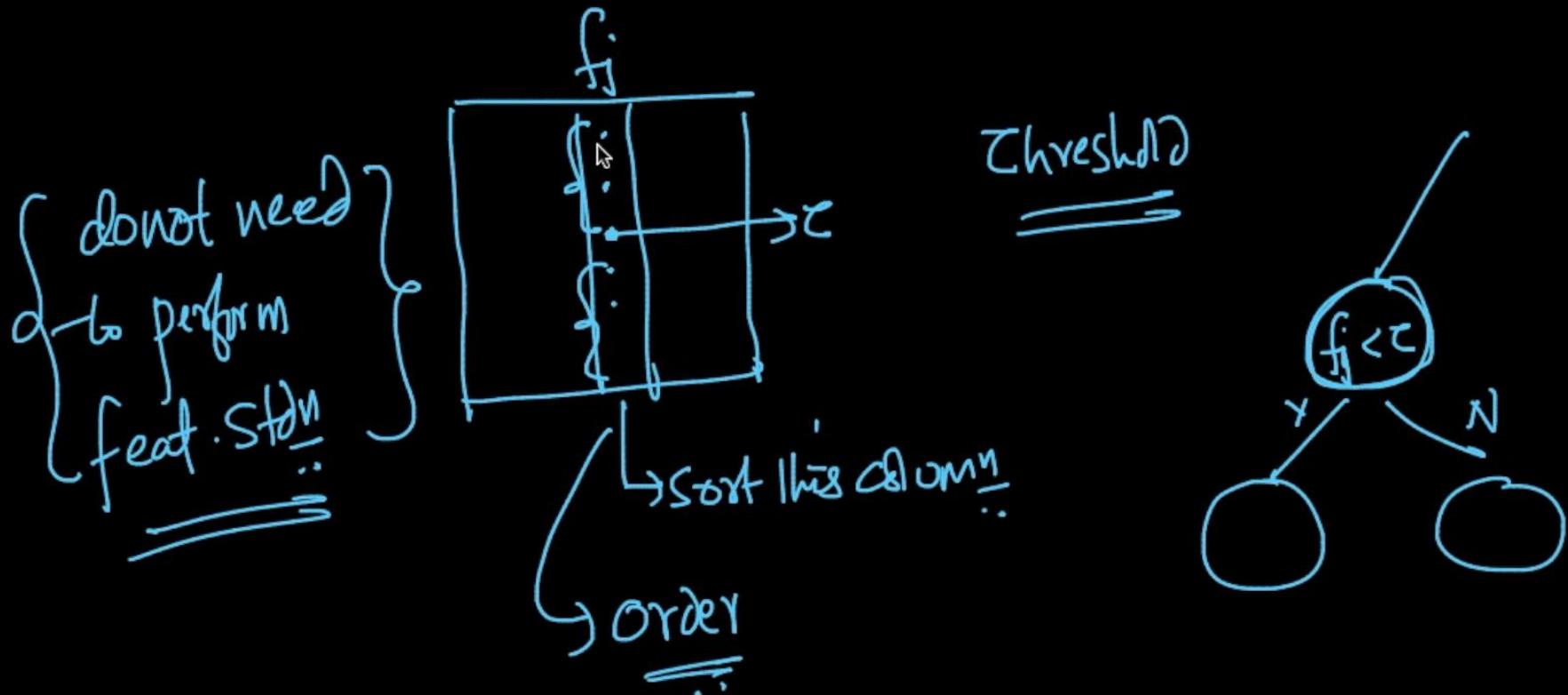
## Feature Standardization:

dist.  $\{$  Logistic regn, SVM, KNN  $\}$   $\rightarrow$  feature std.

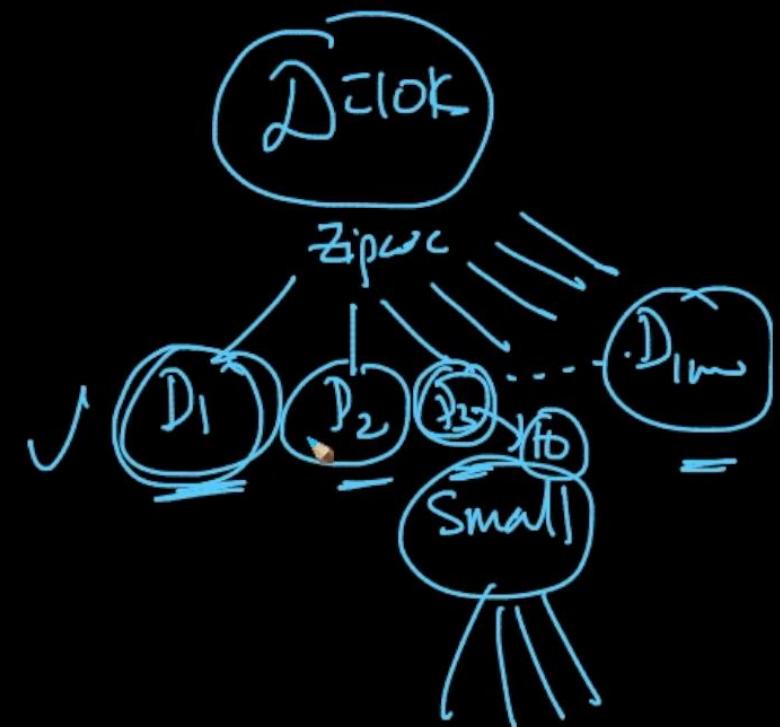
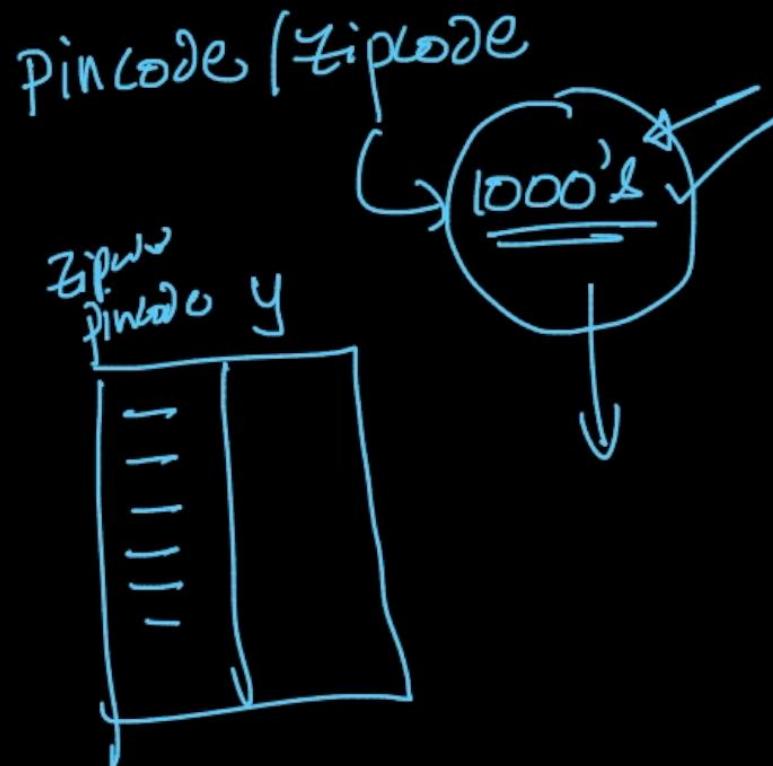


$$x'_{ij} = \frac{x_{ij} - \mu_j}{\sigma_j}$$

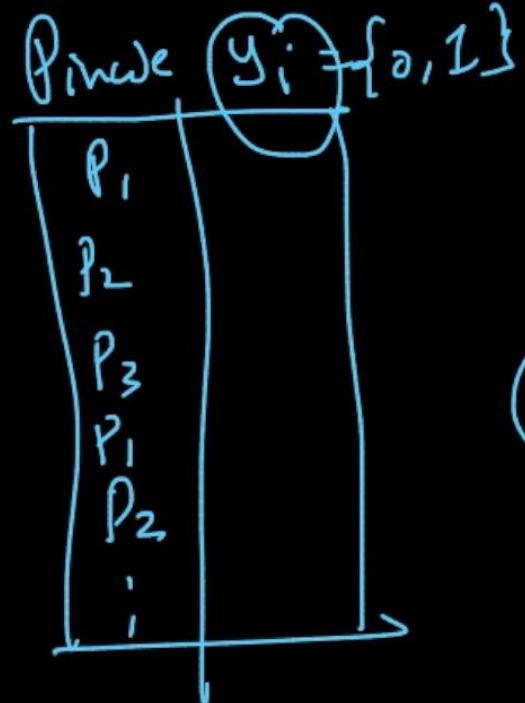
Decision Trees: - not a distance based method



✓ Categorical features with many categories  
↳ levels



hack (feature engg)



20 times  $p_j$   
19 times  $y_i = 1$

$$p_j \rightarrow \boxed{\phantom{00}}$$

$$\begin{aligned} &P(y_i=1 | p_j) = \frac{19}{20} \\ &= \frac{\# y_i=1 \leq p_j}{\# p_j} \end{aligned}$$

Pincode / Zipcode  $\rightarrow$  Categorical

Convert it numerical  
feature

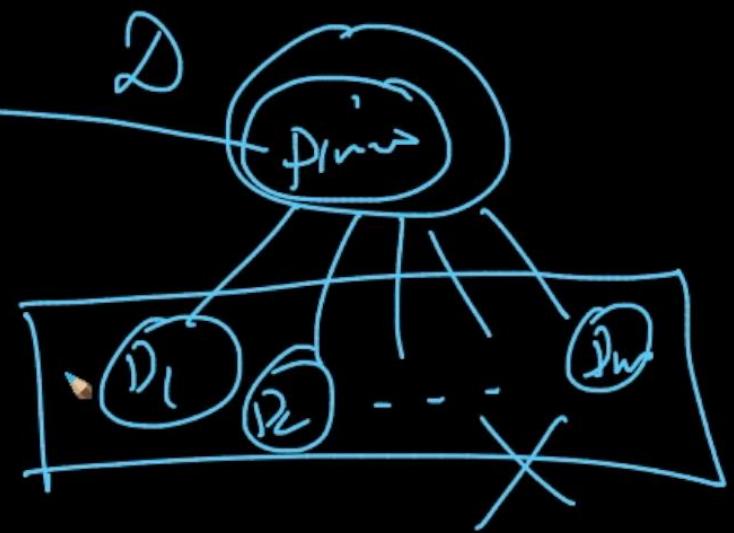
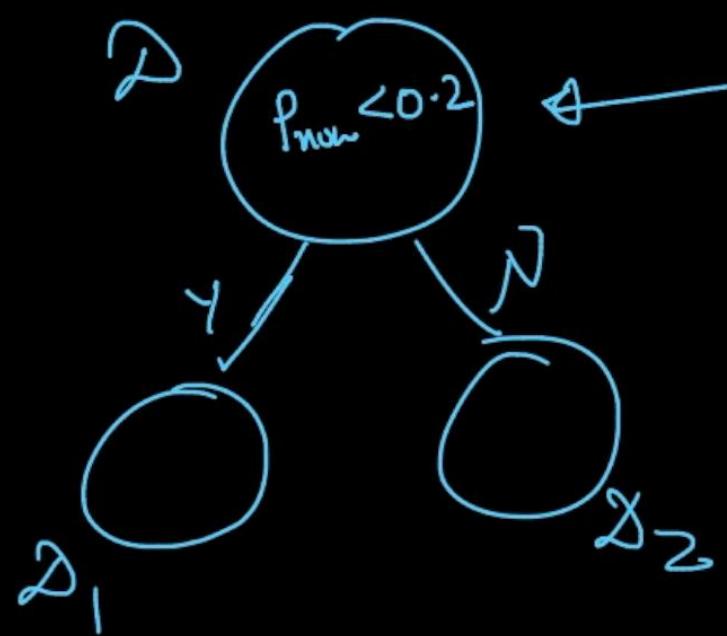
$p_j \rightarrow$  numerical feature =  $P(Y_i=1 | p_j)$

$\mathcal{D} \rightarrow$  remove pincode feature

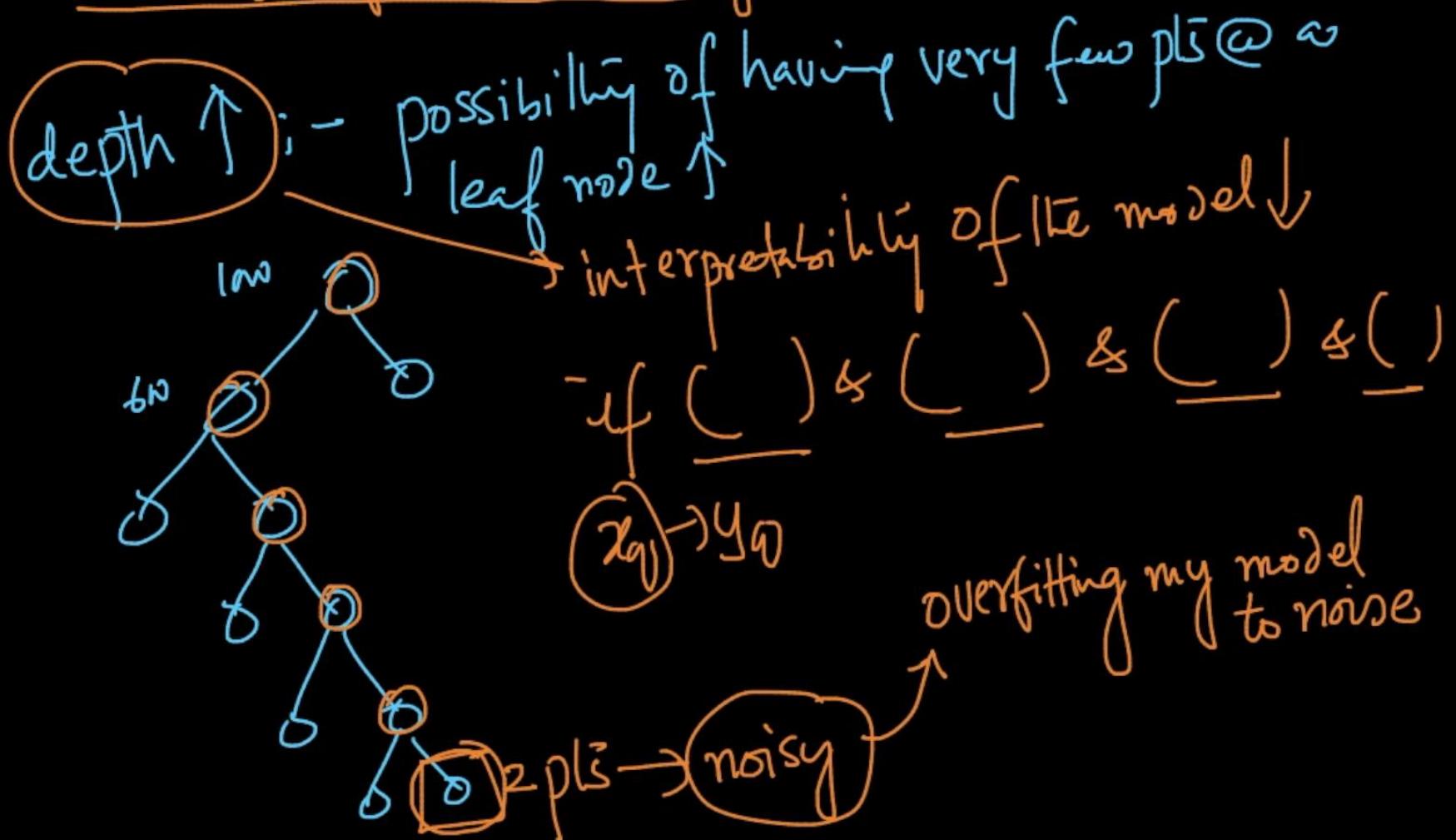
numerical features

$p_i \& p_j$

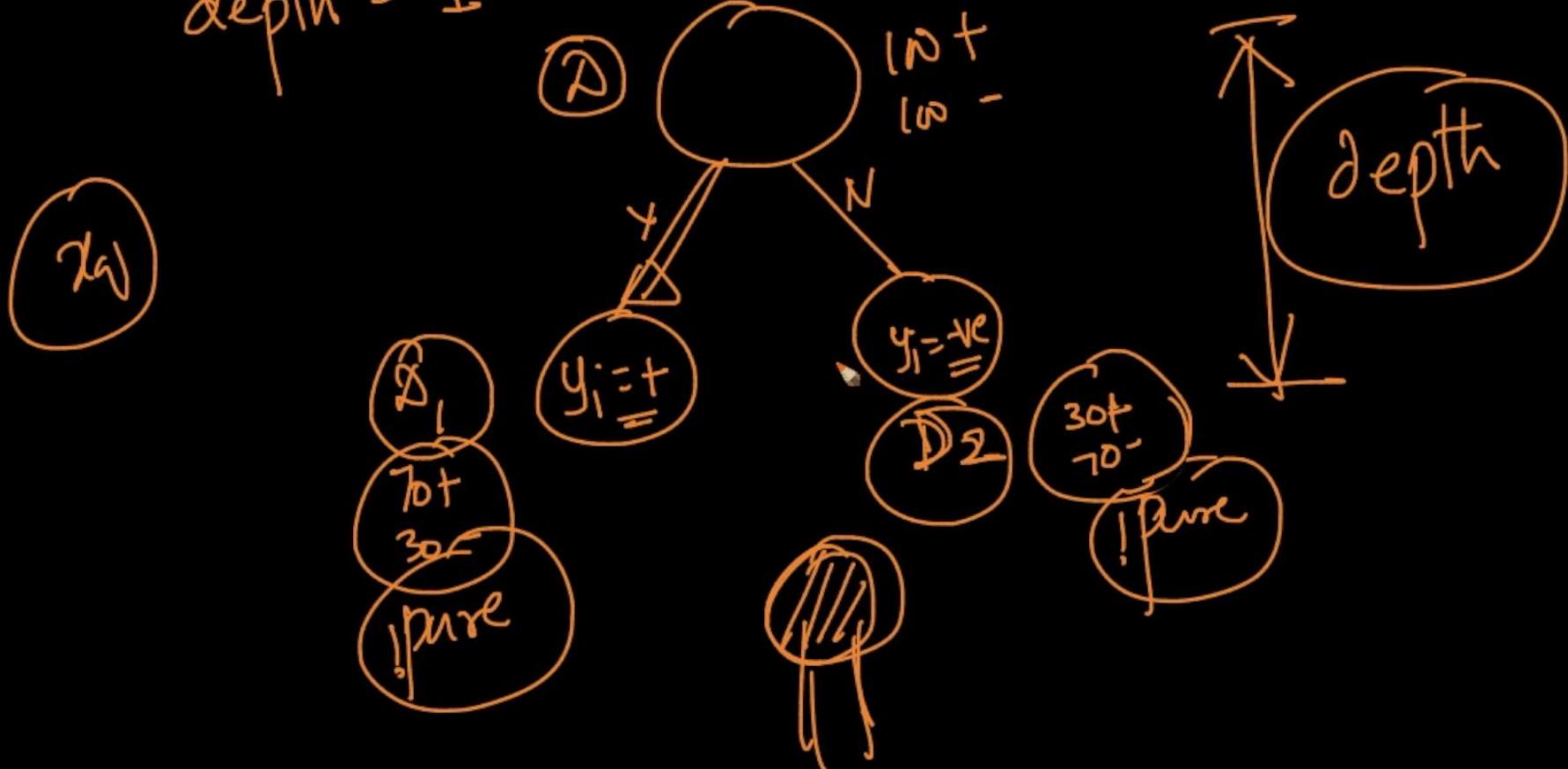




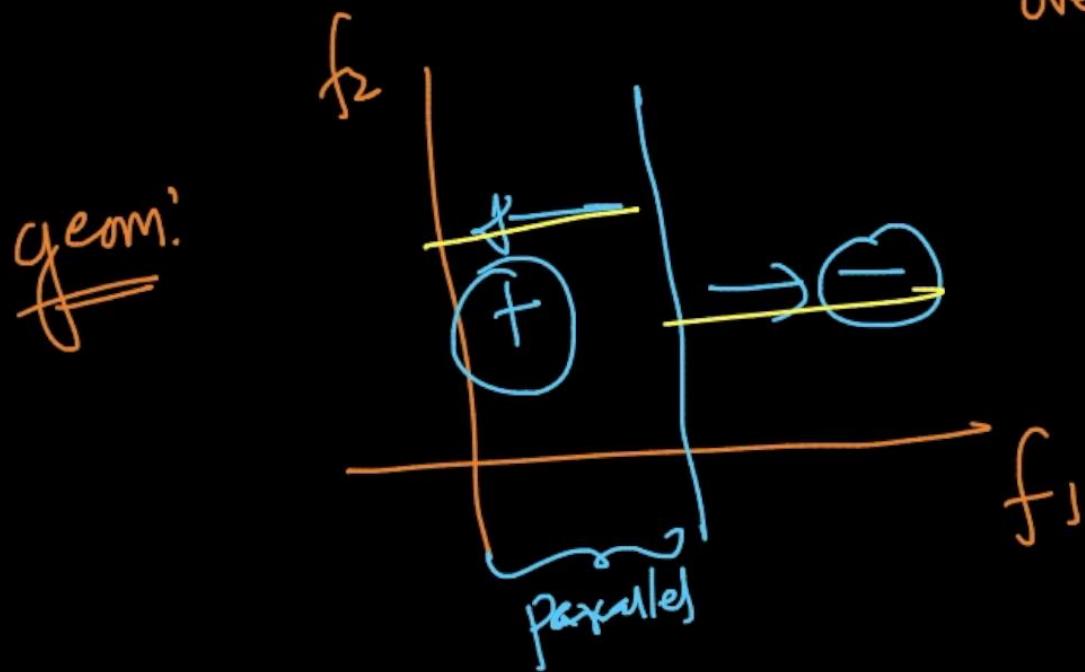
## Overfitting & Underfitting:



depth = 1 = decision stump

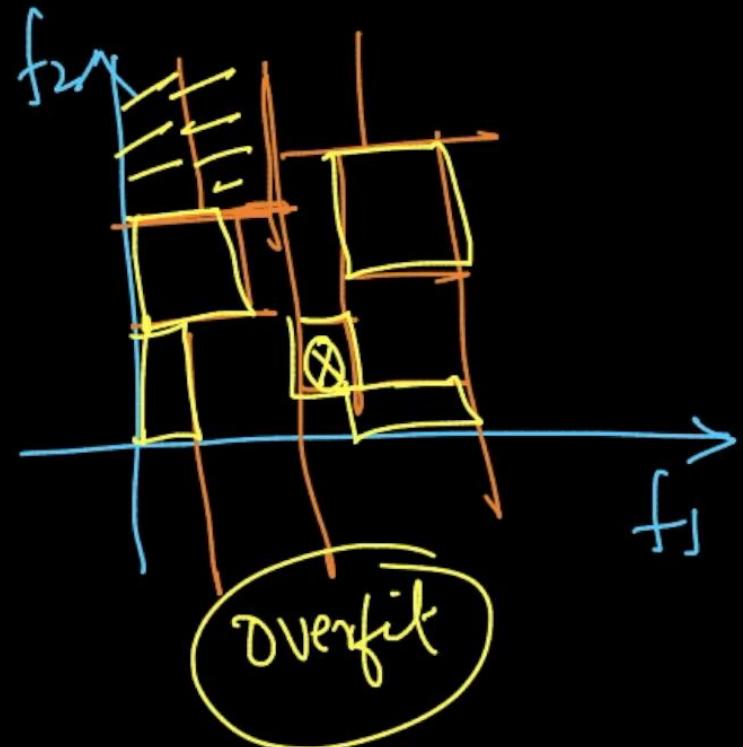


depth :-  $\leq$



too deep  
↓  
overfit

too shallow  
↓  
underfit



## Train & Run time complx:

Train:  $\sim O(n(\lg n)d)$        $n = \# \text{pts} D_{\text{Train}}$   
(Time)                   $d = \text{dim.}$

numerical features: - (threshold)

↳ algorithmic methods

$n \lg n$  → sorting

$O(n \lg n)$

large d



After  
training:-

@Runtime Space:-

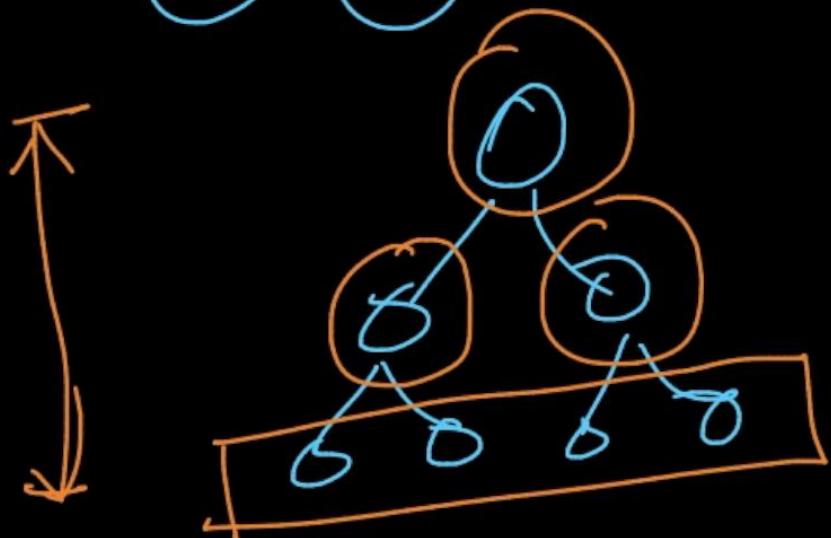
$x_v \rightarrow f_v$

Store my DTree  $\rightarrow O(\text{nodes})$

{ if else }

{ nested if else }

$\uparrow \# \text{ internal-nodes}$   
 $+ \# \text{ leaf nodes.}$

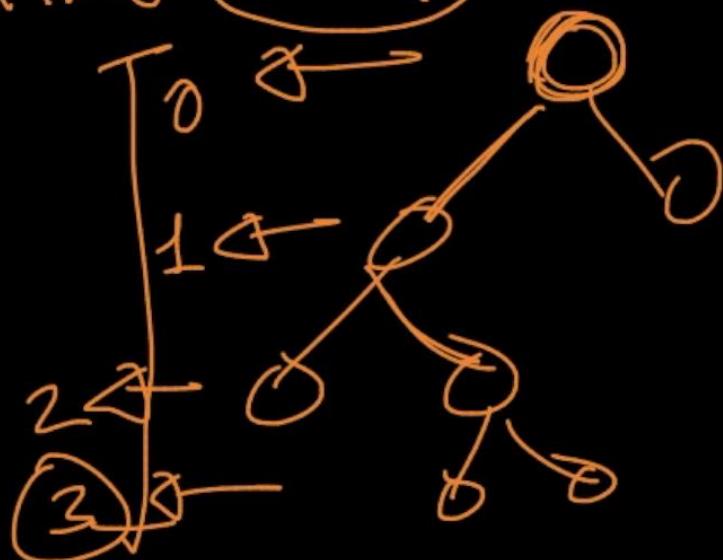


$$\text{depth} = \underline{\underline{2+0}}$$

$\text{depth} \uparrow \Rightarrow \text{interpret} \downarrow$

runtime Space :- reasonable

runtime Time complexity :-  $T_q \rightarrow Y_{qj}$



$$O(\text{Depth})$$

$k = \text{max-depth of any leaf-node}$

at most 3 Comparisons



DT :- large data, dim is small

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

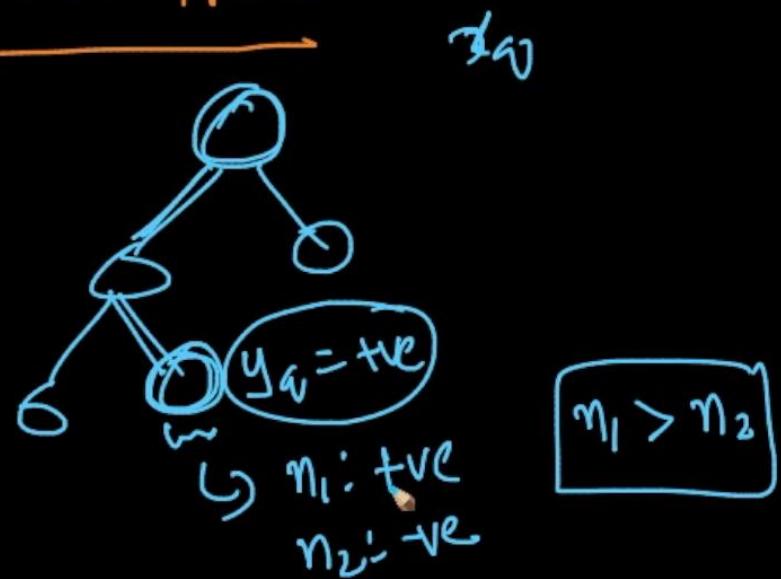
RF, GBDT

↑ popular in internet application



Regression using Decision Trees:

DT  $\rightarrow$  Classification



✓  $|G \rightarrow \text{classification}$

$(f_1, f_2)$

MSE or MAE  $\rightarrow$  regression

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Entropy

$$\hat{y}_i = \underset{D}{\text{mean}}(y_i)$$

$$y_i \in \mathbb{R}$$

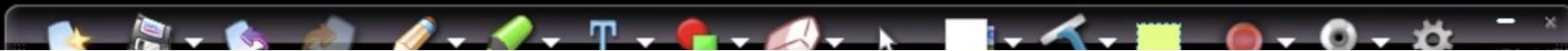
avg = mean / median

$f?$

$$\hat{y}_i = \underset{D_2}{\text{mean}}(y_i)$$

$$\boxed{\omega_1, \text{MSE}_{D_1}, \text{MSE}_{D_2}}$$

$$\hat{y}_i = \underset{D_1}{\text{mean}}(y_i)$$



Classifn:  $f_i$  :- reducing entropy :- "0"

regsn:  $f_i$  :- reduce MSE  $\rightarrow$  "0"

$$MSE(y_i) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

MAE = MAE

$$\text{Median}(|y_i - \hat{y}_i|)$$

MAE



Contents - Google Drive X URI.cs 32-decision-trees.ppt X play tennis decision tree X Decision Tree Learning X 1.10. Decision Trees - X sklearn.tree.Decision X Chekuri Srikan...

← → C i scikit-learn.org/stable/modules/tree.html

1.10.3. Multi-output problems

1.10.4. Complexity

1.10.5. Tips on practical use

1.10.6. Tree algorithms: ID3, C4.5, C5.0 and CART

1.10.7. Mathematical formulation

- 1.10.7.1. Classification criteria
- 1.10.7.2. Regression criteria

*classfn/regress*

↓ depth = underfit  
depth↑ = overfit

Decision Tree Regression

Some advantages of decision trees are:

## Cases:- (DT)

✓ imbalanced data: - balance it → upsampling  
Excessut

90% +ve  
10% -ve ↳ impacts Entropy calc (MSE)

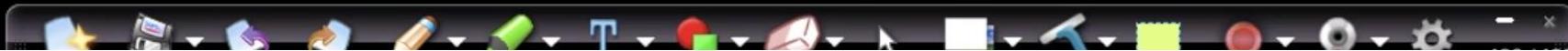
~~large d~~ → @ each node, Split  
↓  
each features (G)  
↓  
Time compx to train DP increases

Cat. feat  $f_1 \rightarrow$  avoid one-hot encoding  
 $\rightarrow$  ID features  $\rightarrow$  dim  $\searrow$

✓ Categorical feat → numerical feat  
lots of levels       $\underbrace{\qquad\qquad}_{\text{f} = \text{C}_1}$  } ✓

$$\sqrt{P(y_i=1 \mid f=C_1)}$$

Similarity Matrix :- DT need the features explicitly



Multi-class classfn: → ~~OVR~~

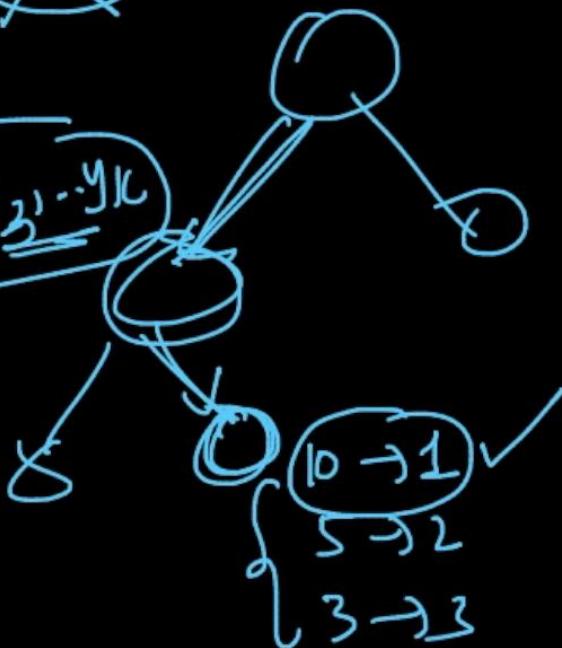
$I_{GJ}$

DT

Entropy

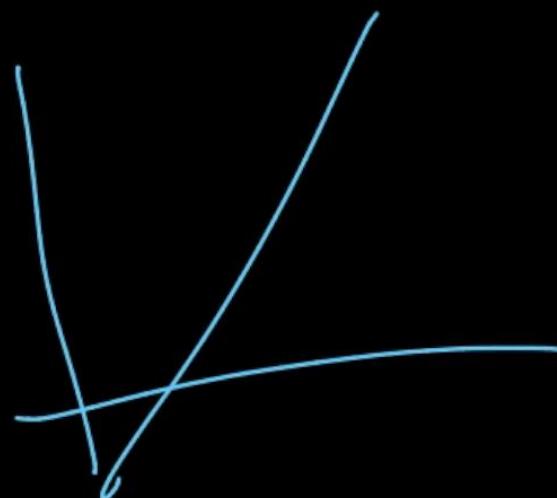
$H \rightarrow y_1, y_2, y_3, \dots, y_K$

$y = \{1, 2, 3\}$



{  
1 → 1  
2 → 2  
3 → 3

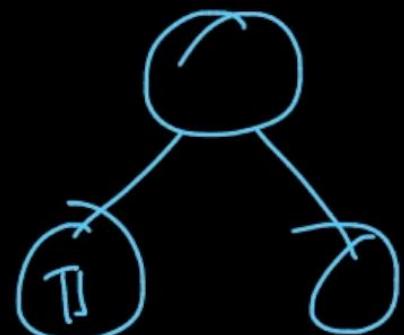
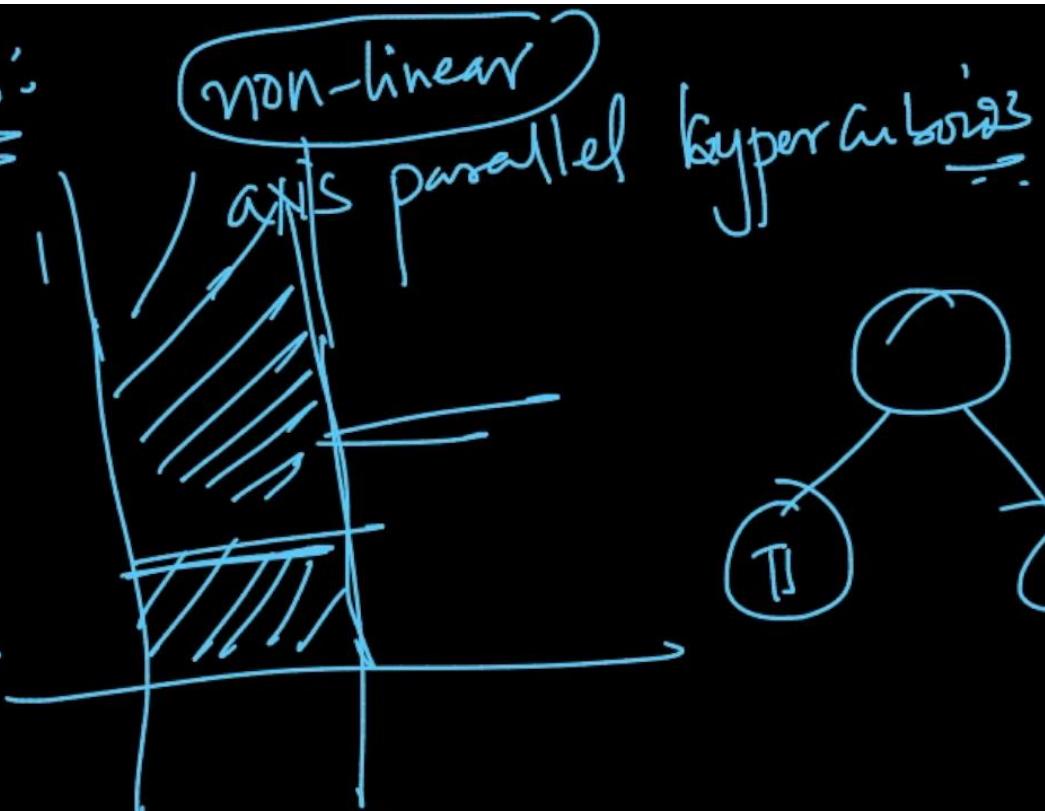
decision-Surfaces:



non-linear

axis parallel

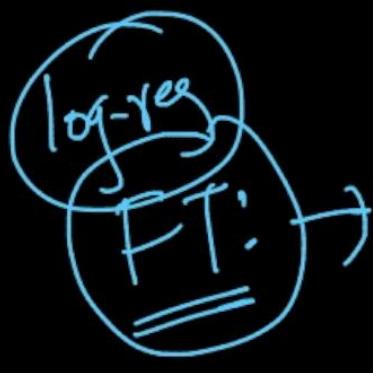
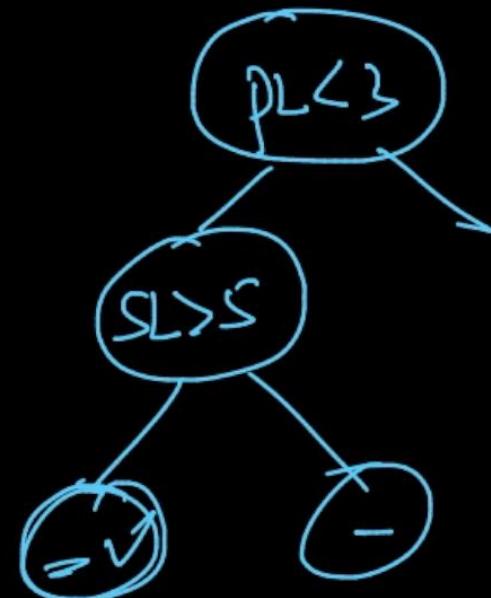
hypercurv<sup>is</sup>



feature-interactions: (logical)

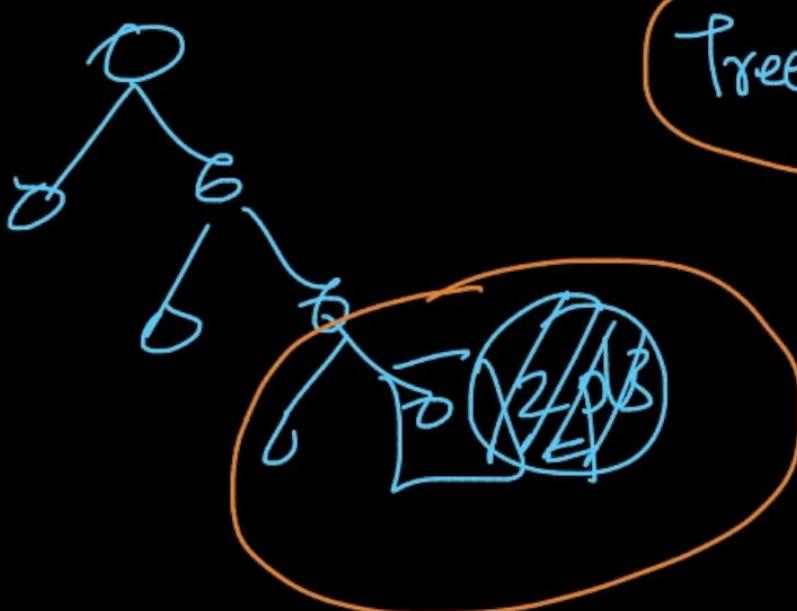
$SL, SW, PL, PW$

$\overbrace{(PL < 3)}^{\sim} \text{ AND } \overbrace{(SL > S)}^{\sim}$

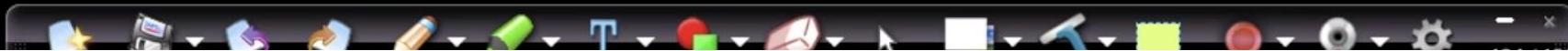


Outliers:

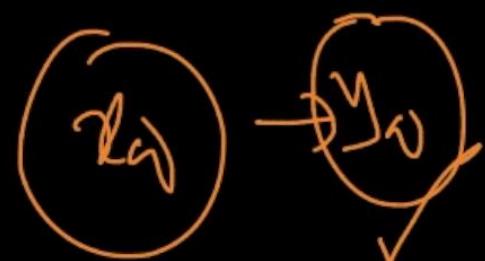
depth↑ ; - outlier will impact



Tree ~~unstable~~



Interpretability:-

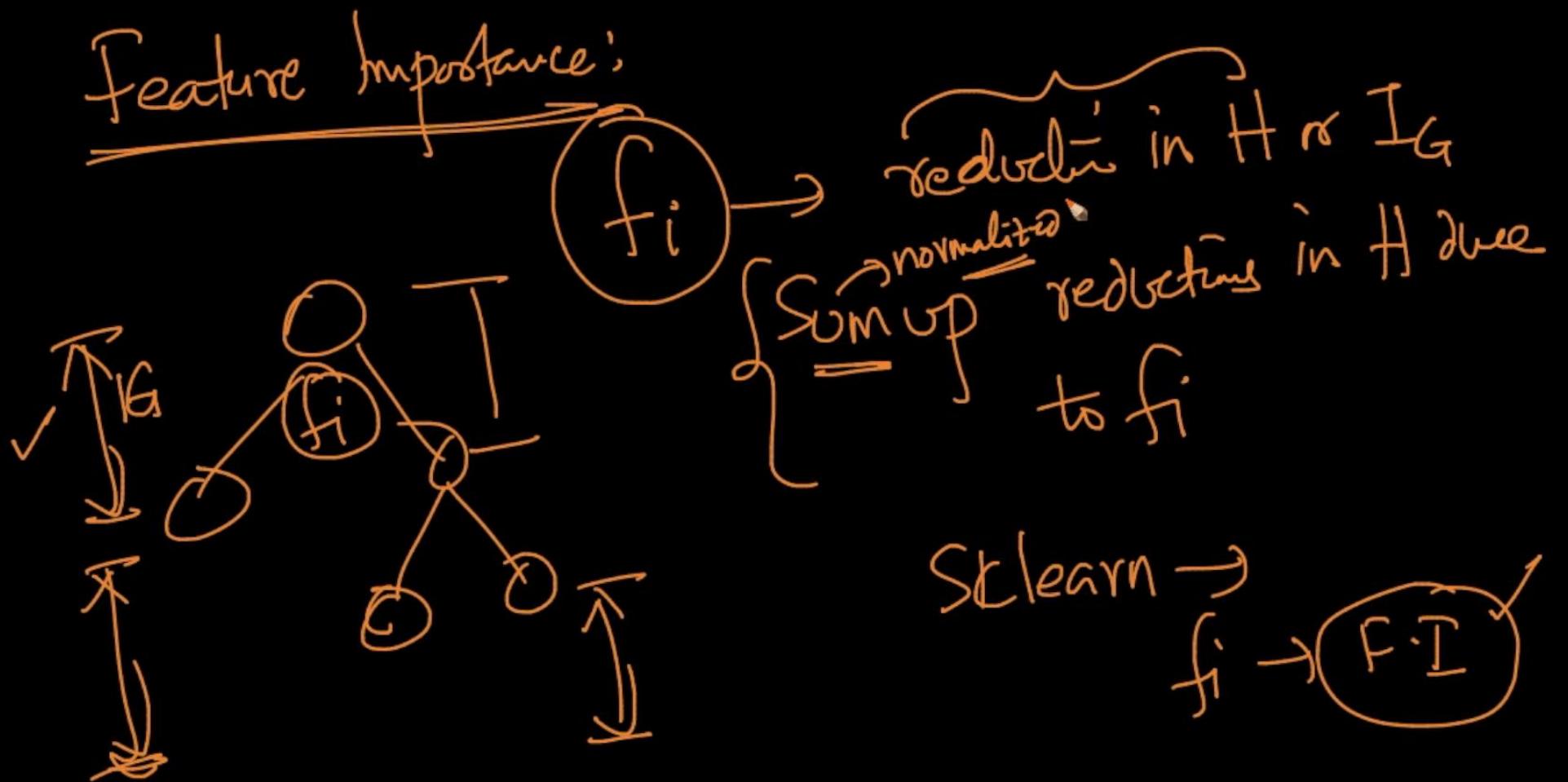


if  
else

if  
—  
—  
—

}





Press Esc to exit full screen

$$D \xrightarrow{\text{Var}} D_1, D_2, \dots, D_k$$
$$IG(Y; \text{Var}) = \sum_{i=1}^k \frac{|D_i|}{|D|} H_{D_i}(Y) - H_D(Y)$$