

Decision Trees (DT) \rightarrow if -- else

kernel
↑

KNN, Naive Bayes, Logistic regression, LinearRe, SVM

instance
based
method

probabilistic
method

geometric, hyperplane

DT: nested if -- else classifier

EDA: Iris Dataset

$y_i \in \{1, 2, 3\}$
(SL, SW, PL, PW)

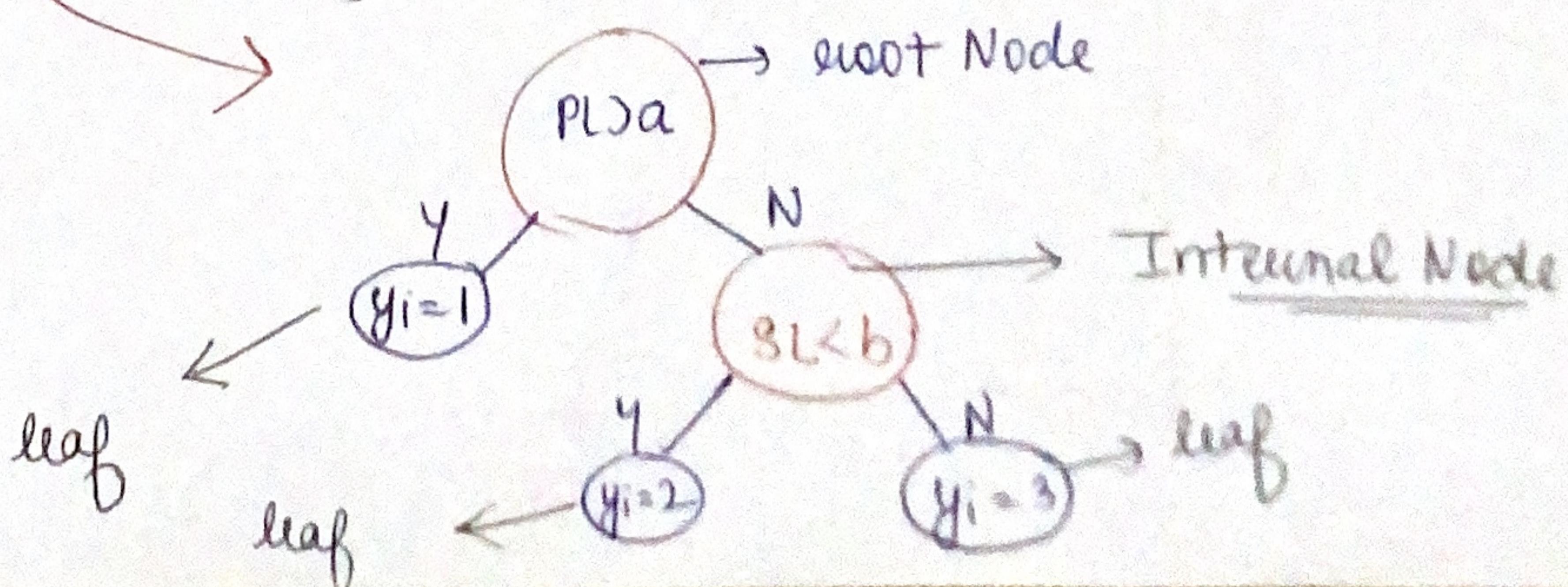
Simple $\left\{ \begin{array}{l} \text{if } PL < 5 \\ \text{then } y_i = 1 \end{array} \right\}$

Model

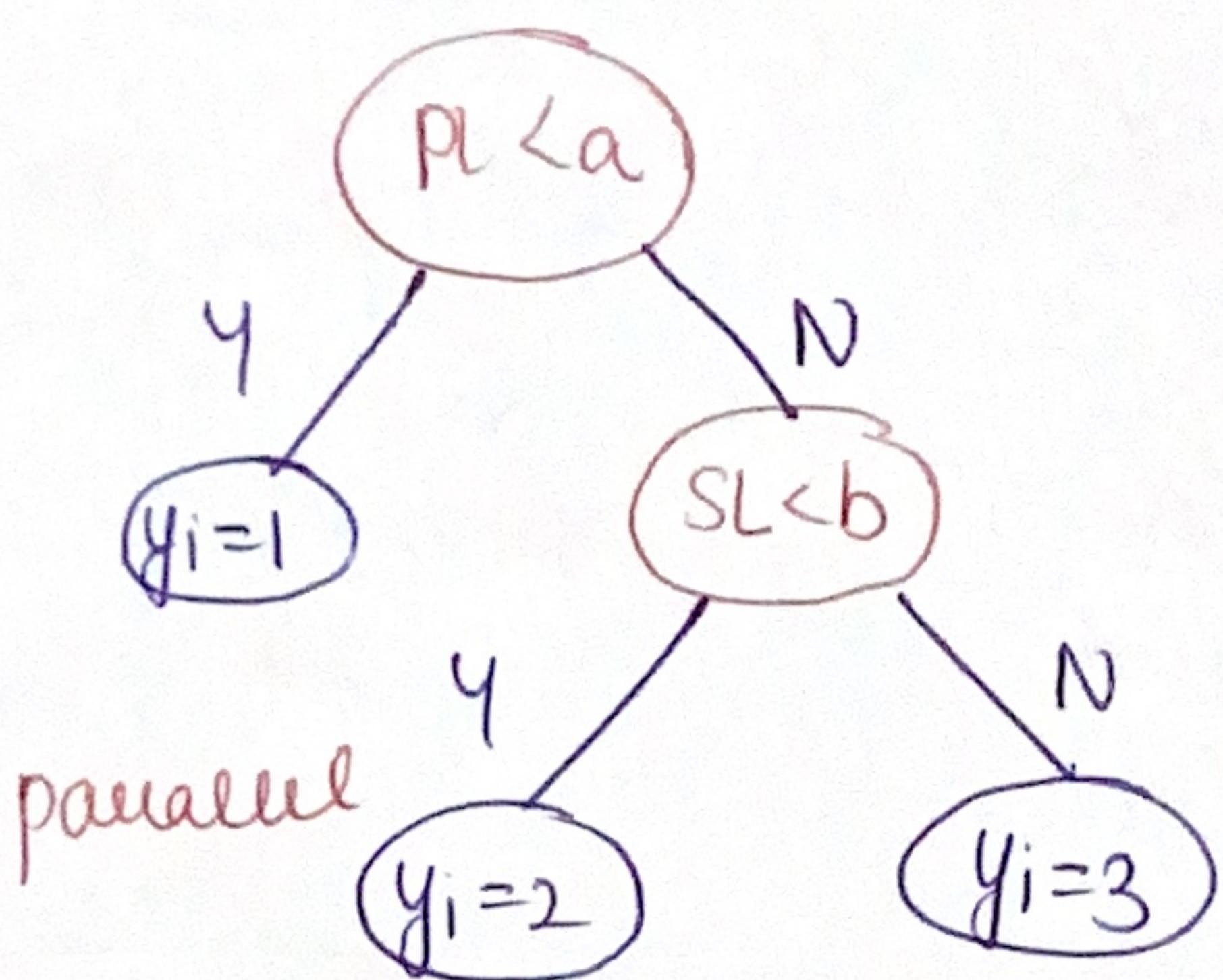
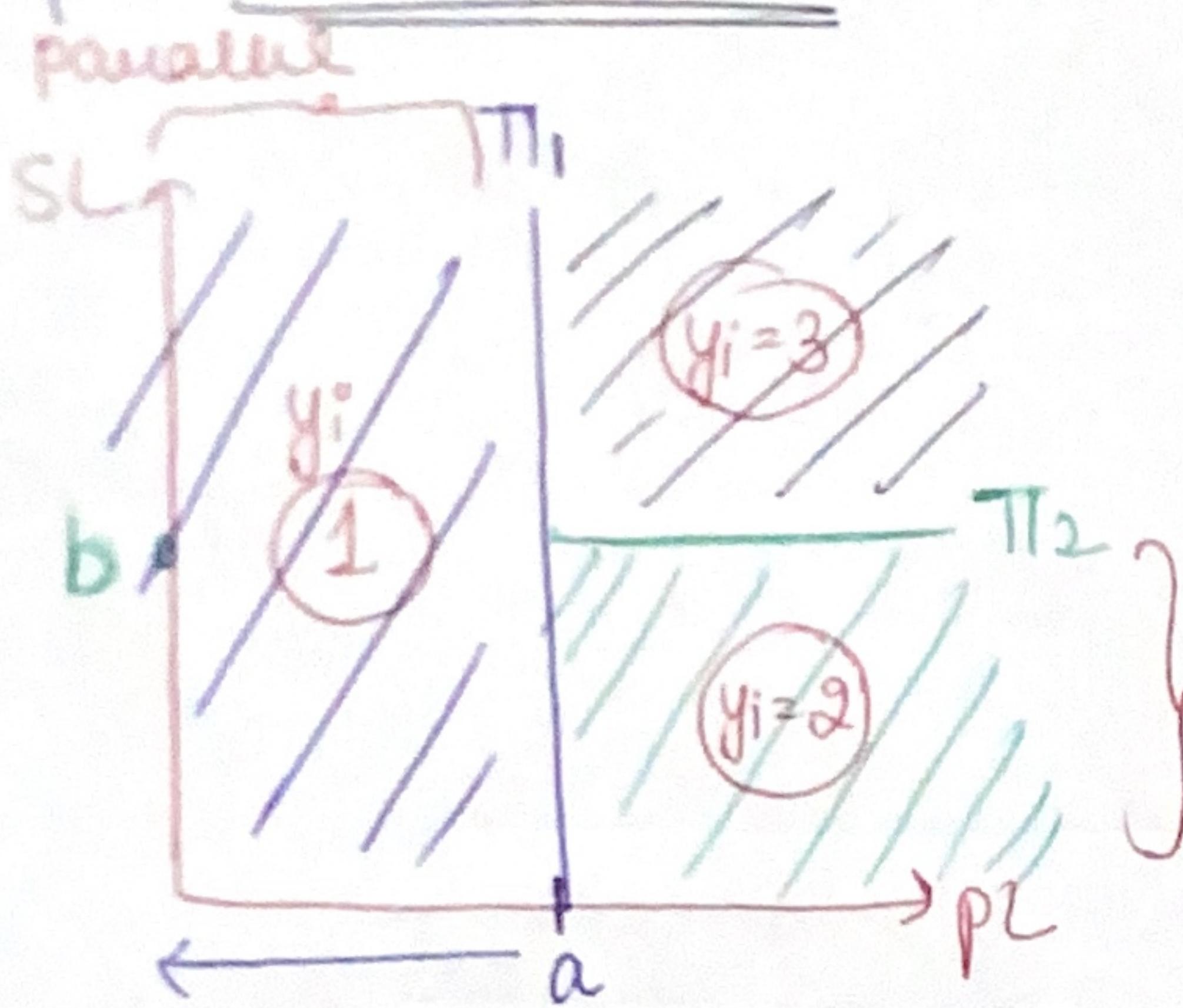
$x_i = \langle SL, SW, PL, PW \rangle$
 $\left[\begin{array}{l} \text{if } PL < a \\ \quad y_i = 1 \\ \text{else if } SL < b \\ \quad \text{class} = 2 \\ \text{else} \\ \quad \text{class} = 3 \end{array} \right]$ \rightarrow Nested if else conditions

Diagram
(Tree)

(Tree)



Geometric Institution D.T: set of axis parallel hyperplanes



decision₁ :- H_1 ,

decision₂ :- H_2

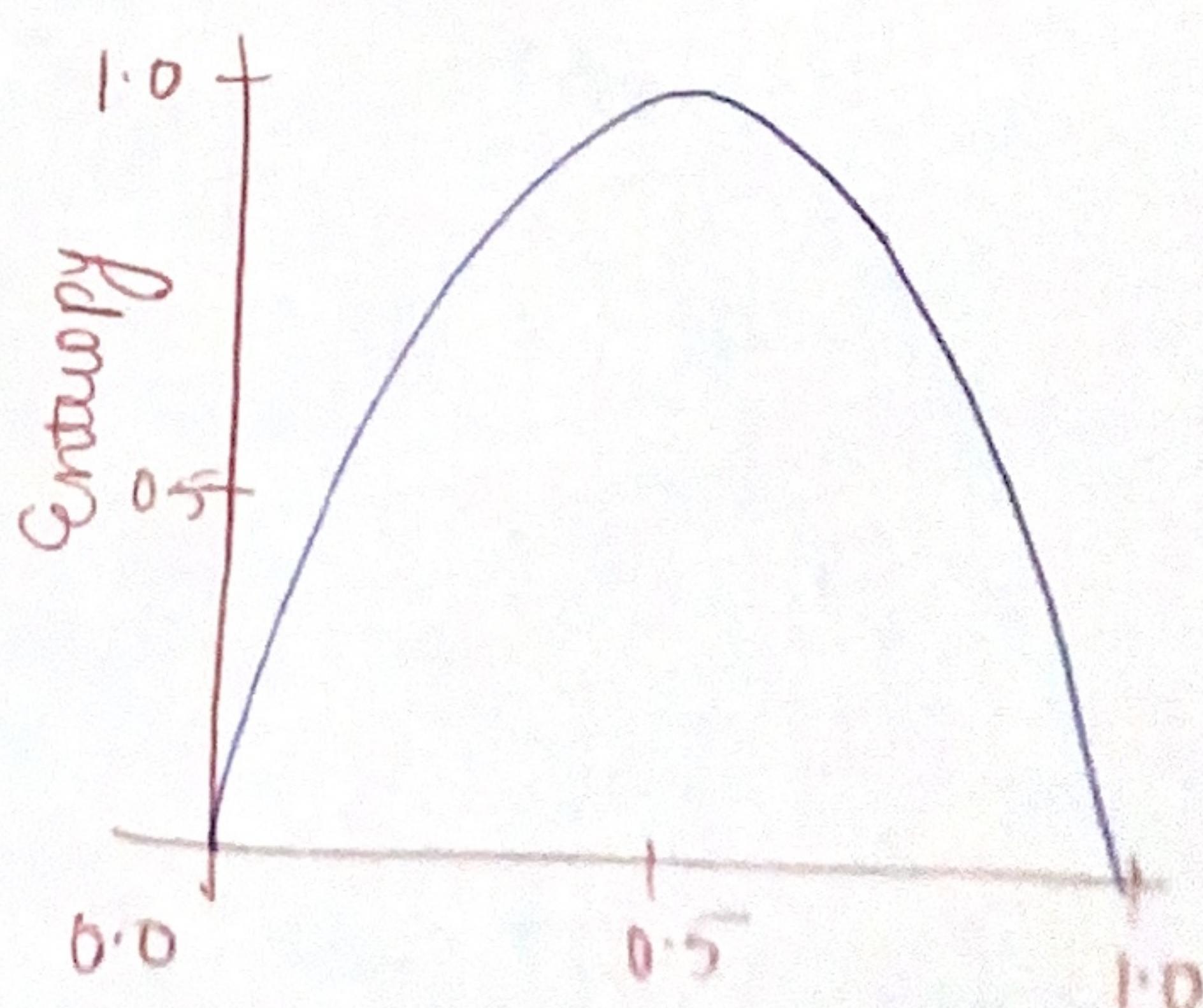
all of your hyperplanes are axis parallel
|| to any one axis.

D.T :- Nested if else \rightarrow programmatic
axis parallel hyperplanes \rightarrow geometric

Entropy

The key to decision tree induction is notion of entropy

Entropy = measure of randomness



Obs: Entropy is max if we have 50-50% of split among +ve and -ve values.

Obs: Entropy is zero if we have all +ve or all -ve examples

Entropy

- Let $D = \{(\bar{x}_1, y_1), \dots, (\bar{x}_l, y_l)\} \subseteq A^n \times \{+1, -1\}^l$

A dataset D of l labelled examples. Each eg has features $\bar{x} \in A^n$ and binary label $y \in \{+1, -1\}$

- l_+ = number of positive examples in D

l_- = number of -ve " "

$$\text{Total } l = l_+ + l_-$$

$p_+ = \frac{l_+}{l}$ $p_- = \frac{l_-}{l}$ are class probabilities .

$$\text{Entropy}(D) = -\frac{l_+}{l} \log_2\left(\frac{l_+}{l}\right) - \frac{l_-}{l} \log_2\left(\frac{l_-}{l}\right)$$

$$D = -p_+ \log_2(p_+) - p_- \log_2(p_-)$$

Building a DT: Entropy

$D_{Train} \rightarrow DT$

$\text{M.V } Y \rightarrow y_1, y_2, y_3, \dots, y_k$

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

$b=2 \checkmark$
or
 $b=e = 2.718$

$p(y_i)$

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

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

$\frac{\# \text{ wpts}}{\# \text{ wpts}}$

\uparrow
 $p(y_+)$

\uparrow
 $p(y_-)$

$\% \text{ age of wpts}$
in B

Properties: $Y \rightarrow y_+, y_-$ (2 classes)

Case 1 $\emptyset \rightarrow y_+ \rightarrow 99\%$

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

$\rightarrow y_- \rightarrow 1\%$

Case 2 $\emptyset \rightarrow y_+ \rightarrow 50\%$

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

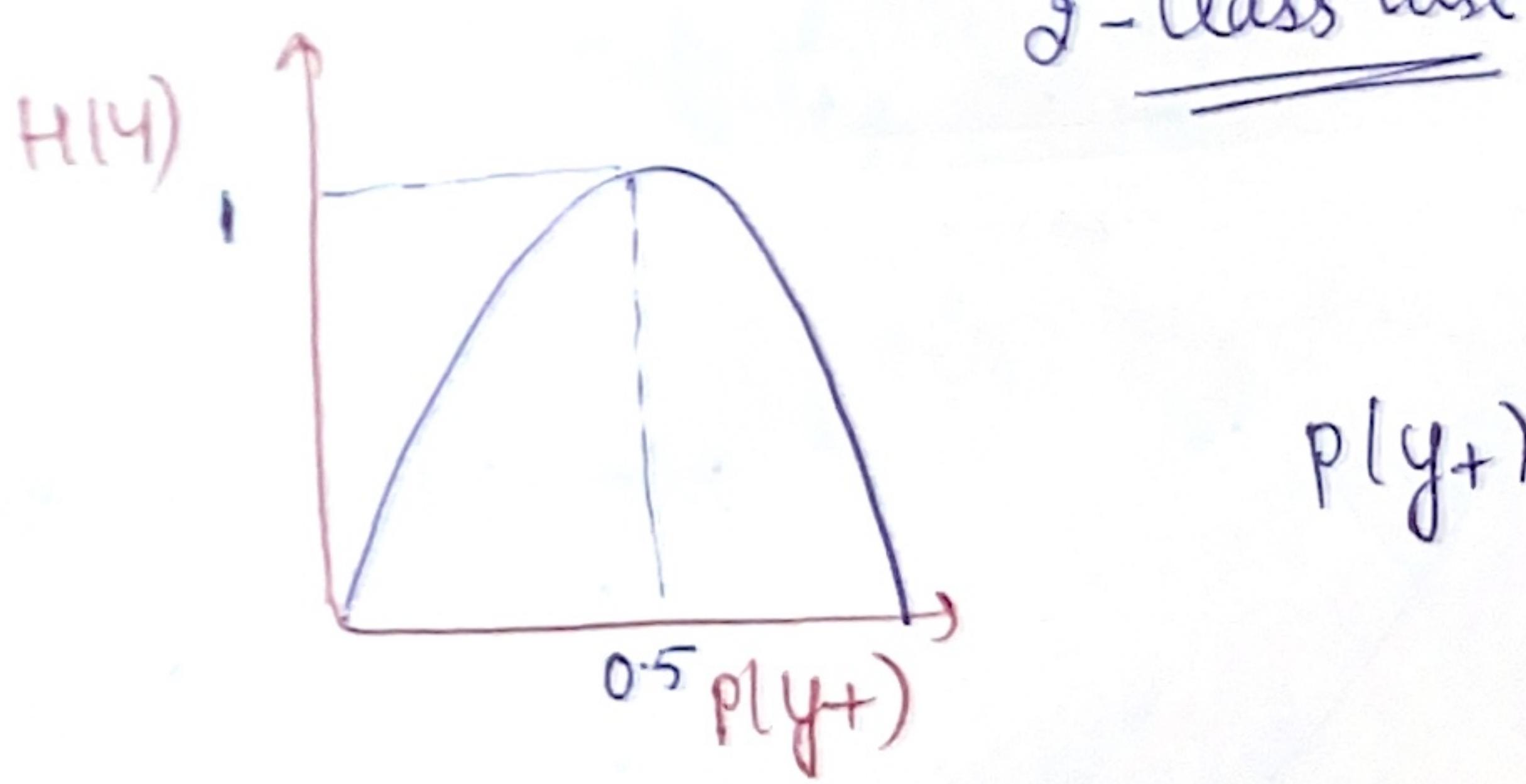
$\rightarrow y_- \rightarrow 50\%$

Case 3

$y_+ \rightarrow 0\%$

$$H(Y) = 0$$

$y_- \rightarrow 100\%$



$$p(y+) = 1 - p(y_i)$$

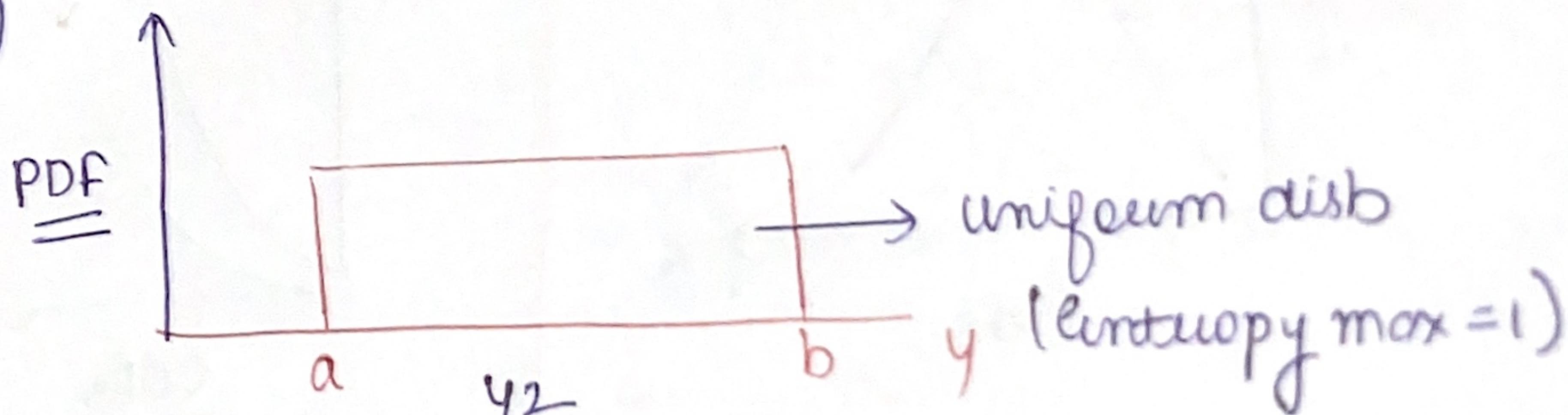
$Y \rightarrow y_1, y_2, \dots, y_k$

equi-probability \rightarrow entropy is maximum

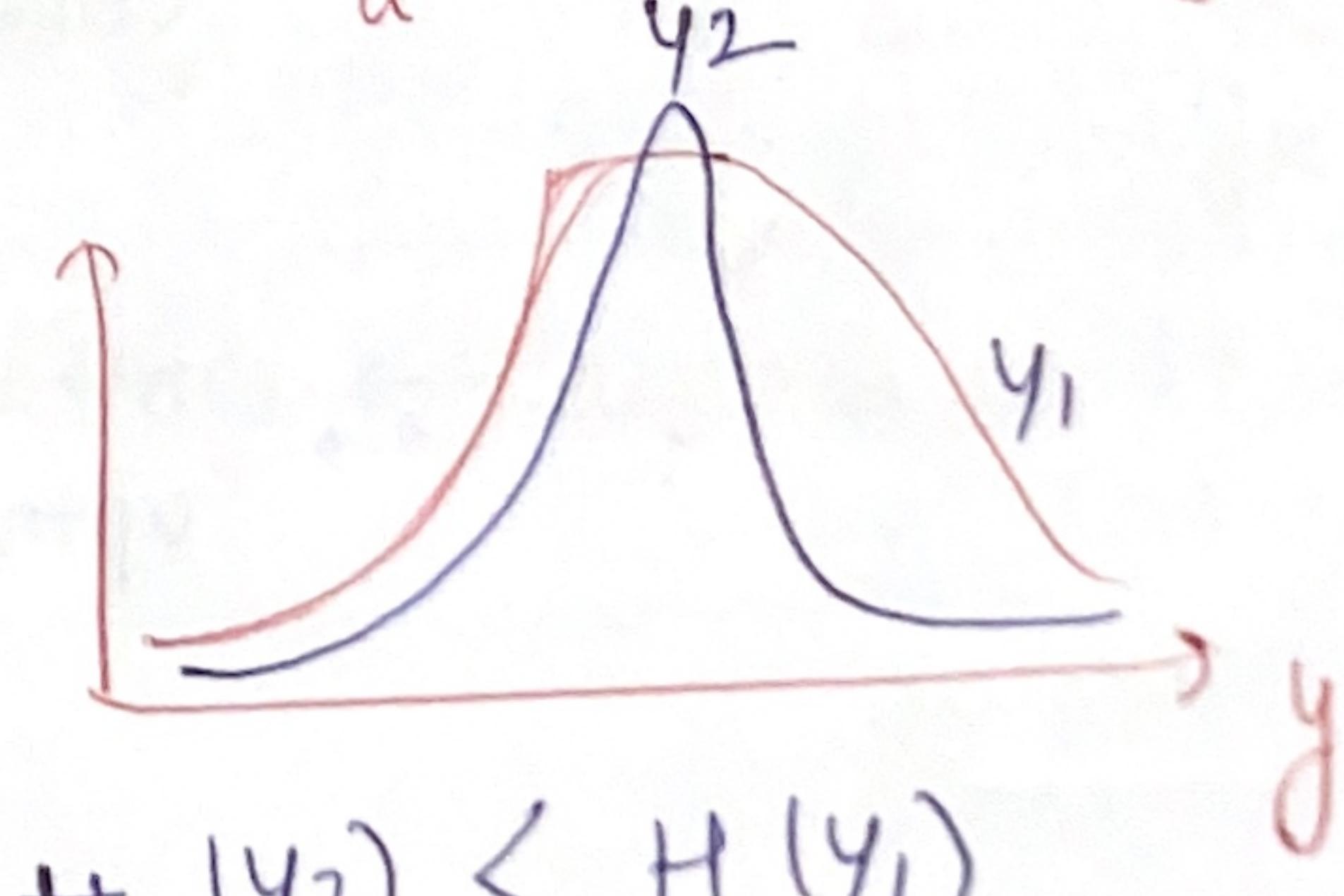
$y_1 \rightarrow$ most probable } entropy is minimum
 $y_2, y_3, \dots \rightarrow 0$

0.94

(PDF)



pdf



$y_1 \rightarrow$ less peaked
 $y_2 \rightarrow$ peaked

$$\#(y_2) < H(y_1)$$

Kullback - Leibler (KL) Divergence



PDFs / PMFs

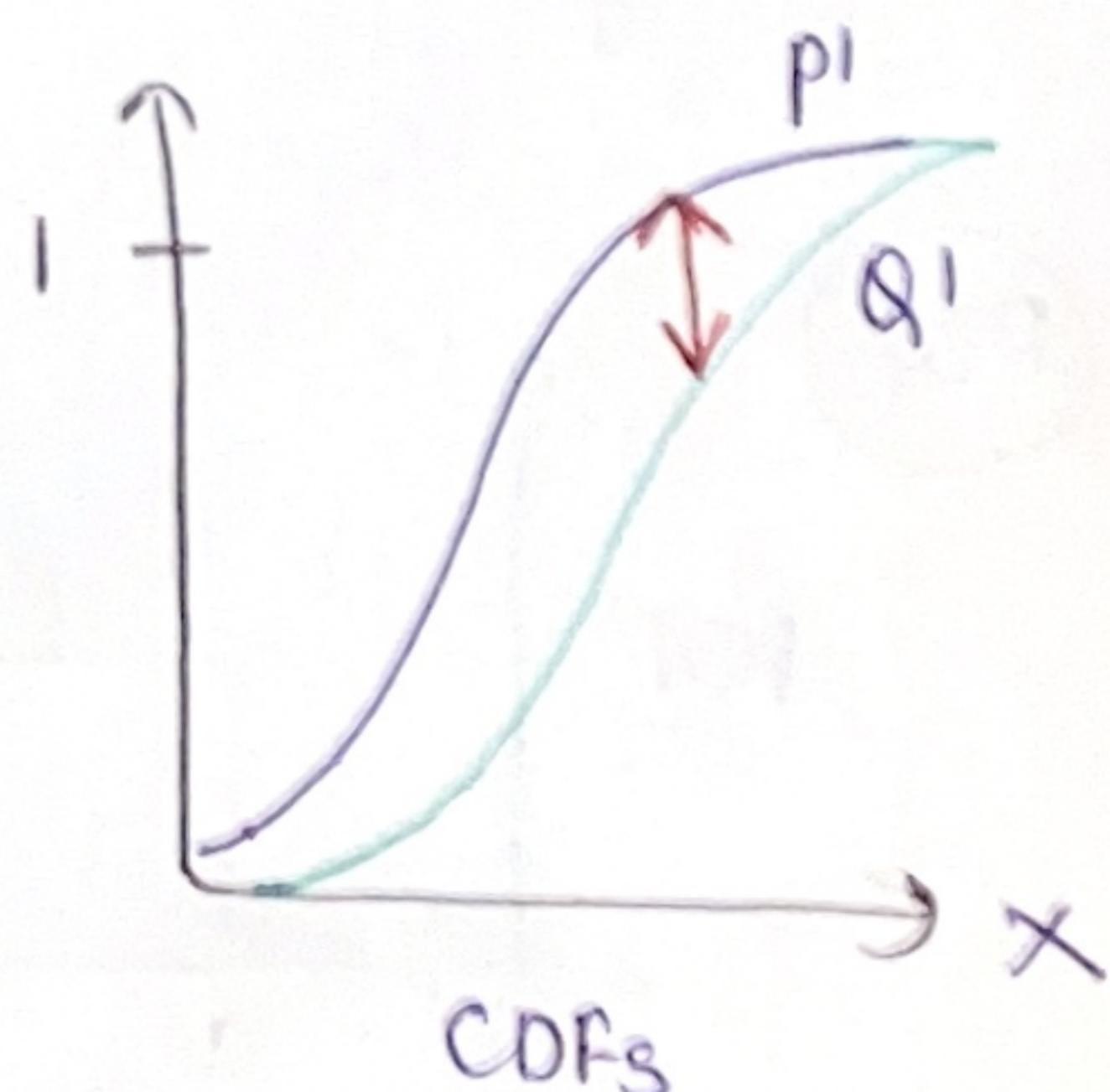
$\text{dist}(P, Q) \rightarrow$ Small
 \downarrow \rightarrow Large
 must be differentiable to
 use as loss function

One Idea: KS-Statistic = $\sup \{ |P^*(x) - Q^*(x)| \}$



PDFs

$\text{dist}(P, Q)$



CDFs

derivative \rightarrow can't use in
 optimization prob

Info-Theoretic measure

$$D_{KL}(P||Q) = \sum_n p(n) \log \left(\frac{p(n)}{q(n)} \right)$$

← $D_{KL}(P||Q)$
 Info
 Divergence
 discrete
 \leftrightarrow Σ $p(n)$
 or
 $\int p(x) \log \left(\frac{p(x)}{q(x)} \right)$
 continuous
 \leftrightarrow \int

a.k.a. relative entropy

$$D_{KL}(P||Q) = \sum_n P(n) \log \left(\frac{P(n)}{Q(n)} \right)$$

$\log \left(\frac{a}{b} \right)$
 $= \log a - \log b$

= $\sum_n P(n) \log P(n) - \sum_n P(n) \log Q(n)$

diff ↘

Information Gain (IG)

$$IG(Y, \text{outlook}) = \underbrace{\left(\frac{5}{14} \times 0.97 \right) + \left(\cancel{\frac{3}{14} \times 0} \right) + \left(\frac{5}{14} \times 0.97 \right)}_{\downarrow \text{weighted entropy after } D_1, D_2, D_3} - 0.94$$
$$= \left(\frac{5}{7} \times 0.97 \right) - 0.94 = 1G$$

$$\textcircled{D} Y \xrightarrow{\text{var}} \overset{D_1}{y_1}, \overset{D_2}{y_2}, \dots, \overset{D_K}{y_{1C}}$$

$$\boxed{IG(Y, \text{var}) = \sum_{i=1}^K \frac{|D_i|}{|D|} \times H_{D_i}(Y) - H_D(Y)}$$

Gini Impurity \curvearrowleft Similar to Entropy

$I_g(Y) \neq I_g(Y)$

Gini Impurity

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

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

$$Y \rightarrow y_+ \\ y_-$$

Case 2 $p(y_+) = 1$

$$p(y_-) = 0$$

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

$$H(Y) = 0$$

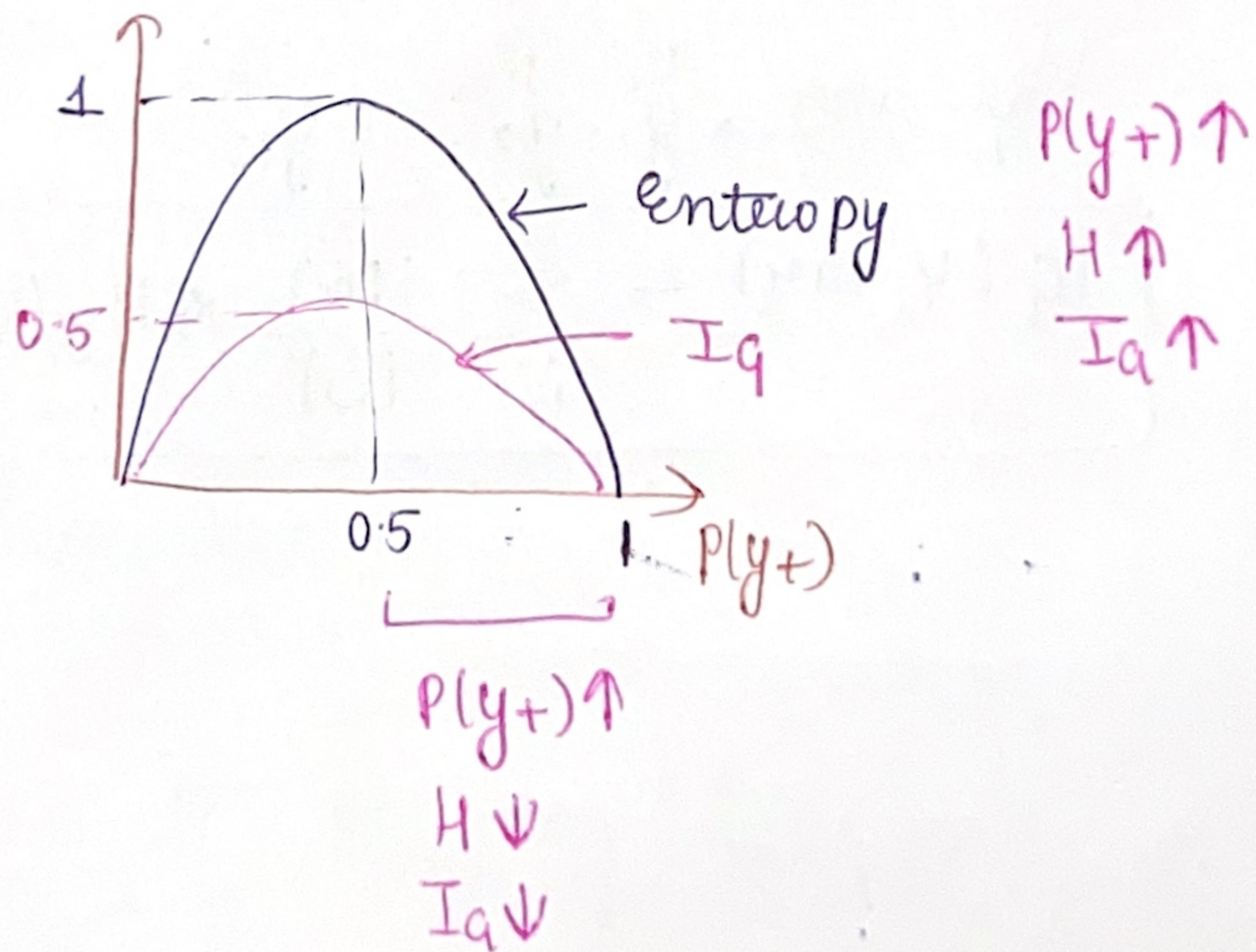
Case 1 $p(y_+) = 0.5$

$$p(y_-) = 0.5$$

$$I_g(Y) = 1 - (0.25 + 0.0025) \\ = 0.5$$

$$H(Y) = 1$$

2 - Category Case y_+, y_- $p(y_+) = 1 - p(y_-)$



$I_q(4)$

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



no log

more computationally easy
and efficient to calculate

$H(4)$

$$-P(y_+) \log_2 P(y_+) - P(y_-) \log_2 (P(y_-))$$

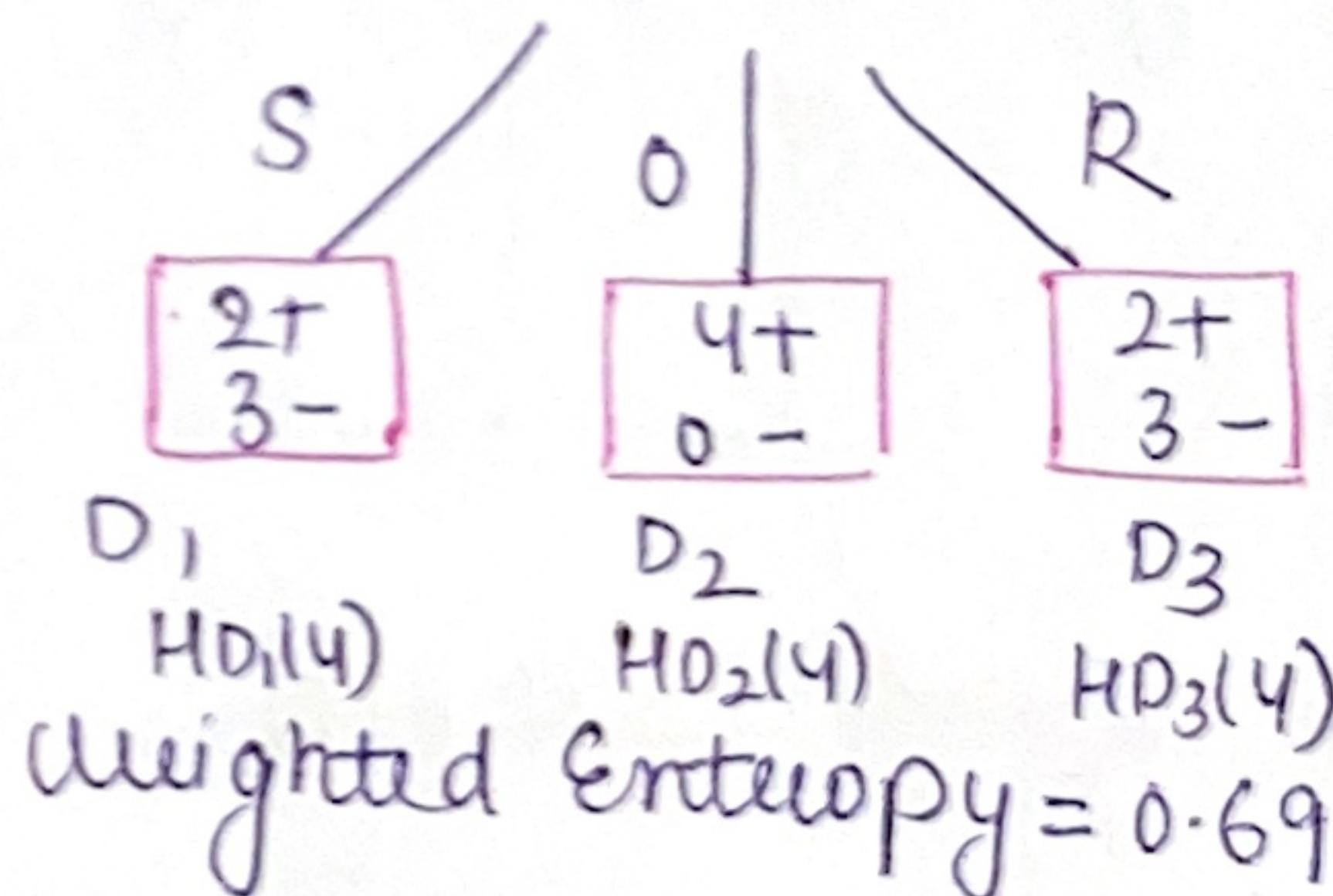
↓
log

Construct a DT: H, I_G, I_q

① $D \rightarrow q^+, 5^-$

$$H_D(4) = 0.94$$

② $\xrightarrow{D=0.94}$ outlook $0.94 - 0.69$



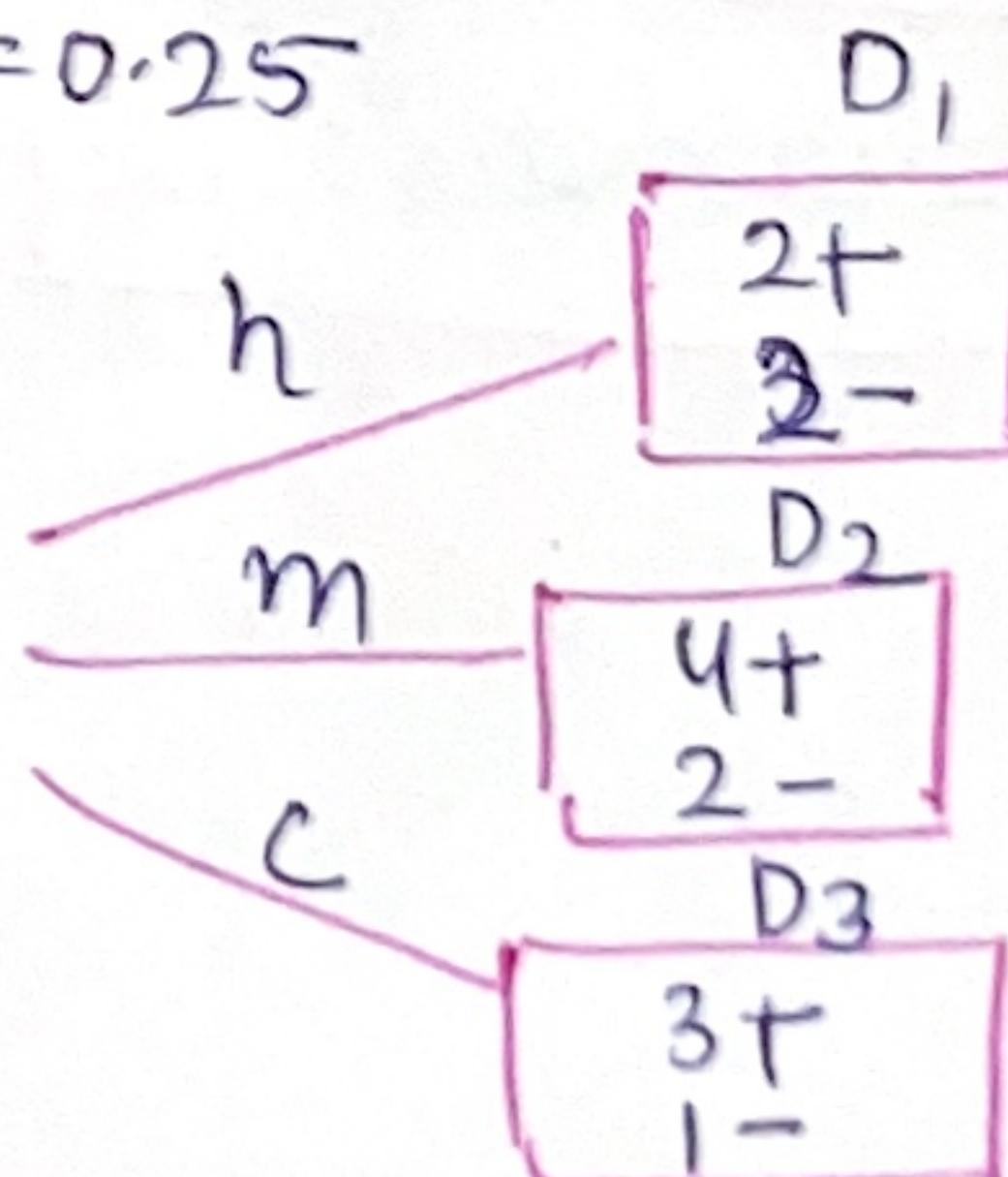
$$D_1, H_{D_1}(4) \quad D_2, H_{D_2}(4) \quad D_3, H_{D_3}(4)$$

$$\text{Weighted Entropy} = 0.69$$

$$0.94 - 0.69 = 0.25$$

$$I_G(4, \text{outlook}) = 0.25$$

③ Temperature



$$I_G(4, \text{Temp}) =$$

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

$$\left\{ \begin{array}{l} IG(Y, \text{outlook}) = 0.85 \\ IG(Y, \text{temp}) = \\ IG(Y, \text{Humidity}) = \\ IG(Y, \text{windy}) = \end{array} \right.$$

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

- ① pure node - Stop growing the node
- ② can't grow the tree anymore because of lack of pts.
- ③ if we are too deep.

depth of tree $\uparrow \Rightarrow$ overfit \uparrow (few pts)

depth is small \Rightarrow underfit

DT: hyperparameter: depth

CV

Splitting numerical features

Construct a DT: splitting a node \rightarrow IG

IG: - Entropy

gini impurity \rightarrow computationally efficient

Discussed so far on categorical features we have seen till now

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

\rightarrow Split based on cat. value was easy

f_2 : 3 categories



① Sort the numerical feature in ascen order

②

$$f_1 < 2.2$$

$$f_1 < 2.6$$

$$f_1 < 3.5$$

$$f_1 < 4.6$$

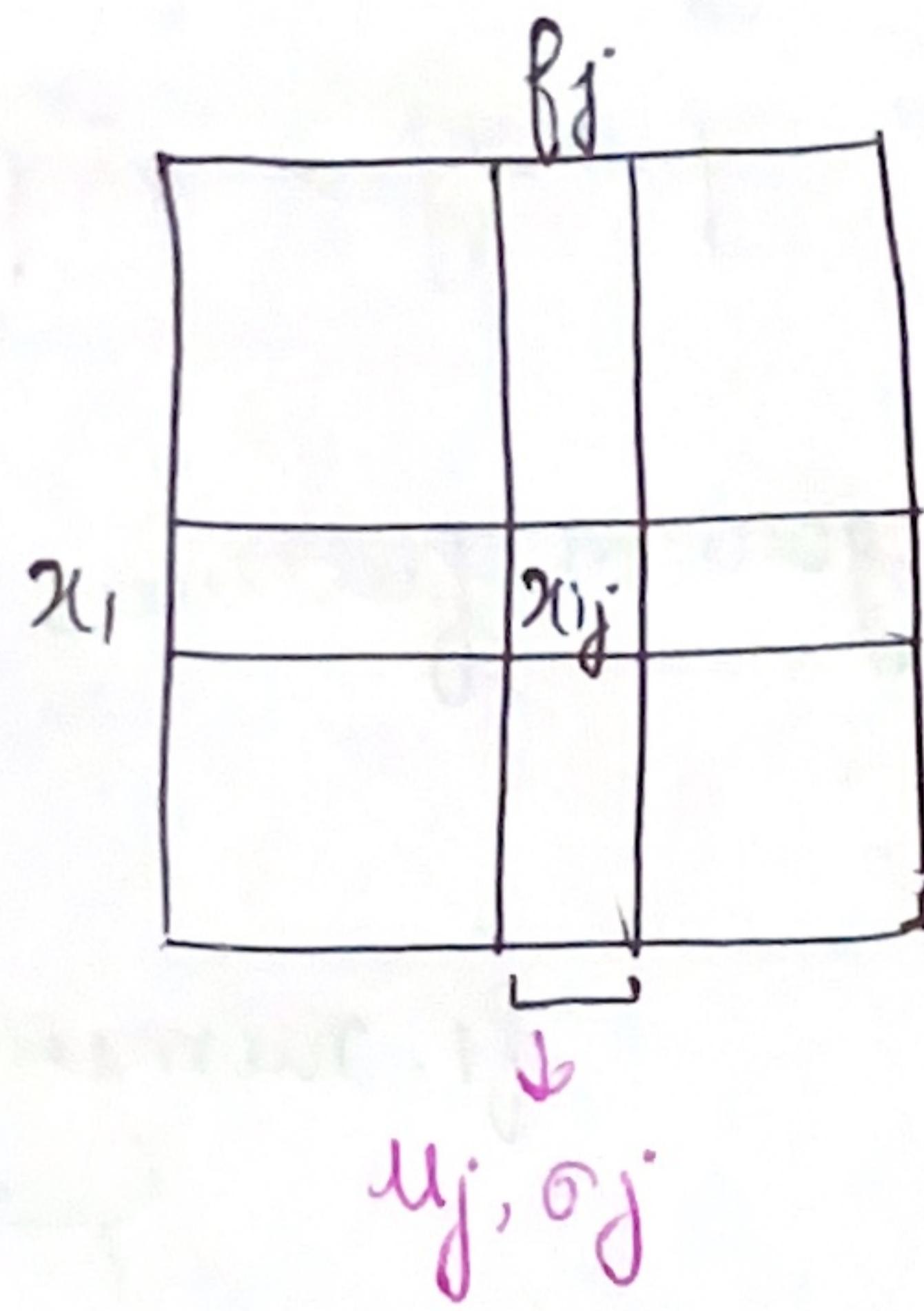
$$f_1 < 5.3$$

Numerical features

$$f_1 < c_1 \begin{cases} D_1 \\ D_2 \end{cases} f_2 < c_2 \begin{cases} D_1 \\ D_2 \end{cases} f_3 < c_3 \begin{cases} D_1 \\ D_2 \end{cases}$$

Feature Standardization

{ logistic reg :
SVM
KNN } → feature stdn

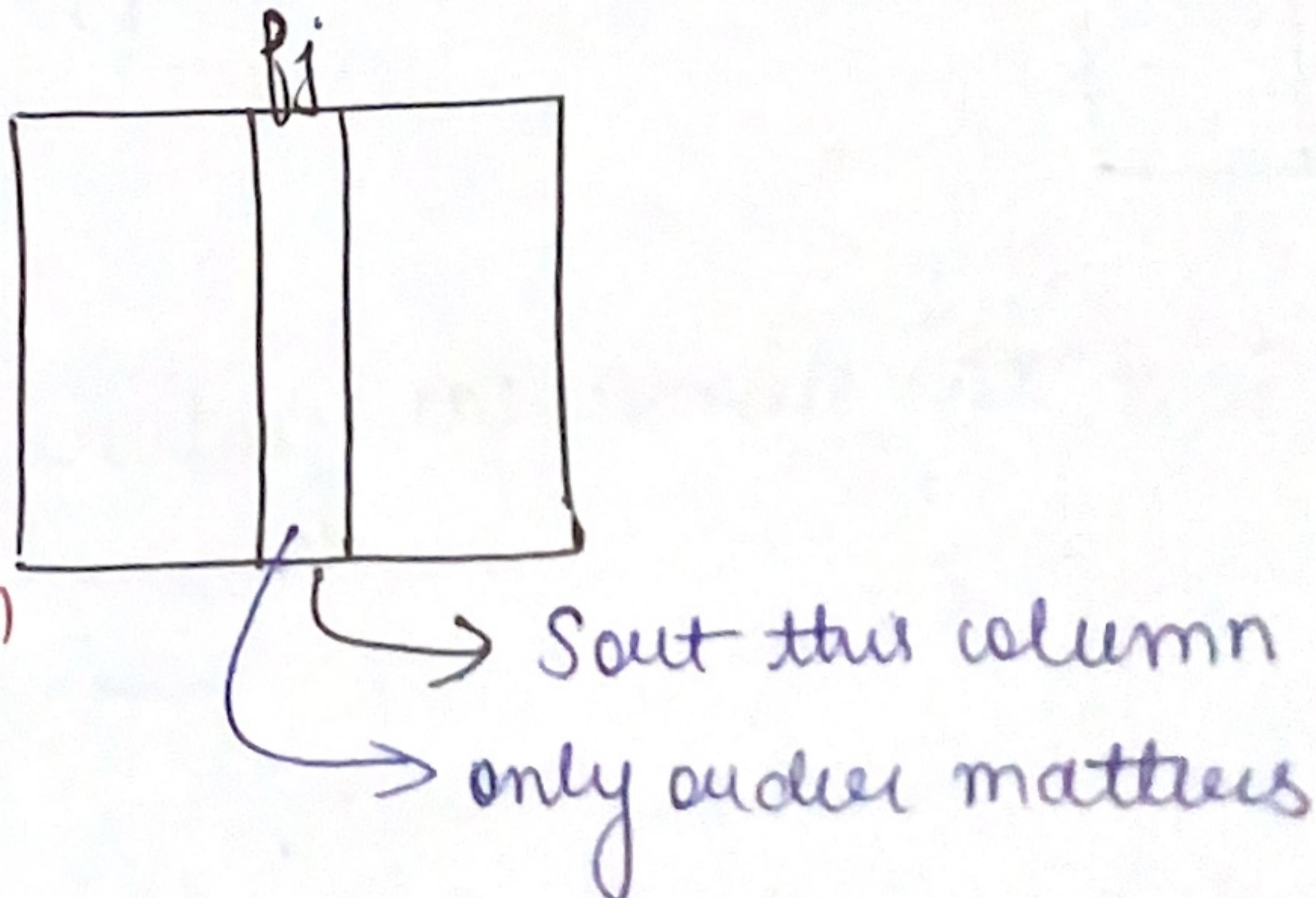


$$x'_{ij} = \frac{x_{ij} - \mu_j}{\sigma_j}$$

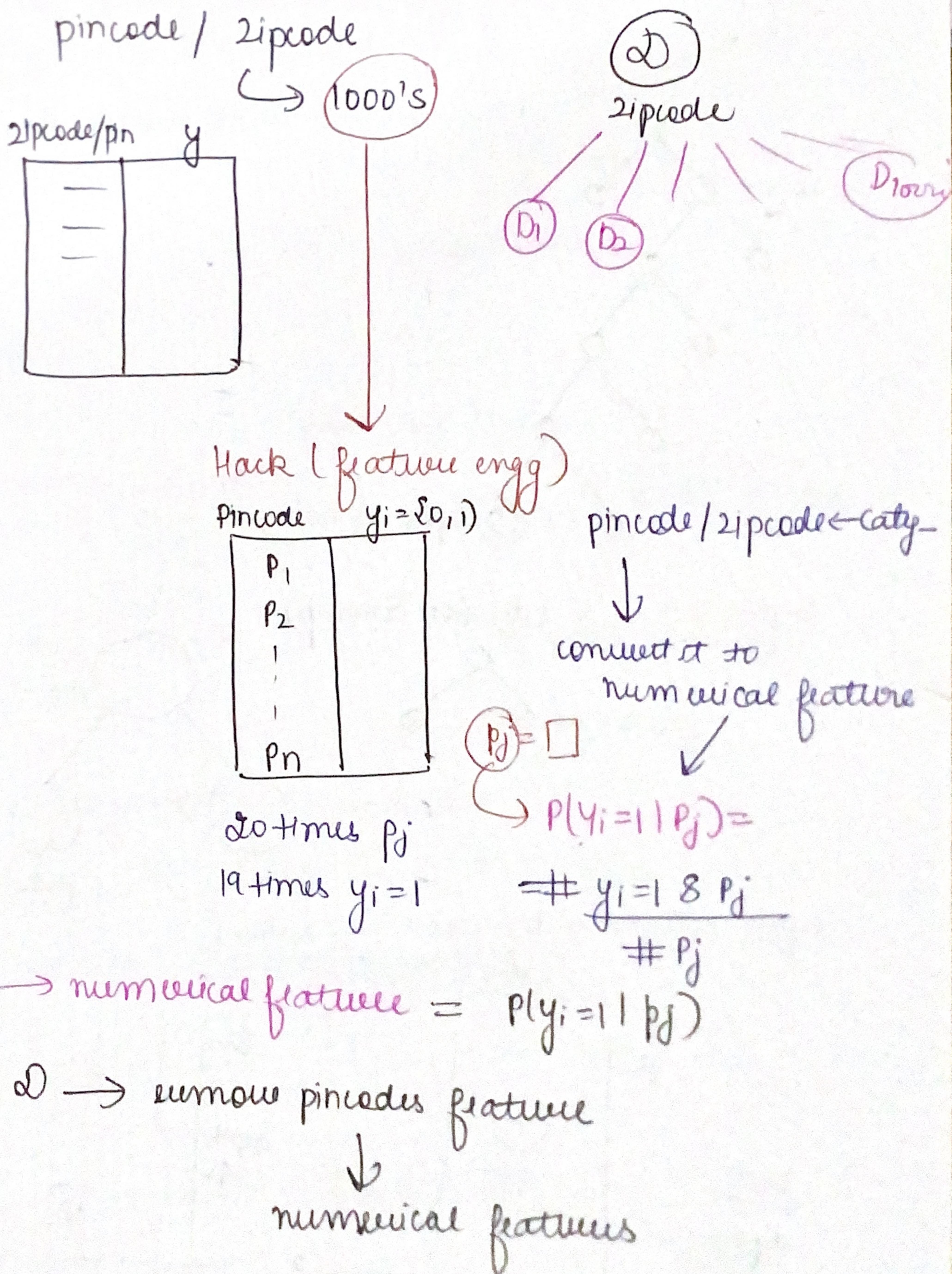
Decision Trees

not a distance based method

{ do not need
to perform
feature
standardization }

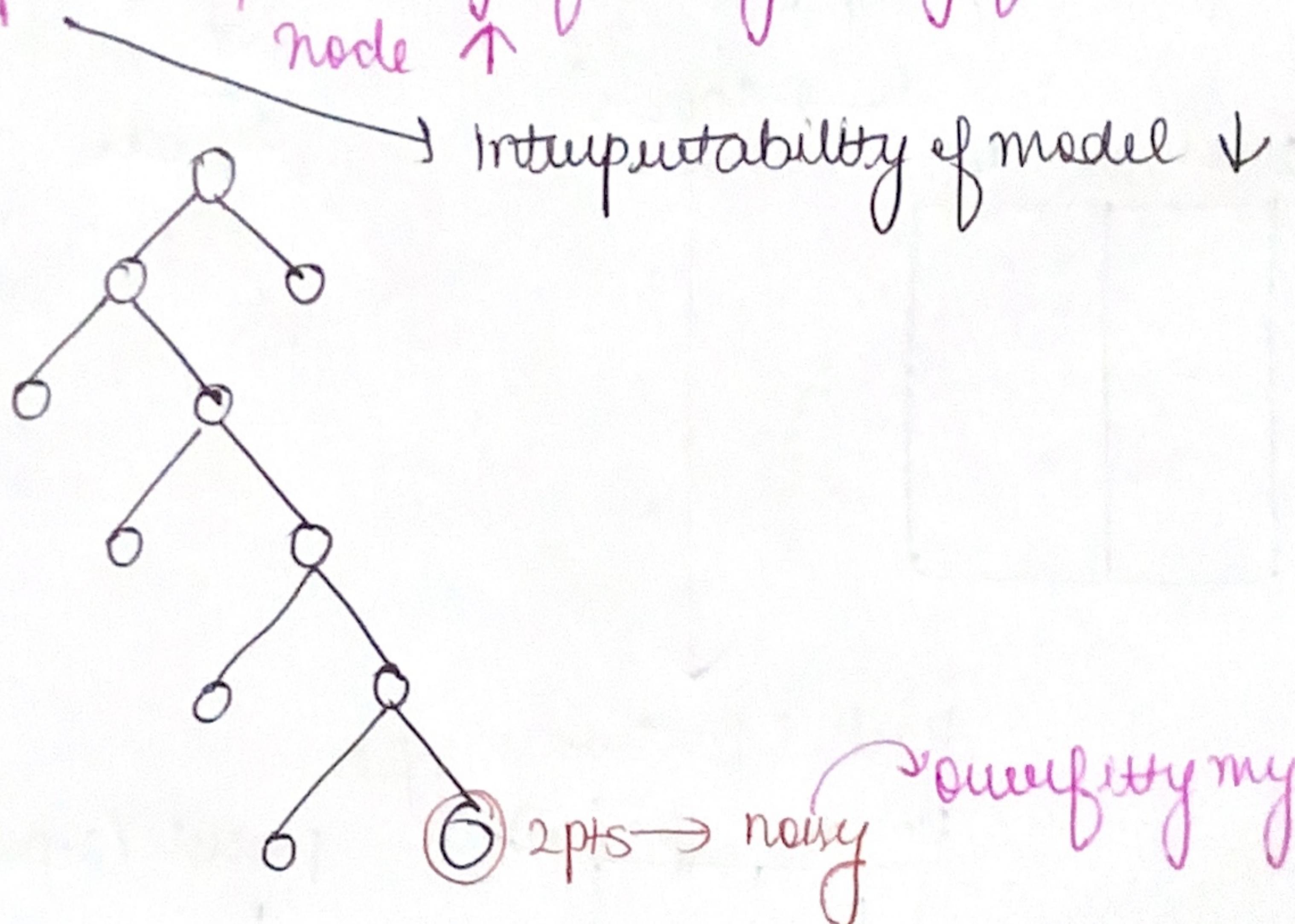


Categorical features with many categories
↳ levels.

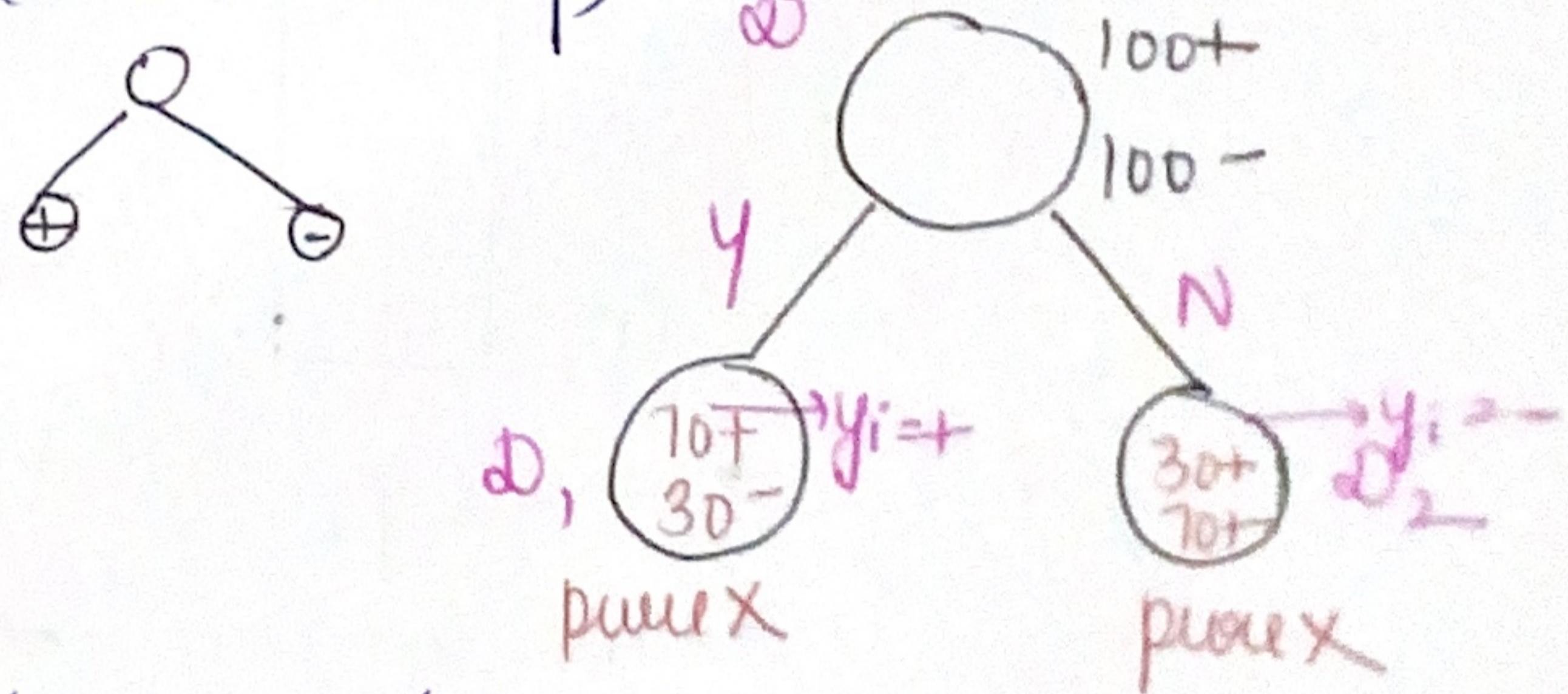


Overshifting & Undershifting

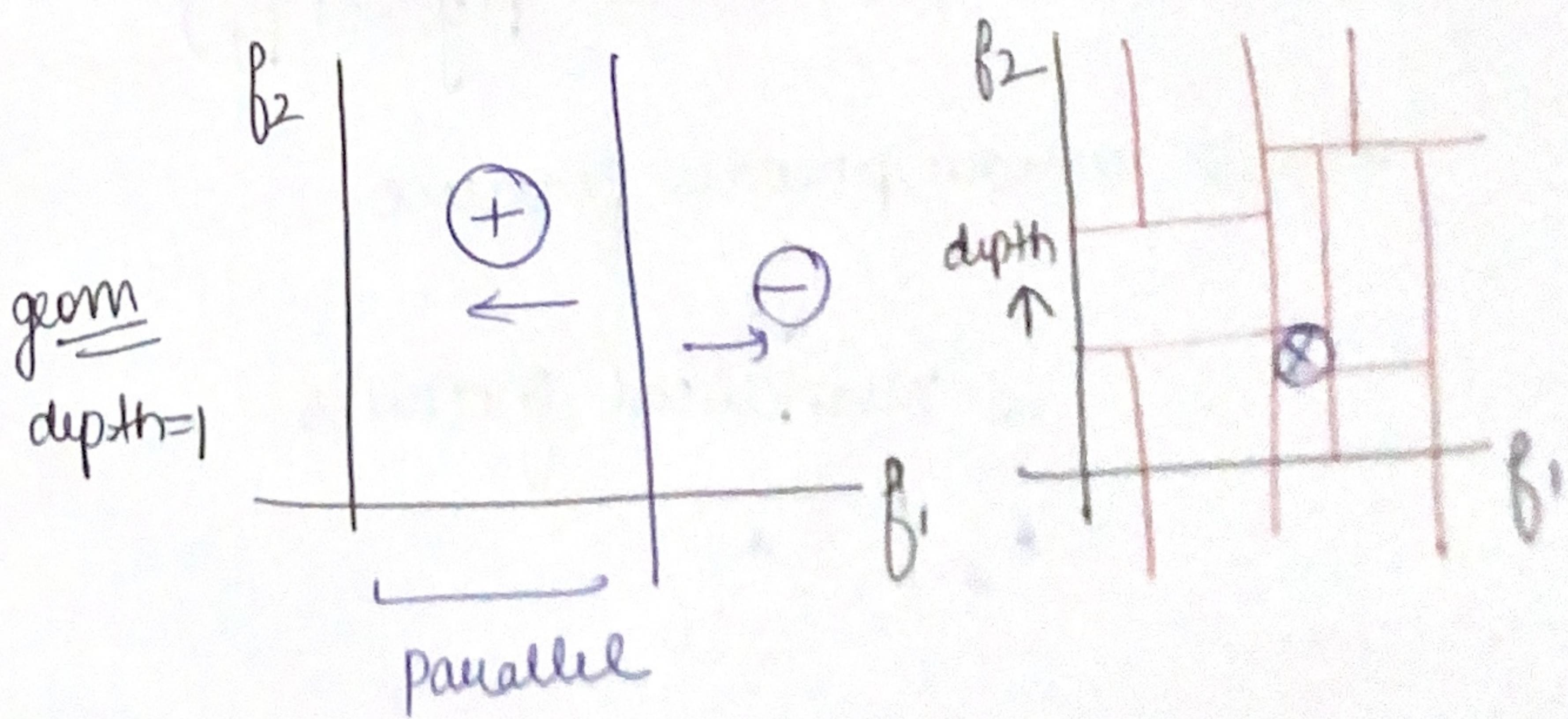
depth ↑: possibility of having very few pts at leaf
node ↑



depth ↓ := 1 = (decision Stump)



depth determined by Cross Validation



Train 8 Run Time Complx

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

numerical features: (threshold)
↳ algorithmic methods

After Training:

at Runtime Space

$x_g \rightarrow y_g$

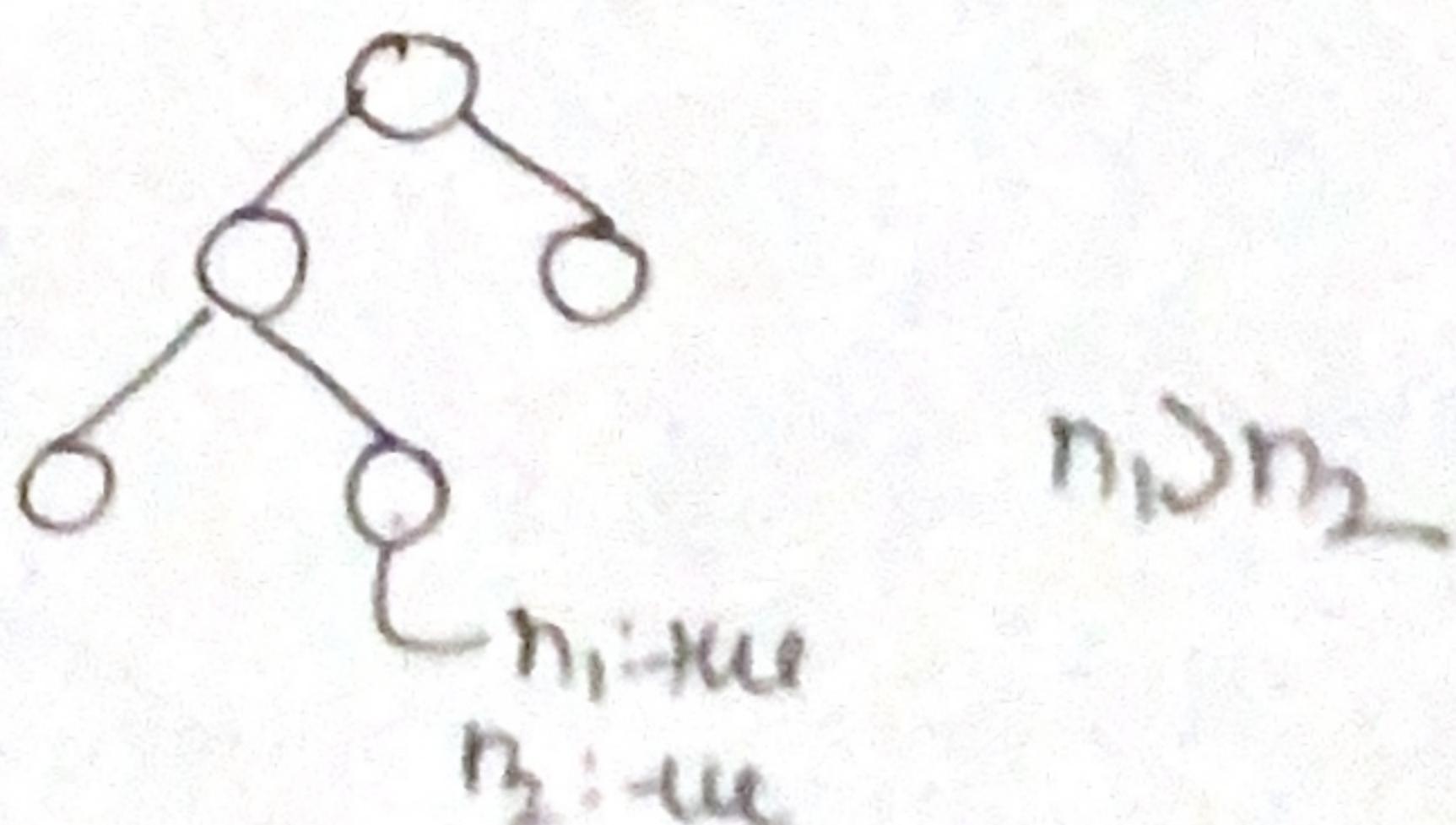
Store my Rule

{ ↳ if else } nested if else
[# internal nodes
+
leaf nodes

DT: large data, dim is small
low latency $\rightarrow O(\text{depth})$

Regression Using Decision Trees:

DT \rightarrow classfn



IG \rightarrow Clasfn only, Regression X

MSE or MAE \rightarrow Regression

Mean Squared Error (MSE)

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

Entropy

$$g_i = \text{mean}_{D_i}(y_i)$$

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

aug = mean/median

weighted sum of MSE

$$(w_1 \text{MSE}_{D_1} + \text{MSE}_{D_2} w_2)$$



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

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

Clasfn \rightarrow f_i :- reducing entropy : "0"

regress \rightarrow f_i : reduce MSE \rightarrow "0"

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

$$\text{MAE} = \text{MAD} \rightarrow "0"$$

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