

# Decision Tree → Classification & Regression

Healthcare  
Recall

threshold →

0.2 (scenario)

FN reduce  
ROC-AUC

threshold →

0.8 (scenario)

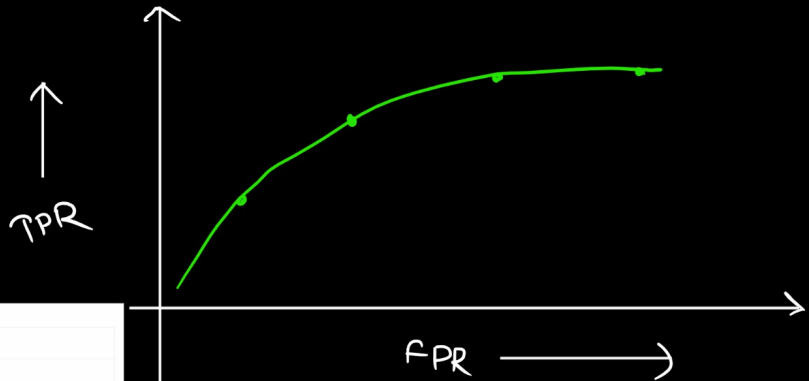
FP reduce

$$FPR = \frac{FP}{FP + TN}$$

$$TPR = \frac{TP}{TP + FN}$$

widely  
asked  
interview  
question

	Predicted Spam (1)	Predicted Not Spam (0)
Actual Spam (1)	40 (True Positives)	10 (False Negatives)
Actual Not Spam (0)	5 (False Positives)	45 (True Negatives)



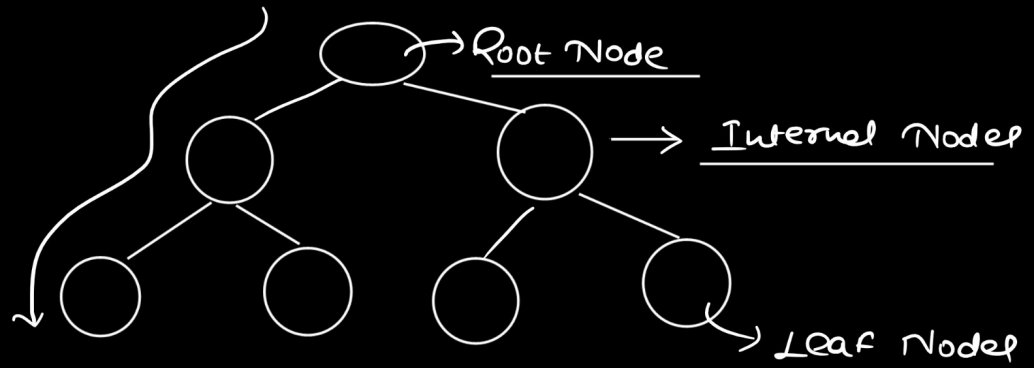
Threshold	TPR (Recall)	FPR
0.2 ✓	0.95	0.40
0.4 ✓	0.90	0.20
0.6 ✓	0.85	0.10
0.8 ✓	0.75	0.05
1.0 ✓	0.00	0.00

0.2 > — 1  
0.2 ≤ — 0

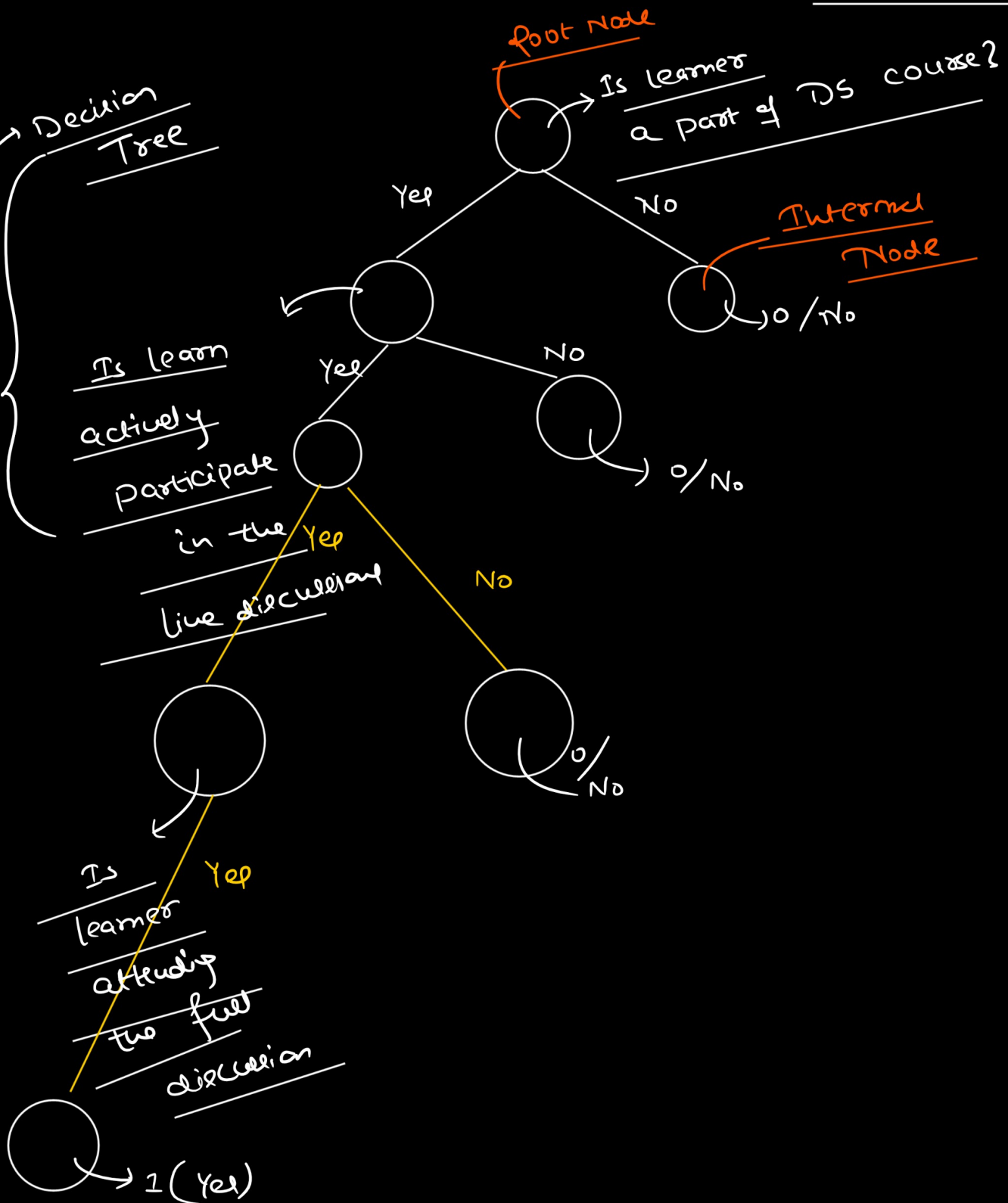
CM

Prob  
0.5 > —→ 1  
0.5 ≤ —→ 0

# Tree Data Structure (Non-linear)



## Decision Tree



# Gini Index & Entropy

Impurity measure

$$Gini = 1 - \sum_{i=1}^c p_i^2$$

$f_1$   $f_2$  target

Class	Gender	Stay in hostel
9	M ✓	Yes
10	F	No
8	F	Yes
8	F	No
9	M ✓	Yes
10	M	No
11	F	Yes
11	M	Yes
8	F	Yes
9	M ✓	No
11	M	No
11	M	Yes
10	F	No
10	M	Yes

$$Gini(class) = 0.404$$

$$Gini(Gender) = 0.482$$

→ weighted avg

$$Gini(class) = \frac{3}{14} * \frac{4}{9} + \frac{3}{14} * \frac{4}{9} + \frac{4}{14} * \frac{3}{8} + \frac{4}{14} * \frac{3}{8} \geq 0.404$$

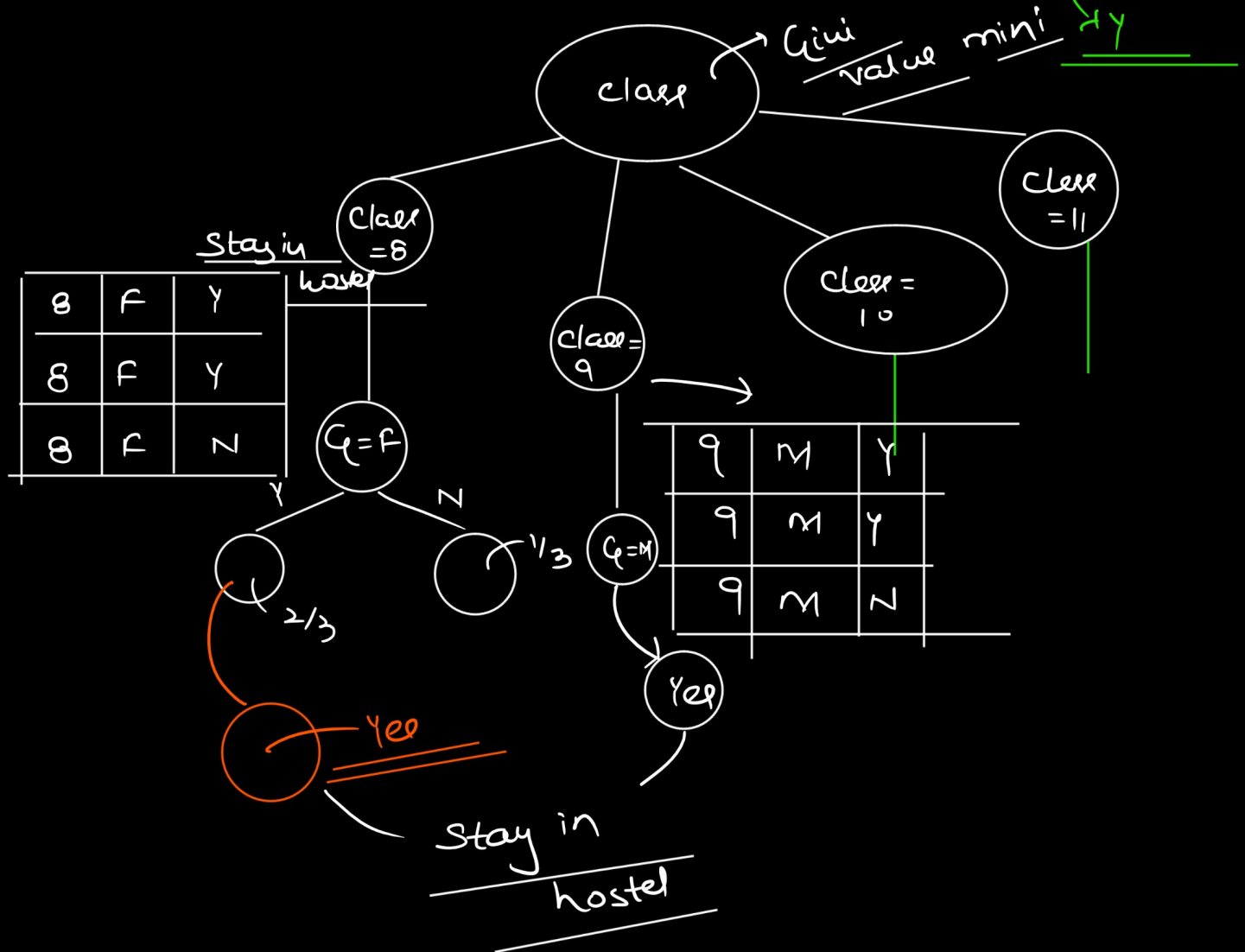
class	stay in hostel	$P(Y)$	$P(N)$	Gini
8	$Y=2$ $N=1$	$\frac{2}{3}$	$\frac{1}{3}$	$1 - \left( \left( \frac{2}{3} \right)^2 + \left( \frac{1}{3} \right)^2 \right) = \frac{4}{9}$
9	$Y=2$ $N=1$	$\frac{2}{3}$	$\frac{1}{3}$	$\frac{4}{9}$
10	$Y=1$ $N=3$	$\frac{1}{4}$	$\frac{3}{4}$	$1 - \left( \left( \frac{1}{4} \right)^2 + \left( \frac{3}{4} \right)^2 \right) = \frac{3}{8}$
11	$Y=3$ $N=1$	$\frac{3}{4}$	$\frac{1}{4}$	$\frac{3}{8}$

Gender	Stay in hostel	P(Y)	P(N)	Gini
M	Y=5 N=3	5/8	3/8	$1 - \left( \left( \frac{5}{8} \right)^2 + \left( \frac{3}{8} \right)^2 \right) = 0.468$
F	Y=3 N=3	3/6	3/6	$1 - \left( \left( \frac{3}{6} \right)^2 + \left( \frac{3}{6} \right)^2 \right) = 0.5$

$$\underline{Gini(Gender) = \frac{8}{14} * 0.468 + \frac{6}{14} * 0.5}$$

$$= 0.482$$

9 M → Yes  
8 F → Y/N



Qini

Index

Interpretation

→ [0, 0.5]

Pure | Impure

$$1 - \left( (0)^2 + (1)^2 \right)$$

$$\Rightarrow 0$$



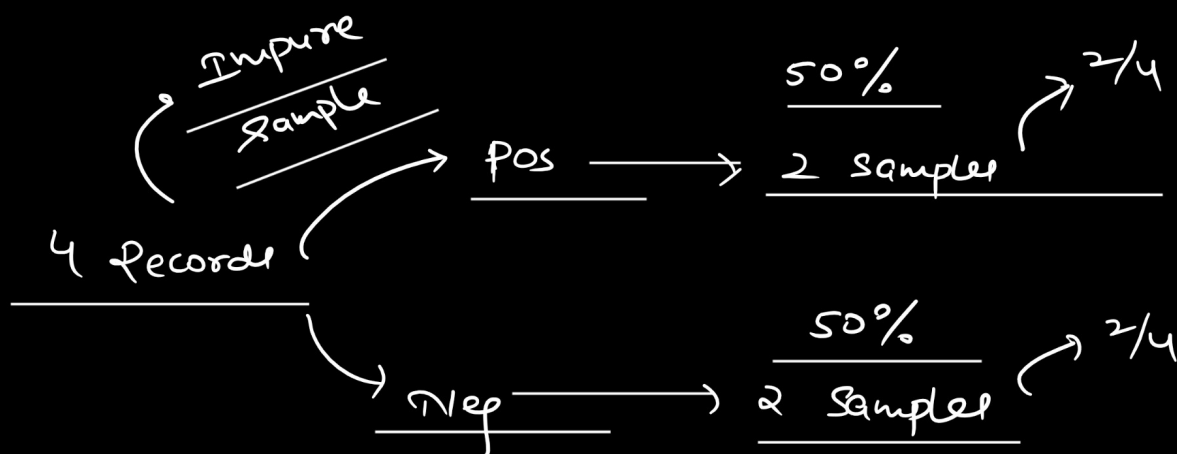
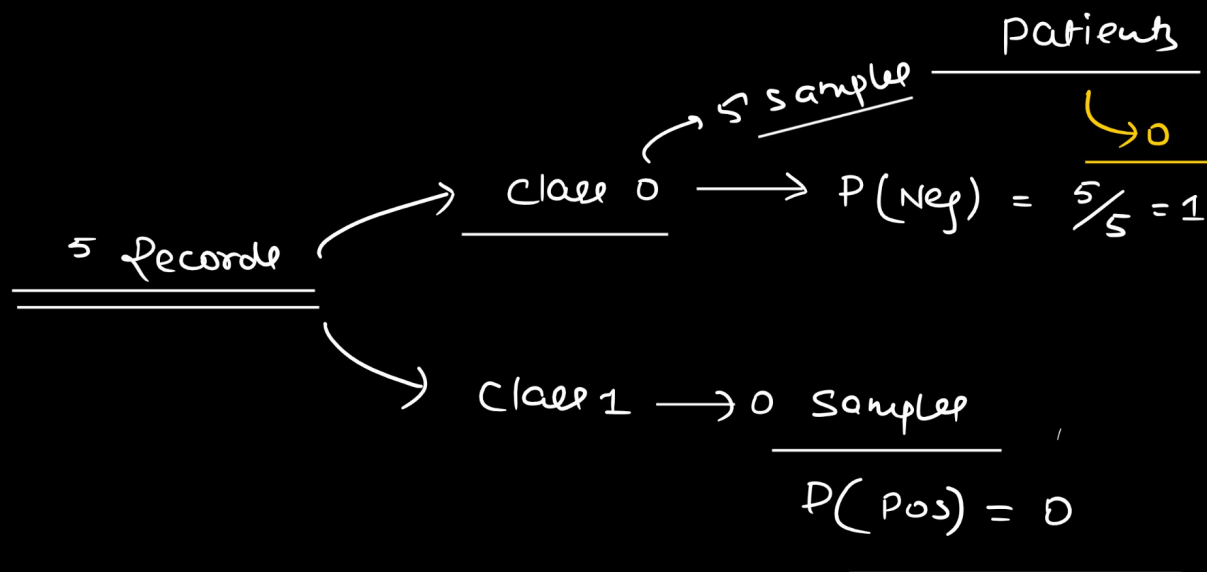
Highest level  
of purity → 1

only diabetic patients or only non-diabetic

$$1 - \left( \left( \frac{2}{4} \right)^2 + \left( \frac{2}{4} \right)^2 \right)$$

$$\Rightarrow 0.5$$

→ Highest level of  
impurity



Note →

① 0 means perfectly pure

② Closer to 1 means very impure

By default, Gini is the criterion

→ Computationally, it's very  
easy

$$\text{Entropy} \Rightarrow - \sum_{i=1}^c p_i \log_2 p_i$$

→ Information Gain

Class	Gender	Stay in hostel
9	M	Yes
10	F	No
8	F	Yes
8	F	No
9	M	Yes
10	M	No
11	F	Yes
11	M	Yes
8	F	Yes
9	M	No
11	M	No
11	M	Yes
10	F	No
10	M	Yes

$$\text{Entropy (class)} = \underline{0.857} \text{ (More Purity)}$$

$$\text{Entropy (Gender)} = \underline{0.974}$$

$$\text{Entropy (target = stay in hostel)} =$$

$$- \frac{8}{14} \log_2 \frac{8}{14} - \frac{6}{14} \log_2 \frac{6}{14}$$

$$\Rightarrow \underline{\underline{0.985}}$$

← class

$$IG \Rightarrow \text{Entropy (target)} - \text{Entropy (class)} \Rightarrow 0.985 - 0.857 = \underline{\underline{0.128}}$$

$$IG (\text{Gender}) = 0.985 - 0.974 = 0.011 \leftarrow$$

class	stay in hostel	$P(Y)$	$P(N)$	Entropy
8	$Y=2 \quad N=1$	$2/3$	$1/3$	$-2/3 \log 2/3 - 1/3 \log 1/3 = 0.918$
9	$Y=2 \quad N=1$	$2/3$	$1/3$	$0.918$
10	$Y=1 \quad N=3$	$1/4$	$3/4$	$-1/4 \log 1/4 - 3/4 \log 3/4 = 0.811$
11	$Y=3 \quad N=1$	$3/4$	$1/4$	$0.811$

$$\text{Entropy}(\text{class}) = \frac{3}{14} * 0.918 + \frac{3}{14} * 0.918 + \frac{4}{14} * 0.811 +$$

$$\frac{4}{14} * 0.811 = \underline{\underline{0.857}}$$

Gender	stay in hostel	$P(Y)$	$P(N)$	Entropy
M	$Y=5 \quad N=3$	$5/8$	$3/8$	$-5/8 \log 5/8 - 3/8 \log 3/8 = 0.954$
F	$Y=3 \quad N=3$	$3/6$	$3/6$	$-3/6 \log 3/6 - 3/6 \log 3/6 = 1$

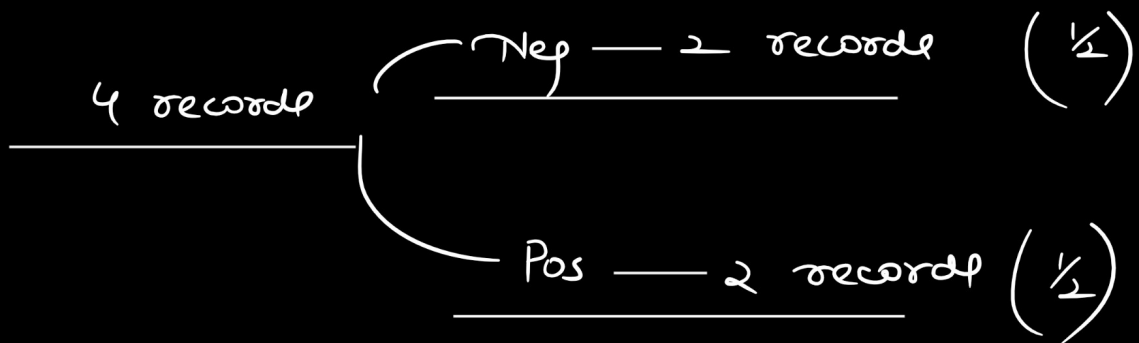
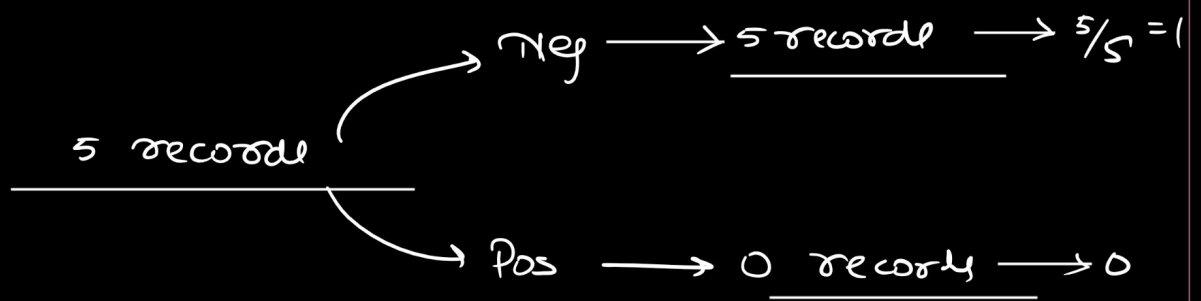
$$\text{Entropy}(\text{Gender}) = \frac{8}{14} * 0.954 + \frac{6}{14} * 1$$

$$\Rightarrow \underline{\underline{0.974}}$$

Note: which features give the most information

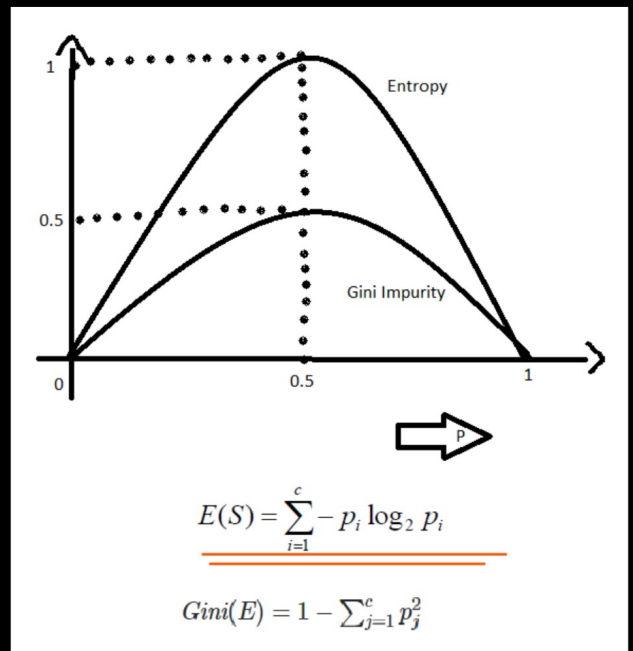
↳ Information  
Gain

	<u>[0, 1]</u>	
	<u>Pure</u>	<u>Impure</u>
<u>Entropy</u> →	$-0 \log 0 - 1 \log 1$	$-\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2}$
	$\Rightarrow 0$	$\Rightarrow \underline{\underline{1}}$





$Gini \rightarrow 0.5$  Highest Impurity  
 $Entropy \rightarrow 1$



Tasks

- use case
- ① When should we go for "Entropy" as the criterion?
  - ② Decision Tree Regression