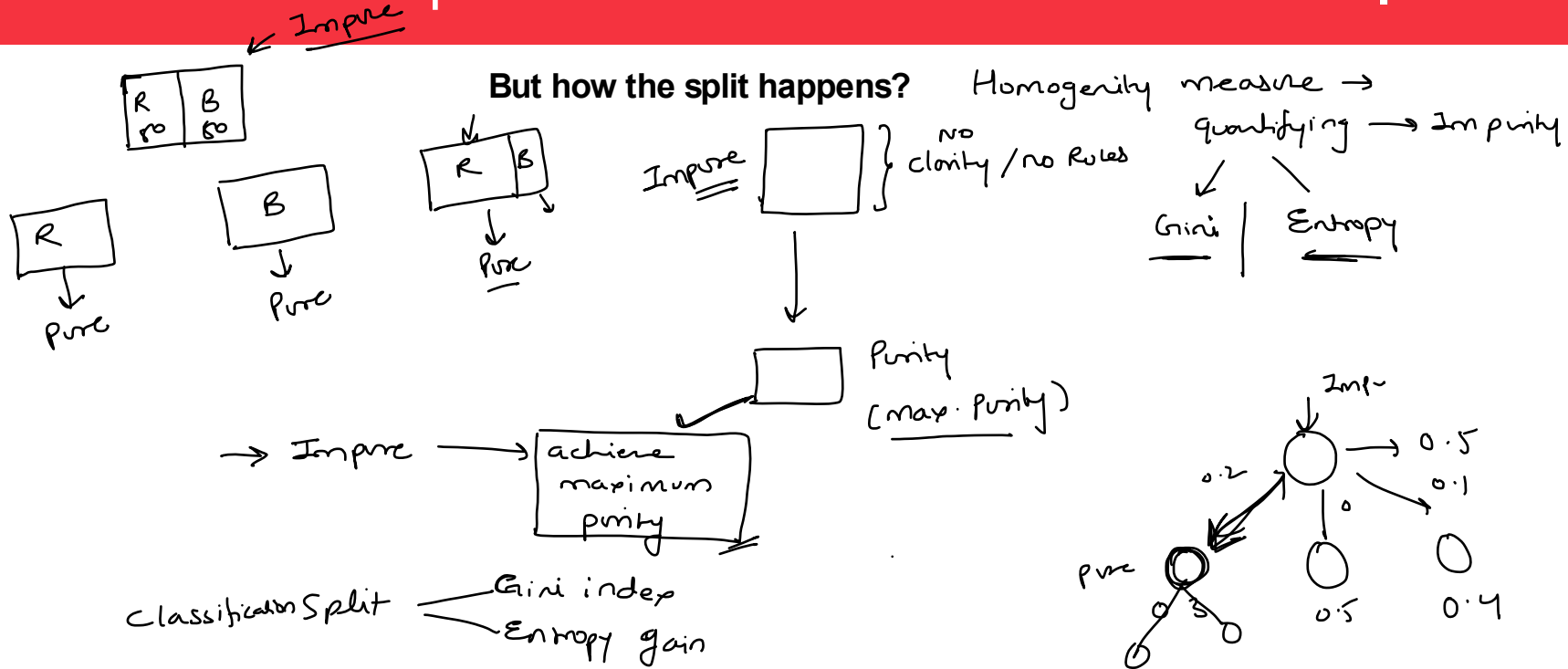


# Tree Models: Concepts and Doubts

# Tree Models: Concepts



# Tree Models: Concepts

Yes	9
No	5

$$\sum p_i \leq 1$$

upGrad

$$\Delta G_B - G_A = \text{maximum}$$

	Yes	No	
Rainy	2	3	5
overcast	4	0	4
Sunny	3	2	5

Gini Index

Before split

$$\sum p_i(1-p_i) = 1 - \sum p_i^2$$

$$\begin{aligned} G(\text{play golf}) &= 1 - \left[ \left(\frac{9}{14}\right)^2 + \left(\frac{5}{14}\right)^2 \right] \\ &= 1 - [0.41 + 0.12] \\ &= 0.47 \end{aligned}$$

$$\begin{aligned} G(\text{outlook}) &= \\ &= \frac{5}{14} \left[ 1 - \left\{ \left(\frac{2}{5}\right)^2 + \left(\frac{3}{5}\right)^2 \right\} \right] + \frac{4}{14} \left[ 1 - \left\{ \left(\frac{4}{4}\right)^2 \right\} \right] + \frac{5}{14} \left[ 1 - \left\{ \left(\frac{3}{5}\right)^2 + \left(\frac{2}{5}\right)^2 \right\} \right] \end{aligned}$$

$$\Rightarrow \frac{5}{14} [1 - 0.52] + \frac{4}{14} \times 0 + \frac{5}{14} \times [1 - 0.52]$$

$$\Rightarrow 0.171 + 0 + 0.171 = 0.342$$

$$\Delta (0.47 - 0.34) = 0.13 \checkmark \text{ maximum}$$

Outlook	Temp.	Humidity	Windy	Play Golf
Rainy	Hot	High	False	No —
Rainy	Hot	High	True	No —
Overcast	Hot	High	False	Yes —
Sunny	Mild	High	False	Yes —
Sunny	Cool	Normal	False	Yes —
Sunny	Cool	Normal	True	No —
Overcast	Cool	Normal	True	Yes —
Rainy	Mild	High	False	No —
Rainy	Cool	Normal	False	Yes —
Sunny	Mild	Normal	False	Yes —
Rainy	Mild	Normal	True	Yes —
Overcast	Mild	High	True	Yes —
Overcast	Hot	Normal	False	Yes —
Sunny	Mild	High	True	No —

# Poll Question

	c <sub>0</sub>	c <sub>1</sub>	
family	1	3	= 4
sports	8	0	= 8
Luxury	1	7	= 8

	c <sub>0</sub>	c <sub>1</sub>	
Male	6	4	= 10
female	4	6	= 10

In the following questions we will try to calculate Gini Index

**Question-1:** Compute the Gini Index of the overall training example.

**Question-2:** Calculate the Gini Index of the gender attribute

**Question-3:** Calculate the Gini Index of the car\_type attribute

$$G(\text{Class}) = 1 - \left[ \left( \frac{10}{20} \right)^2 + \left( \frac{10}{20} \right)^2 \right] = 0.5$$

$$G(\text{Gender}) = \frac{10}{20} \left[ 1 - \left\{ \left( \frac{6}{10} \right)^2 + \left( \frac{4}{10} \right)^2 \right\} \right] + \frac{10}{20} \left[ 1 - \left\{ \left( \frac{6}{10} \right)^2 + \left( \frac{4}{10} \right)^2 \right\} \right]$$

$$= 0.48$$

	c <sub>0</sub>	c <sub>1</sub>	
Small	3	2	5
medium	3	4	7
Large	2	2	4
XL	2	2	4

Customer ID	Gender	Car Type	Shirt Size	Class
1	M	Family	Small	C0
2	M	Sports	Medium	C0
3	M	Sports	Medium	C0
4	M	Sports	Large	C0
5	M	Sports	Extra Large	C0
6	M	Sports	Extra Large	C0
7	F	Sports	Small	C0
8	F	Sports	Small	C0
9	F	Sports	Medium	C0
10	F	Luxury	Large	C0
11	M	Family	Large	C1
12	M	Family	Extra Large	C1
13	M	Family	Medium	C1
14	M	Luxury	Extra Large	C1
15	F	Luxury	Small	C1
16	F	Luxury	Small	C1
17	F	Luxury	Medium	C1
18	F	Luxury	Medium	C1
19	F	Luxury	Medium	C1
20	F	Luxury	Large	C1

$$n(\text{cor-type}) = \frac{4}{20} \left[ 1 - \left\{ \left( \frac{1}{4} \right)^2 + \left( \frac{3}{4} \right)^2 \right\} \right] + \frac{8}{20} \left[ 1 - \left( \frac{8}{8} \right)^2 \right] + \frac{8}{20} \left[ 1 - \left\{ \left( \frac{1}{8} \right)^2 + \left( \frac{7}{8} \right)^2 \right\} \right]$$

$$= 0.1625$$

$$n(\text{shirt size}) =$$

$$\frac{5}{20} \left[ 1 - \left\{ \left( \frac{3}{5} \right)^2 + \left( \frac{2}{5} \right)^2 \right\} \right] + \frac{7}{20} \left[ 1 - \left\{ \left( \frac{3}{7} \right)^2 + \left( \frac{4}{7} \right)^2 \right\} \right] + \frac{4}{20} \left[ 1 - \left\{ \left( \frac{2}{4} \right)^2 + \left( \frac{2}{4} \right)^2 \right\} \right]$$

$$+ \frac{4}{20} \left[ 1 - \left\{ \left( \frac{2}{4} \right)^2 + \left( \frac{2}{4} \right)^2 \right\} \right]$$

$$= 0.49$$

$$\Delta(0.5 - 0.48) = 0.02$$

$$\Delta(0.5 - 0.1625) = 0.33$$

$$\Delta(0.5 - 0.49) = 0.01$$


} → ✓

$$Entropy = (-P_i \log_2 P_i)$$

**Question-4:** Calculate the Gini Index of the shirt\_size attribute

**Question-5:** Among all the attributes, which attribute is better to use for the first split in the formation of the decision tree?

- 0.52
- 0.55
- 0.54



Customer ID	Gender	Car Type	Shirt Size	Class
1	M	Family	Small	C0
2	M	Sports	Medium	C0
3	M	Sports	Medium	C0
4	M	Sports	Large	C0
5	M	Sports	Extra Large	C0
6	M	Sports	Extra Large	C0
7	F	Sports	Small	C0
8	F	Sports	Small	C0
9	F	Sports	Medium	C0
10	F	Luxury	Large	C0
11	M	Family	Large	C1
12	M	Family	Extra Large	C1
13	M	Family	Medium	C1
14	M	Luxury	Extra Large	C1
15	F	Luxury	Small	C1
16	F	Luxury	Small	C1
17	F	Luxury	Medium	C1
18	F	Luxury	Medium	C1
19	F	Luxury	Medium	C1
20	F	Luxury	Large	C1

# Tree Models: Concepts

- Use multiple trees
- Bagging
- Add bias at each tree and at each node.



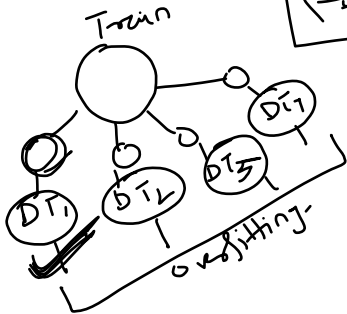
Decision Tree

← overfit your training data

- High variance model
- low bias

Bootstrapping

~~Bagging~~ - Level-1 Adding bias



- Level-2 Bias
- 

## Random Forest

Bagging → Bootstrapped → Creating samples with repetition  
 → Aggregation → Combining.

100  
 9  
~~Total~~

1	2	3
4	5	6
7	8	9

1	1	2	6	8	9	9	5	5	B.S1 ✓
4	5	5	4	3	2	1	1	9	B.S2 ✓
7	6	5	4	7	2	1	3	1	B.S3 ✓

BS50

BS1

BS2

BS3

BS7

DT50

DT1

DT2

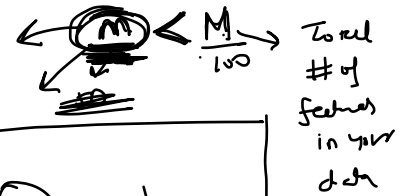
DT3

DT4



Those trees for which, it was not the part of the B-Sample

→ At each node a random sample of feature are taken for creating the split



Agg.  $\left[ \frac{25}{5} \leq \frac{20}{5} \rightarrow \text{true} \right]$

1 → true

# Tree Models: Concepts

oob error =  $\frac{\# \text{ of incorrect class}}{\text{total no. of data}}$  **upGrad**

$\frac{4}{10} = 0.4$

		Training data		Random Forest (oob error)						
Prediction	Actual		BS <sub>1</sub> DT <sub>1</sub>	BS <sub>2</sub> DT <sub>2</sub>	BS <sub>3</sub> DT <sub>3</sub>	BS <sub>4</sub> DT <sub>4</sub>	BS <sub>5</sub> DT <sub>5</sub>	BS <sub>6</sub> DT <sub>6</sub>	#	
true	true	1	X	true	X	true	-ve	X		
-ve	-ve	2	-ve	X	-ve	X	X	true		
true	-ve	3	X	X	true	X	true	true		
-ve	true	4	-ve	X	-ve	-ve	X	X		
true	true	5	X	true	X	true	X	true		
-ve	-ve	6	-ve	X	true	X	-ve	X		
-ve	true	7	true	X	X	X	-ve	-ve		
-ve	-ve	8	X	-ve	X	-ve	-ve	X		
-ve	true	9	-ve	X	true	X	-ve	X		
-ve	true	10	X	true	X	X	-ve	-ve		





# Thank You!