

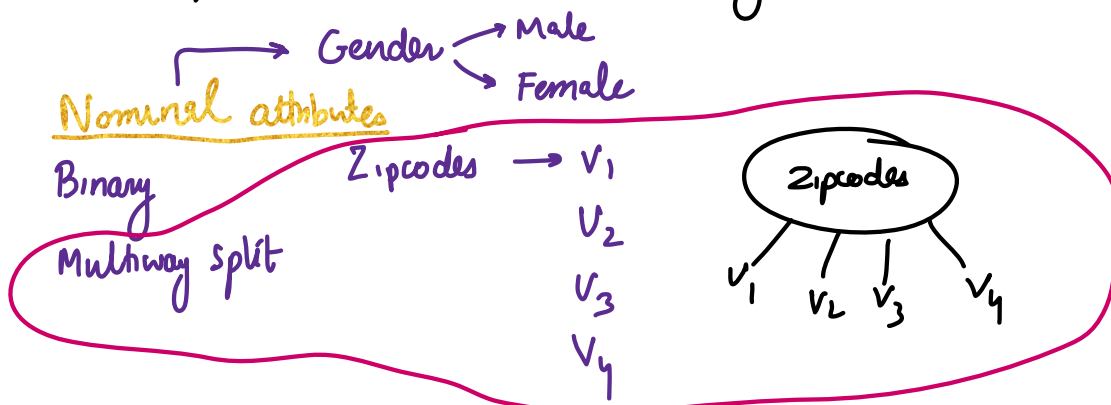
# Lecture 20

Tuesday, October 25, 2022 12:07 PM

Till hunt's algo in Midsem  
current slides till pg. 45

## Split Conditions & outcomes

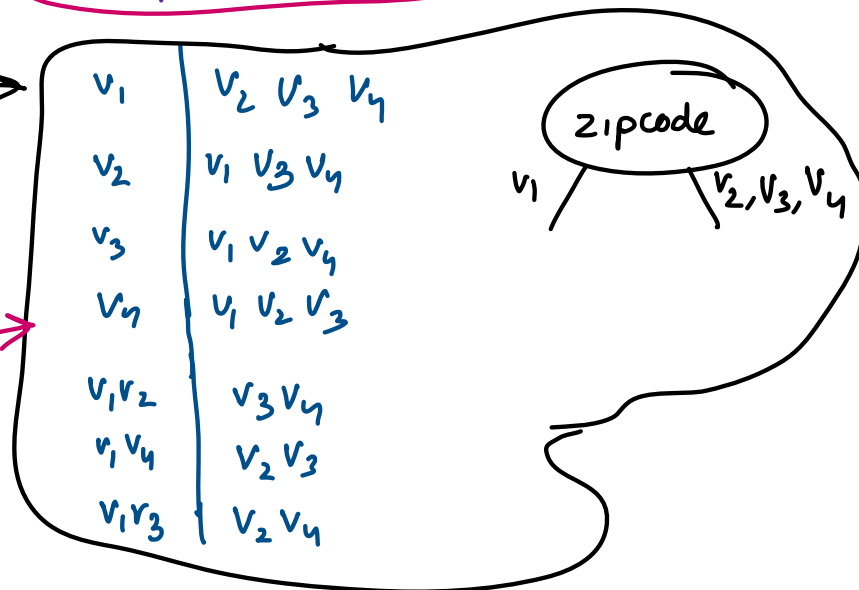
Ordinal, Nominal, Continuous, Binary



Binary split

$2^{k-1}$   
 $k = \text{unique numbers}$   
 $k = 4$

$$2^{4-1} - 1 = 7$$



Ordinal attributes

Binary

Multway split

Service → very good, good, bad, very bad

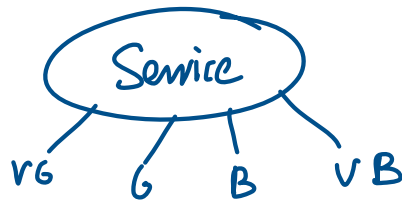
Order property must be maintained while grouping values for binary split.

VG	G, B, VB	} possible splits, maintain order property ✓✓
VG, G	B, VB	
VG, G, B	VB	



VG, B      G, VB → doesn't maintain order XX

for multiway split, its the same



## Binary Attributes

binary splits { no multiway, cuz all we have is 2 }

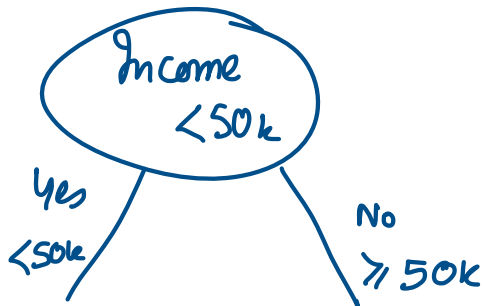


## Continuous Attributes

Binary split

Multway split

Split value  $\rightarrow (v)$



10k  $\rightarrow$  100k

$< 10k, 10k-25k, 25k-50k, 50k-80k, > 80k$

② split it into the ranges & select the one with most efficiency.

① we get  $v$  value by checking it with all the values in the column & find which one splits it the best by using certain measures

### Measures of Node Impurity

class distribution  $\begin{cases} \text{before splitting} \\ \text{after splitting} \end{cases}$

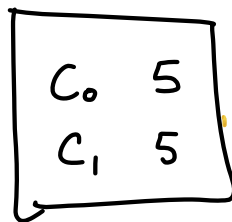
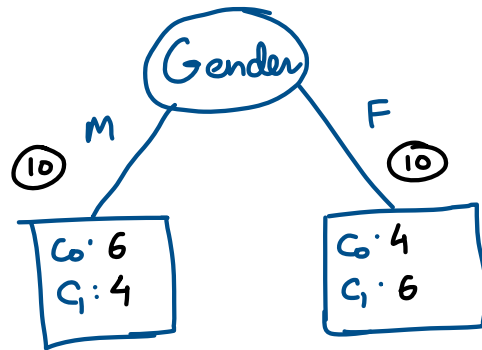
class  $\begin{cases} c_0 \\ c_1 \end{cases}$



Measures of Node Impurity				
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

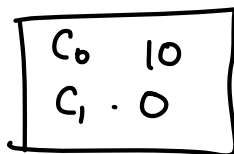
Class distribution before Splitting:  
10 records of class 0  
10 records of class 1

① Gender



⊛ This is uniform class distribution

MOST IMPURE SPLIT



⊛ Most skewed distribution

PUREST SPLIT

{ IDEAL }

is that a word?

What we need is the measure of the skewedness of the split

## Impurity Measure

$t \rightarrow$  node

$i \rightarrow$  class

$P(i/t) \rightarrow$  fraction of records belonging to  $i$  @ node  $t$

$C \rightarrow$  no. of classes

⊛ Entropy =  $-\sum_{i=0}^{C-1} P_i \log_2 P_i$

1  $\leq C \leq 2$

} calculated in example.

$$(*) \text{ Gini} = 1 - \sum_{i=0}^{C-1} p_i^2$$

$$(*) \text{ Classification Error} = 1 - \max_i \{p_i\}$$

in example,  
don't panic



$$\log_2(5/6) = \frac{\log(5/6)}{\log 2}$$

## Measures of Node Impurity



Node $N_1$	Count
Class=0	0
Class=1	6

$$\text{Gini} = 1 - (0/6)^2 - (6/6)^2 = 0$$

$$\text{Entropy} = -(0/6) \log_2(0/6) - (6/6) \log_2(6/6) = 0$$

$$\text{Error} = 1 - \max[0/6, 6/6] = 0$$

Node $N_2$	Count
Class=0	1
Class=1	5

$$\text{Gini} = 1 - (1/6)^2 - (5/6)^2 = 0.278$$

$$\text{Entropy} = -(1/6) \log_2(1/6) - (5/6) \log_2(5/6) = 0.650$$

$$\text{Error} = 1 - \max[1/6, 5/6] = 0.167$$

Node $N_3$	Count
Class=0	3
Class=1	3

$$\text{Gini} = 1 - (3/6)^2 - (3/6)^2 = 0.5$$

$$\text{Entropy} = -(3/6) \log_2(3/6) - (3/6) \log_2(3/6) = 1$$

$$\text{Error} = 1 - \max[3/6, 3/6] = 0.5$$

Node 1

$$C_0 \rightarrow 0$$

$$C_1 \rightarrow 6$$

$$p_0 \rightarrow 0/6 = 0$$

$$p_1 \rightarrow 6/6 = 1$$

$$\text{Entropy} = - [0 \log 0 + 1 \log 1] = 0$$

$$\text{Gini} = 1 - [0^2 + 1^2] = 0$$

$$\underline{Gini} = 1 - [0^2 + 1^2] = 0$$

$$\underline{\text{Classification Error}} = 1 - \max\{0, 1\} = 1 - 1 = 0$$

Node 2

$$C_0 \rightarrow 1 \quad P_0 = 1/6$$

$$C_1 \rightarrow 5 \quad P_1 = 5/6$$

$$\underline{\text{Entropy}} = - \left[ \frac{1}{6} \log 6 + \frac{5}{6} \log \frac{5}{6} \right] = 0.65$$

$$\underline{Gini} = 1 - \left[ \left(\frac{1}{6}\right)^2 + \left(\frac{5}{6}\right)^2 \right] = 1 - \frac{26}{36} = 0.278$$

$$\underline{\text{Classification Error}} = 1 - \max\{1/6, 5/6\} = 1/6$$

Homework { Verify with slides }

Node 3