

Data MiningProblem 1 :-

1.1

$$\begin{aligned}
 E(P) &= \left(\frac{-5}{12} \log_2 \left(\frac{5}{12} \right) - \left(\frac{7}{12} \right) \log_2 \left(\frac{7}{12} \right) \right) \\
 &= +.525 - .453 + .525 \\
 &= .978
 \end{aligned}$$

Survived	5
Not Survived	7

$$\begin{aligned}
 E(\text{Survived}, \text{Gen}) &= \frac{6}{12} \left(\frac{-3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \left(\frac{3}{6} \right) \right) \\
 &\quad + \left(\frac{6}{12} \left(\frac{4}{6} \right) + \frac{-2}{6} \log_2 \frac{2}{6} - \frac{4}{6} \log_2 \left(\frac{4}{6} \right) \right) \\
 &= \frac{6}{12} (-.2 + .918) \\
 &= \frac{6}{12} (.718) = 0.359
 \end{aligned}$$

Gender

```

graph TD
    Gender --> Male
    Gender --> Female
    Male --> Yes1[Yes: 2]
    Male --> No1[No: 4]
    Female --> Yes2[Yes: 3]
    Female --> No2[No: 3]
  
```

$$E(\text{Survived}, \text{Gen})$$

$$\begin{aligned}
 \text{Information gain} &= E(\text{Survived}) - E(\text{Survived}, \text{Gen}) \\
 &= .024
 \end{aligned}$$

E(Survived, age)

Cheatsheet	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Age	8	12	20	24	26	32	34	43	45	47	55	61			
$\leq 6 <$	$\leq 10 <$	$\leq 16 <$	$\leq 22 <$	$\leq 25 <$	$\leq 29 <$	$\leq 32 <$	$\leq 34 <$	$\leq 41 <$	$\leq 45 <$	$\leq 51 <$	$\leq 55 <$	$\leq 70 <$			
Yes	0	5	1	4	2	3	3	2	3	2	3	2	4	1	5
No	0	7	0	7	1	6	1	6	2	5	3	4	3	5	0
Enslav.	.978	.86	.73	.91	.92	.908	.95								
Avg	0	.11	.24	.062	.16	.072									

Calculations :-

$$\leq 8 < \left(\frac{0}{12} + \left(-\frac{0}{12} \log_2 \left(\frac{0}{12} \right) - \frac{12}{12} \log_2 \left(\frac{12}{12} \right) \right) + \frac{12}{12} \left(-\frac{5}{12} \log_2 \left(\frac{5}{12} \right) - \frac{7}{12} \log_2 \left(\frac{7}{12} \right) \right) \right) = 0$$

$$\leq 16 < \left(\frac{1}{12} \left(-\frac{1}{1} \log_2 \left(\frac{1}{1} \right) - \frac{1}{1} \log_2 \left(\frac{1}{1} \right) \right) + \frac{12}{12} \left(-\frac{4}{11} \log_2 \left(\frac{4}{11} \right) - \frac{7}{12} \log_2 \left(\frac{7}{12} \right) \right) \right) = .86$$

$$\leq 20 < \left(\frac{2}{12} \left(-\frac{2}{2} \log_2 \left(\frac{2}{2} \right) - \frac{10}{12} \log_2 \left(\frac{10}{12} \right) \right) + \frac{10}{12} \left(-\frac{3}{10} \log_2 \left(\frac{3}{10} \right) - \frac{7}{10} \log_2 \left(\frac{7}{10} \right) \right) \right) = .73$$

$$\leq 24 < \left(\frac{3}{12} \left(-\frac{2}{3} \log_2 \left(\frac{2}{3} \right) - \frac{3}{3} \log_2 \left(\frac{1}{3} \right) \right) + \frac{9}{12} \left(-\frac{3}{9} \log_2 \left(\frac{3}{9} \right) - \frac{6}{9} \log_2 \left(\frac{6}{9} \right) \right) \right) = .91$$

$$\leq 28 < \left(\frac{4}{12} \left(-\frac{3}{4} \log_2 \left(\frac{3}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right) + \frac{8}{12} \left(-\frac{2}{8} \log_2 \left(\frac{2}{8} \right) - \frac{6}{8} \log_2 \left(\frac{6}{8} \right) \right) \right) = .91$$

E(Survived, age)

Cheat sheet	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Age	8	12	20	24	26	32	34	43	45	47	55	65			
$\leq 6 <$	$\leq 10 <$	$\leq 16 <$	$\leq 20 <$	$\leq 25 <$	$\leq 29 <$	$\leq 32 <$	$\leq 34 <$	$\leq 40 <$	$\leq 44 <$	$\leq 48 <$	$\leq 51 <$	$\leq 56 <$	$\leq 60 <$	$\leq 65 <$	
Yes	0	5	142	323	332	323	323	323	323	323	44	4150	50		
No	0	7	070216162534435252616170												
Entropy	.178	.86	.73	.91	.92	.908	.95								
Gain	0	.11	.24	.062	.16	.072									

Calculation :-

$$\leq 6 < = \left(\frac{1}{12} + \left(-\frac{5}{12} \log_2 \left(\frac{5}{12} \right) - \frac{7}{12} \log_2 \left(\frac{7}{12} \right) \right) \right)$$

$$= 0$$

$$\leq 8 < = \left(\frac{1}{12} \left(-\frac{1}{1} \log_2 \left(\frac{1}{1} \right) - \frac{5}{12} \log_2 \left(\frac{5}{12} \right) \right) + \frac{11}{12} \left(-\frac{4}{11} \log_2 \left(\frac{4}{11} \right) - \frac{7}{12} \log_2 \left(\frac{7}{12} \right) \right) \right)$$

$$= .86$$

$$\leq 16 < = \left(\frac{2}{12} \left(-\frac{2}{2} \log_2 \left(\frac{2}{2} \right) - \frac{10}{12} \log_2 \left(\frac{10}{12} \right) \right) + \frac{10}{12} \left(-\frac{3}{10} \log_2 \left(\frac{3}{10} \right) - \frac{7}{10} \log_2 \left(\frac{7}{10} \right) \right) \right)$$

$$= .73$$

$$\leq 20 < = \left(\frac{3}{12} \left(-\frac{2}{3} \log_2 \left(\frac{2}{3} \right) - \frac{3}{3} \log_2 \left(\frac{1}{3} \right) \right) + \frac{9}{12} \left(-\frac{3}{9} \log_2 \left(\frac{3}{9} \right) - \frac{6}{9} \log_2 \left(\frac{6}{9} \right) \right) \right)$$

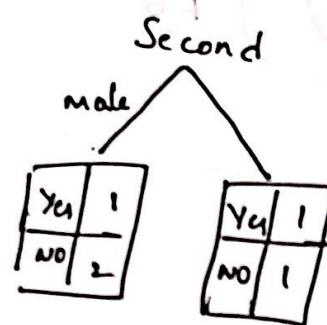
$$\leq 24 < = \left(\frac{4}{12} \left(-\frac{3}{4} \log_2 \left(\frac{3}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right) + \frac{8}{12} \left(-\frac{2}{8} \log_2 \left(\frac{2}{8} \right) - \frac{6}{8} \log_2 \left(\frac{6}{8} \right) \right) \right)$$

$$= .91$$

Performing Split again

on second level of Jtree

E(P)



$$\text{Reftch} = \frac{3}{5} \left(-\frac{1}{3} \log \left(\frac{1}{3} \right) \right) + \frac{2}{5} \left(-\frac{1}{2} \log \left(\frac{1}{2} \right) \right) + \frac{2}{5} \left(-\frac{1}{2} \log \left(\frac{1}{2} \right) \right) + \frac{1}{5} \left(-\frac{1}{2} \log \left(\frac{1}{2} \right) \right)$$
$$= .970$$

0.8	0.4	0.68	0.8	-0.68
43	25	71	3.813	0.7
222	2.882	222	2.882	222
0.8	1.4	3.8	1.4	0.8
0.1	0.1	0.8	0.1	0.1

$$E(\text{Survived}, \text{Gender}) = \frac{2}{4} \left(-\frac{1}{2} \log \left(\frac{1}{2} \times \frac{1}{2} \right) \right) + 0$$

$$= \frac{2}{4} \left(-\frac{1}{2} \log \left(\frac{1}{4} \right) \right)$$

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

$$= \underline{-0.5}$$

Information gain = $\cdot 811 - \cdot 5$
 $= \underline{\cdot 311}$

Age

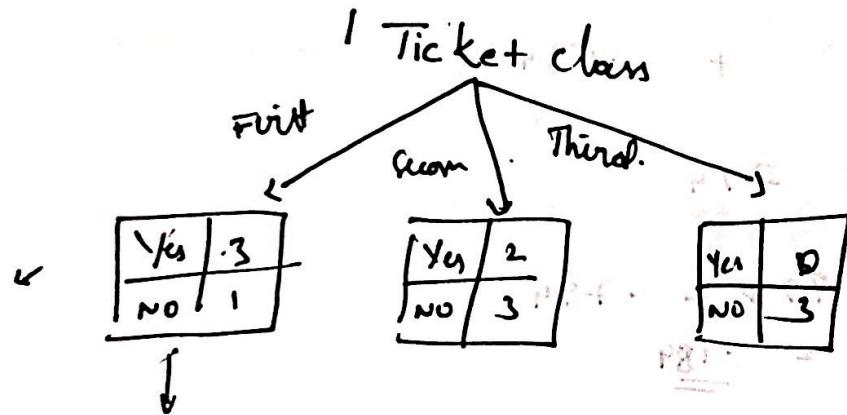
cheat sheet		Yes	No Yes	No	Yes
≤ 80	≤ 24	12	20	24	
< 80 ≤	≤ 60 <	≤ 16 <	≤ 22 <	≤ 26 <	
Yes	0	5	1	4	2
No	0	2	0	2	0
Elaps					
Gain	0	-0.32	-0.3	-0.12	0

$$\text{Information gain} = 0.978 - 0.67 \\ = \underline{\underline{-0.304}}$$

~~We to consider option~~
~~Information gain~~

III). as we are getting higher

Now consider Ticket class is the Parent as it has given highest information among the dependent variables.



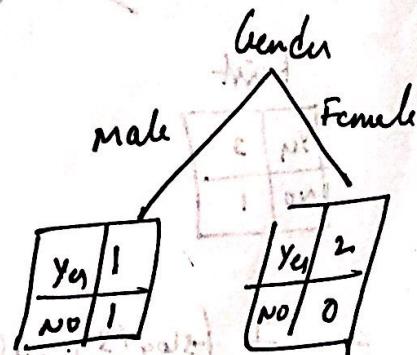
$$\text{Entropy (first)} = \left(-\frac{3}{4} \log \left(\frac{3}{4}\right) + -\frac{1}{4} \log \left(\frac{1}{4}\right) \right)$$

Entropy (Second)

Take gender first part

~~E (curve, ends)~~

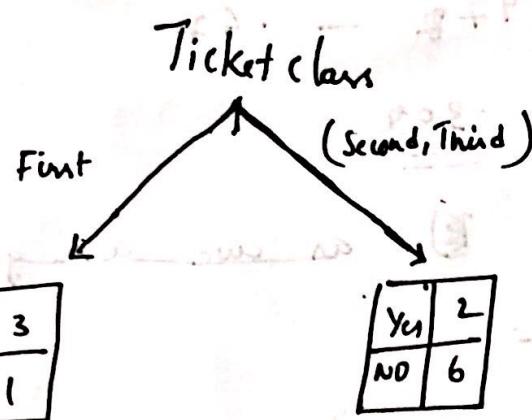
$$E(\text{Survived}, \text{Gender}) = \frac{2}{4} \left(-\frac{1}{2} \log\left(\frac{1}{2}\right) - \frac{1}{2} \log\left(\frac{1}{2}\right) + \right. \\ \left. \frac{2}{4} \left(-\frac{2}{2} \log\left(\frac{2}{2}\right) - 0 \right) \right)$$



$E(\text{survived}, \text{Ticket class})$

option

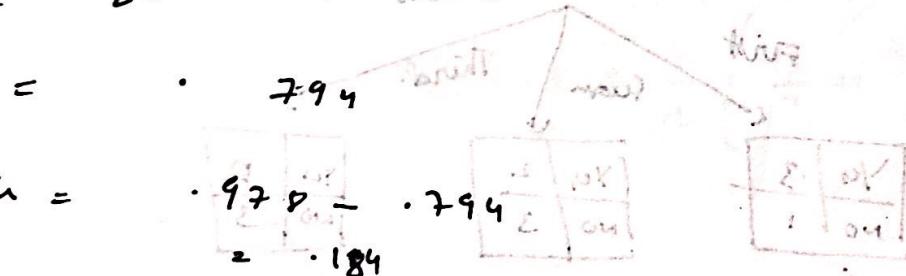
a)



$$E = + \frac{4}{12} \left(-\frac{3}{4} \log_2 \left(\frac{3}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right) + \frac{8}{12} \left(-\frac{2}{8} \log_2 \left(\frac{2}{8} \right) - \frac{6}{8} \log_2 \left(\frac{6}{8} \right) \right)$$

$$= -2.688 + \frac{1}{3} (+0.5) + (-0.3025)$$

$$= -2.6 + 0.534 = -2.066$$



Information gain =

$$\cdot 925 - \cdot 794$$

$$= \underline{\underline{.184}}$$

option

b)



First

	3
Yes	3
No	1

Second

	2
Yes	2
No	3

Third

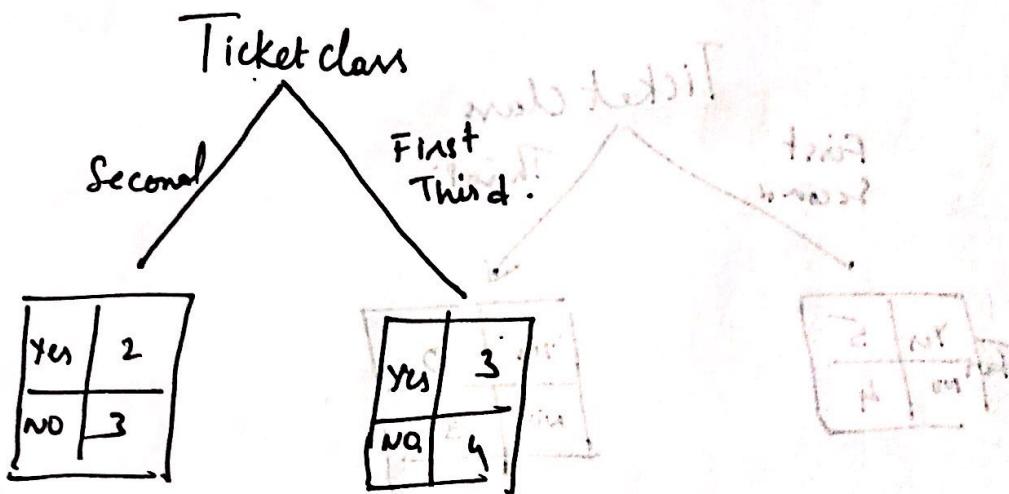
	0
Yes	0
No	3

$$E = \frac{4}{12} \left(-\frac{3}{4} \log_2 \left(\frac{3}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right) + \frac{5}{12} \left(-\frac{2}{5} \log_2 \left(\frac{2}{5} \right) + -\frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right)$$

$$+ (0 \cdot -\frac{3}{3} \log_2 \frac{3}{3})$$

$$= \frac{4}{12} (-0.7267) + \frac{5}{12}$$

$$E(\text{Survived}, \text{Ticket Class})$$



$$E(\text{Survived}, \text{Ticket Class}) =$$

$$+ \frac{5}{12} \left(-\frac{2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right) + \frac{7}{12} \left(-\frac{3}{7} \log_2 \left(\frac{3}{7} \right) - \frac{4}{7} \log_2 \left(\frac{4}{7} \right) \right)$$

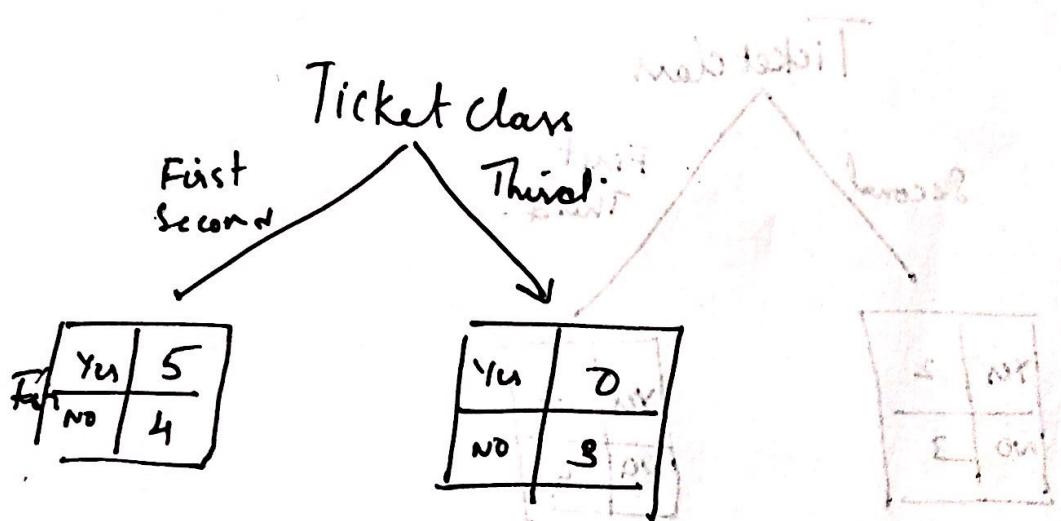
• 973

$$\text{Information gain} = 0.978 - 973$$

$$= \underline{\underline{0.005}}$$

E (Survived, Ticket Class)

(and) ticket class



E (Survived, TicketClass)

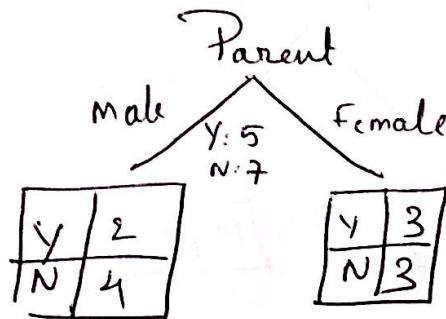
$$\begin{aligned} &= \frac{9}{12} \left(-\frac{5}{9} \log_2 \left(\frac{5}{9}\right) - \frac{4}{9} \log_2 \left(\frac{4}{9}\right) + \right. \\ &\quad \left. \frac{3}{12} \left(0 - \frac{3}{3} \log_2 \left(\frac{3}{3}\right) \right) \right) \\ &= \frac{9}{12} (-.4711 + .52) \\ &\quad (1.04) \times \frac{9}{12} \\ &= -.78 \end{aligned}$$

$$\text{Information gain} = .98 - .78$$

$$= \underline{\underline{.20}}$$

1.2

Splitting Using Gini Index)



$$G_{\text{INI}}(T) = 1 - \sum_j [P(j/T)]^2$$

$$G(\text{male}) = 1 - \left(\frac{2}{6}\right)^2 - \left(\frac{4}{6}\right)^2 = .44$$

$$G(\text{female}) = 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2 = \cancel{.5} \quad .5$$

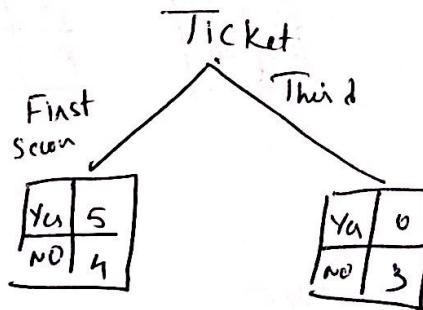
$$G_{\text{INI}}(\text{Parent}) = \frac{6}{12} (0.44) + \frac{6}{12} (0.5) = .47$$

$$1 - \left(\frac{5}{12}\right)^2 - \left(\frac{7}{12}\right)^2 = .486$$

$$\left. \begin{aligned} & - \frac{6}{12} (0.44) + \\ & - \frac{6}{12} (0.5) \end{aligned} \right\} = .47$$

Ticket (box)

option ①



$$G(\text{First, Second}) = 1 - \left(\frac{5}{9}\right)^2 - \left(\frac{4}{9}\right)^2$$

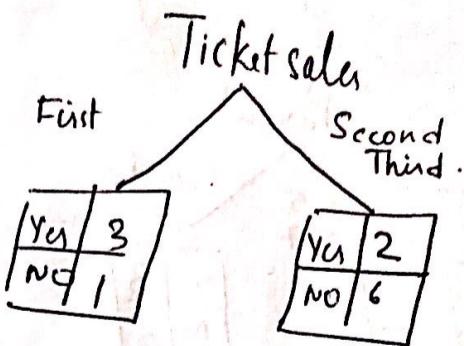
$$G(\text{Third}) = 1 - \left(\frac{0}{3}\right)^2 - \left(\frac{3}{3}\right)^2 = 0$$

$$\left. \begin{aligned} & \frac{9}{12} (0.49) + \frac{3}{12} (0) \\ & = .37 \end{aligned} \right\}$$

option ②

$$G(F_{\text{int}}) = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{7}\right)^2$$

= 0.375

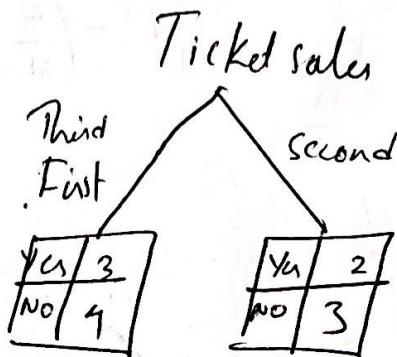


$$6(\text{second, third}) = 1 - \left(\frac{2}{8}\right)^2 - \left(\frac{6}{8}\right)^2 \\ = .375$$

$$G(\text{First}) \cdot \frac{6}{12} + G(\text{Second}) \cdot \frac{4}{12} + G(\text{Third}) \cdot \frac{8}{12} = .375$$

option(3)

$$G(\text{First}, \text{Third}) \\ = 1 - \left(\frac{3}{7}\right)^2 - \left(\frac{4}{7}\right)^2 \\ = 0.48$$



6 (second)

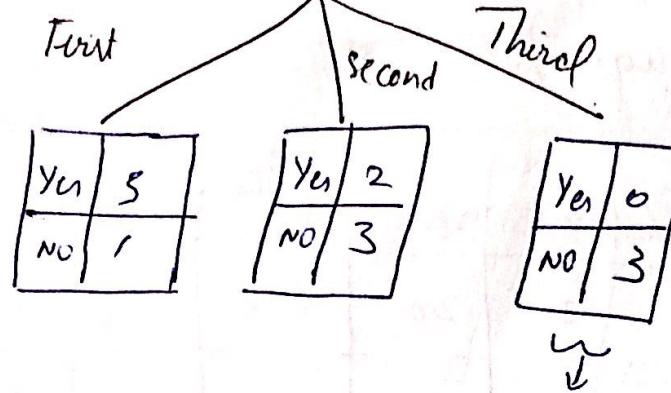
$$= 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 48$$

$$f(\text{First}, \text{Third}) = 10 -$$

$$+ 6(\text{second}) \quad \left(\frac{7}{12} \right) (.48) + \left(\frac{5}{12} \right) (.48)$$

= .48

Ticket sales



$G_{\min(\text{ticket class})} = \frac{4}{12} \left(1 - \left(\frac{3}{4} \right)^2 - \left(\frac{1}{4} \right)^2 \right) + \frac{5}{12} \left(1 - \left(\frac{2}{5} \right)^2 - \left(\frac{3}{5} \right)^2 \right)$

$$+ \frac{3}{12} \left(1 - \left(\frac{0}{3} \right)^2 - \left(\frac{3}{3} \right)^2 \right)$$

$$= \frac{4}{12} (.375) + \frac{5}{12} (.48) = \underline{\underline{.325}}$$

Cheat BS

→ $f(\text{Survived}, \text{Age})$

Cheat sheet	≤ 1	Yes	No	Yes	No	No	No	No	Yes	No	...
Age	8	12	20	24	26	32	34	43	45	47	...
≤ 6 <	≤ 10 <	≤ 16 <	≤ 22 <	≤ 25 <	≤ 28 <	≤ 33 <	≤ 38 <	≤ 44 <	≤ 46 <	≤ 51 <	...
Yes	0	5	11	4	23	2	3	3	2	3	1
No	0	7	0	7	1	6	1	6	2	5	1
Gini	0.486	0.42	0.35	0.44	0.37	0.43	0.44	0.48	0.47	0.481	0.45

Least gini.

→ Least gini from all the attributes belongs to Ticket class

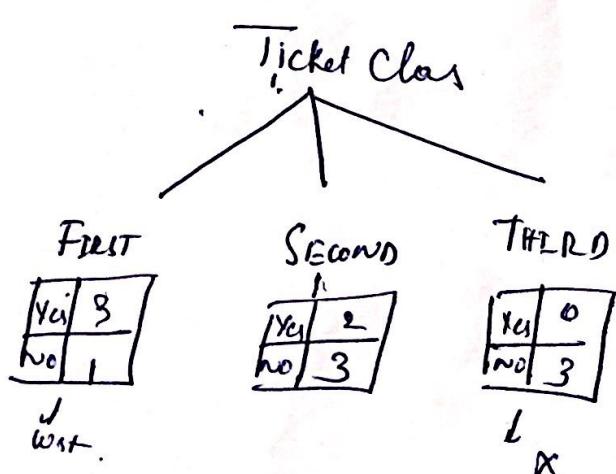
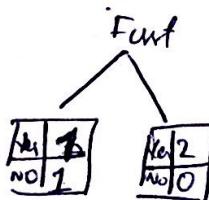
So we make Ticket class as the Parent node.

Next

Gender

$$\begin{aligned} \text{Gini}(\text{male}) &= 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 \\ &= 0.5 \end{aligned}$$

$$\text{Gini } (\text{Female}) = 1 - \left(1 - 0^2\right) = 0$$



$$\text{Gini } (\text{Female}) \quad \text{Gini } (\text{Gender}) = \frac{2}{4} (0.5) = \underline{\underline{0.25}}$$

$G(\text{Age})$

(Cheat sheet) Survived	Y	Y	N	Y		
Age	8	12	20	24		
Split	$\leq 4 <$	$\leq 10 <$	$\leq 16 <$	$\leq 22 <$	$\leq 25 <$	
Yes	0 3 1 2 2 1 2 1 3 0					
No	0 1 0 1 0 1 2 0 1 0					
Gini	.37	.33	.25	.33	.37	



Least Gini

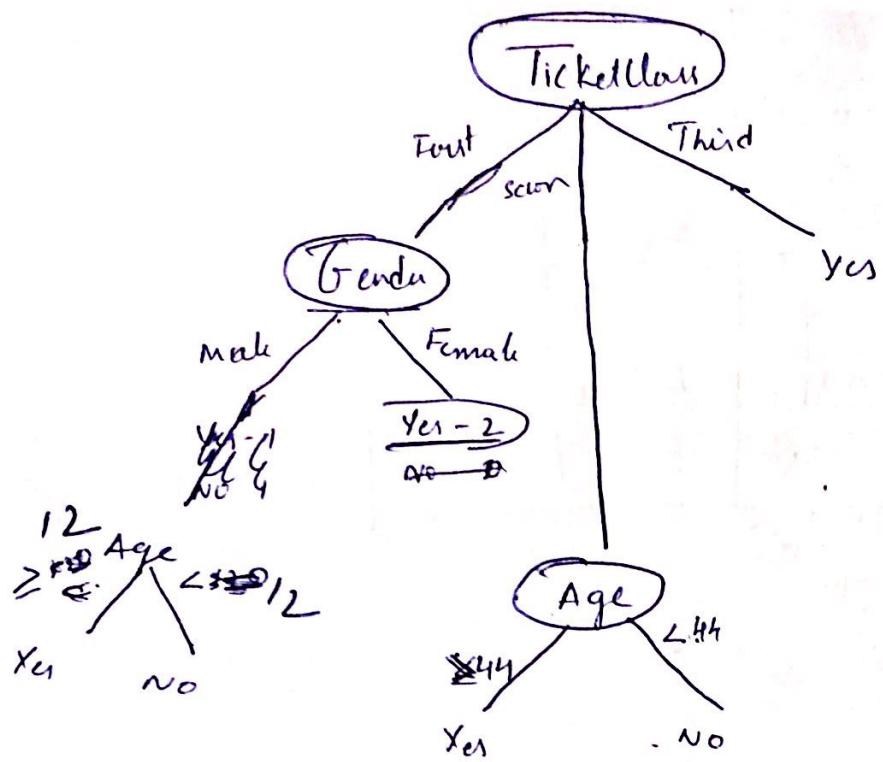
Since we got the gini index equal for both Age and gender
we can choose any of the following and perform split by
using age.

Cheat sheet	NO	NO	NO	YES	YES				
	32	34	43	45	55				
	$\leq 30 <$	$\leq 33 <$	$\leq 38 <$	$\leq 44 <$	$\leq 50 <$				
Y	2 0 2 0 2 0 2 1 1 2 0								
N	0 3 1 2 1 3 0 3 0 3 0								
Gini	.48	.4	.25	0	.3	0.48			



Least

Result



1.3) Error given to each node = .5

$$e(t) = \frac{\sum [e(t_i) + \omega(t_i)]}{\sum n(t_i)}.$$

$$\sum t_i = T$$

$$e'(t) = \frac{e(t) + \omega(t)}{n(t)}$$

$e(t)$ = classification error

n = total number of records

$\omega(t) = 0.25$

$$= \frac{1 + (0.5)(5)}{12} = \frac{5.5}{12}$$

$$e'(t) = \underline{.29}$$