

# Question 1

## CART Algorithm

Department	Junior	Senior	Total
Sales	30	80	110
Systems	23	8	31
Marketing	4	10	14
Secretary	6	4	$\frac{10}{165}$ 165

$$\text{Gini}(\text{dept} = \text{Sales}) = 1 - \left(\frac{30}{110}\right)^2 - \left(\frac{80}{110}\right)^2 = 0.397$$

$$\text{Gini}(\text{dept} = \text{Systems}) = 1 - \left(\frac{23}{31}\right)^2 - \left(\frac{8}{31}\right)^2 = 0.383$$

$$\text{Gini}(\text{dept} = \text{marketing}) = 1 - \left(\frac{4}{14}\right)^2 - \left(\frac{10}{14}\right)^2 = 0.408$$

$$\text{Gini}(\text{dept} = \text{secretary}) = 1 - \left(\frac{6}{10}\right)^2 - \left(\frac{4}{10}\right)^2 = 0.48$$

$$\begin{aligned} \text{Gini}(\text{dept}) &= \frac{110}{165} (0.397) + \frac{31}{165} (0.383) + \frac{14}{165} (0.408) + \frac{10}{165} (0.48) \\ &= 0.400 \end{aligned}$$

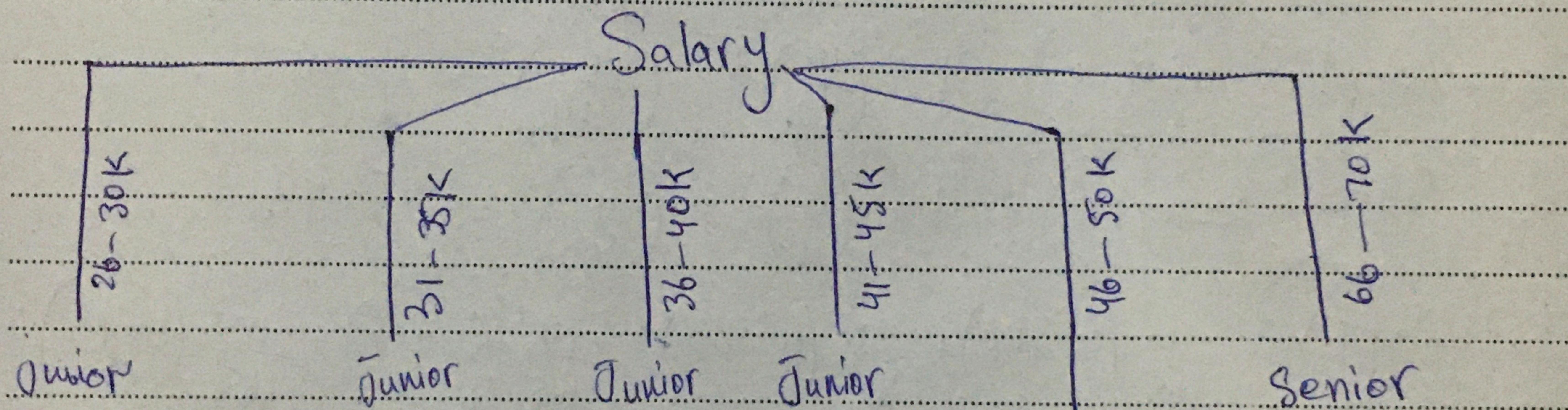
Age	Junior	Senior	Total	Gini
21 - 25	20	0	20	0
26 - 30	49	0	49	0
31 - 35	44	35	79	0.494
36 - 40	0	10	10	0
41 - 45	0	3	3	0
46 - 50	0	4	4	0

$$\text{Gini}(\text{Age}) = \frac{79}{165} (0.494) = 0.237$$

Salary	Junior	Senior	Total	Gini
26 - 30K	46	0	46	0
31 - 35K	40	0	40	0
36 - 40K	40	0	40	0
41 - 45K	4	0	4	0
46 - 50K	23	40	63	0.464
66K - 70K	0	8	8	0

$$\text{Gini}(\text{Salary}) = \frac{63}{165} (0.464) = 0.177$$

Splitting on minimum gini value, which is salary



dept	age	status	count
sales	31-35	senior	<del>30</del> 30
systems	21-25	junior	20
systems	26-30	junior	3
marketing	36-40	senior	10

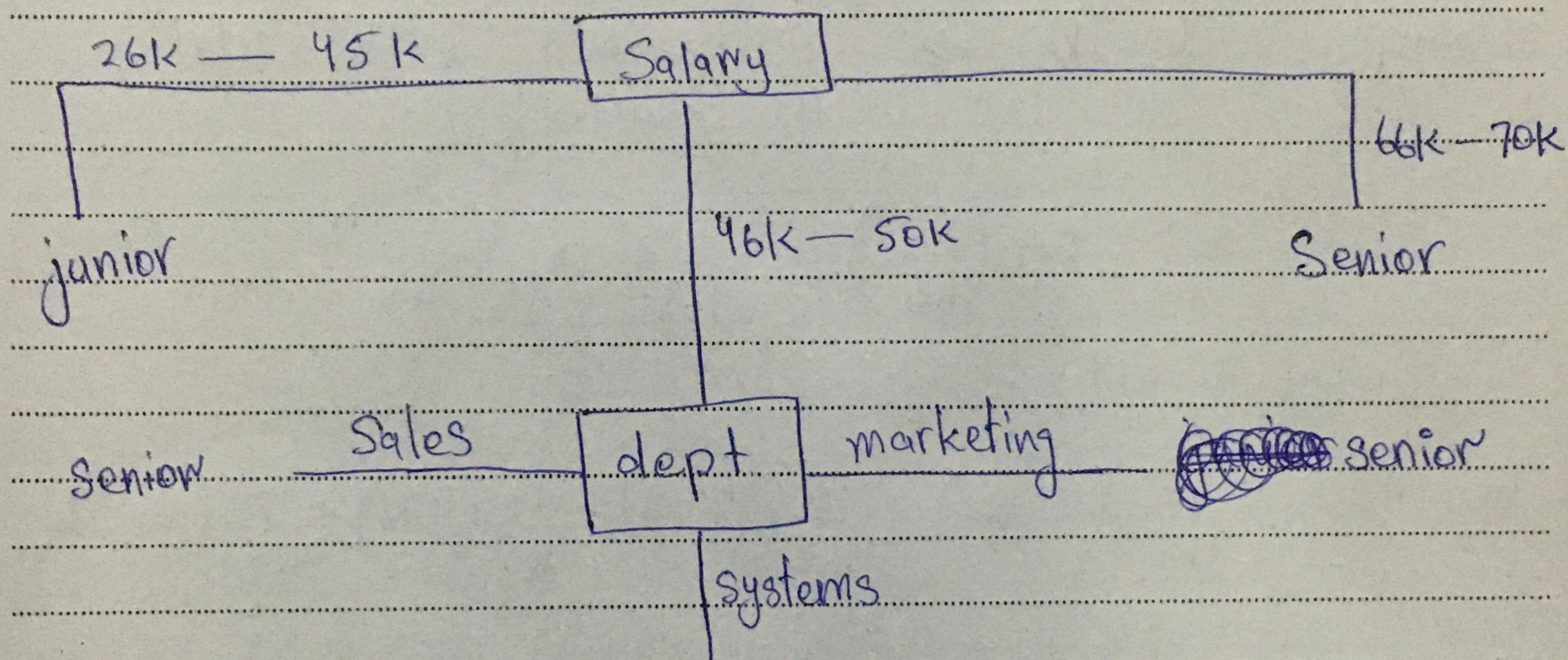
Now, splitting Second table.

dept	Junior	Senior	Total	Gini
Sales	0	30	30	0
Systems	23	0	23	0
Marketing	0	10	10	0

$$\text{Gini}(\text{department}) = 0$$

Age	junior	Senior	Total	Gini Count
21-25	20	0	20	0
26-30	3	0	3	0
31-35	0	30	30	0
36-40	0	10	10	0

$$\text{Gini(Age)} = 0$$



\* Both department and age have gini 0, so chosen based on First come basis.

### IF - THEN RULE

IF Salary  $\geq 26K$  AND SALARY  $< 45K$   
 THEN JUNIOR

IF Salary  $\geq 46K$  AND SALARY  $< 50K$

IF dept = sales than senior

IF dept = systems than junior

IF dept = marketing than senior

IF Salary  $\geq 66K$  and Salary  $< 70K$   
 Then senior

C4.5

Department:

$$\text{Split info} = -\frac{110}{165} \log_2 \left( \frac{110}{165} \right) - \frac{3}{165} \log_2 \left( \frac{3}{165} \right) - \frac{14}{165} \log_2 \left( \frac{14}{165} \right) \\ - \frac{10}{165} \log_2 \left( \frac{10}{165} \right) \\ = 1.388$$

Entropy = Info(D)  $\Rightarrow$  { 52 Seniors } { 113 Juniors } } 165 total

$$= -\frac{52}{165} \log_2 \left( \frac{52}{165} \right) - \frac{113}{165} \log_2 \left( \frac{113}{165} \right) \\ = 0.894$$

$$\text{Info}(A) = \frac{110}{165} (30, 80) + \frac{31}{165} (23, 8) + \frac{14}{165} (4, 10) + \frac{10}{165} (6, 4) \\ = 0.83$$

$$\text{Gain} = \text{Info}(D) - \text{Info}(A) \\ = 0.894 - 0.83 = 0.063$$

$$\text{Gain ratio} = 0.063 / 1.388 = 0.045$$

Age

$$\text{Info}(A) = \frac{20}{165} (20, 0) + \frac{49}{165} (49, 0) + \frac{79}{165} (44, 35) + \frac{10}{165} (0, 10) \\ + \frac{3}{165} (0, 3) + \frac{4}{165} (0, 4) \\ = 0.47$$

$$\text{Splitting info} = 0.3890 + 0.520 + 0.508 + 0.24 + 0.1 + 0.13 = 1.87$$

$$\text{Gain} = 0.424$$

$$\text{Gain Ratio} = 0.424 / 1.87 = 0.224$$

Salary

$$\text{Info}(A) = 0.358$$

$$\text{Gain} = 0.894 - 0.358$$

$$= 0.536$$

$$\text{Splitting info} = 0.513 + 0.13 + 0.13 + 0.5 + 0.2 + 0.4 = 2.01$$

$$\text{Gain Ratio} = 0.266$$

Splitting on highest gain ratio, which is of salary (0.266)

26K - 45K

Salary

66K - 70K

junior

Senior

<u>Sales</u>				
Dept	Age	Status	Count	
Sales	31-35	senior	30	
systems	21-25	junior	20	
systems	26-30	junior	3	
marketing	36-40	senior	10	

$$\text{entropy} = 2 + 2 + 4$$

department :-

$$\text{Splitting} = -\frac{30}{63} \log_2 \left( \frac{30}{63} \right) - \frac{23}{63} \log_2 \left( \frac{23}{63} \right) - \frac{10}{63} \log_2 \left( \frac{10}{63} \right)$$

$$= 1.4611$$

$$\text{Info}(A) = \frac{30}{63} (0, 30) + \frac{23}{63} (23, 0) + \frac{10}{63} (0, 10)$$

$$= 0$$

$$\text{Gain} = 1$$

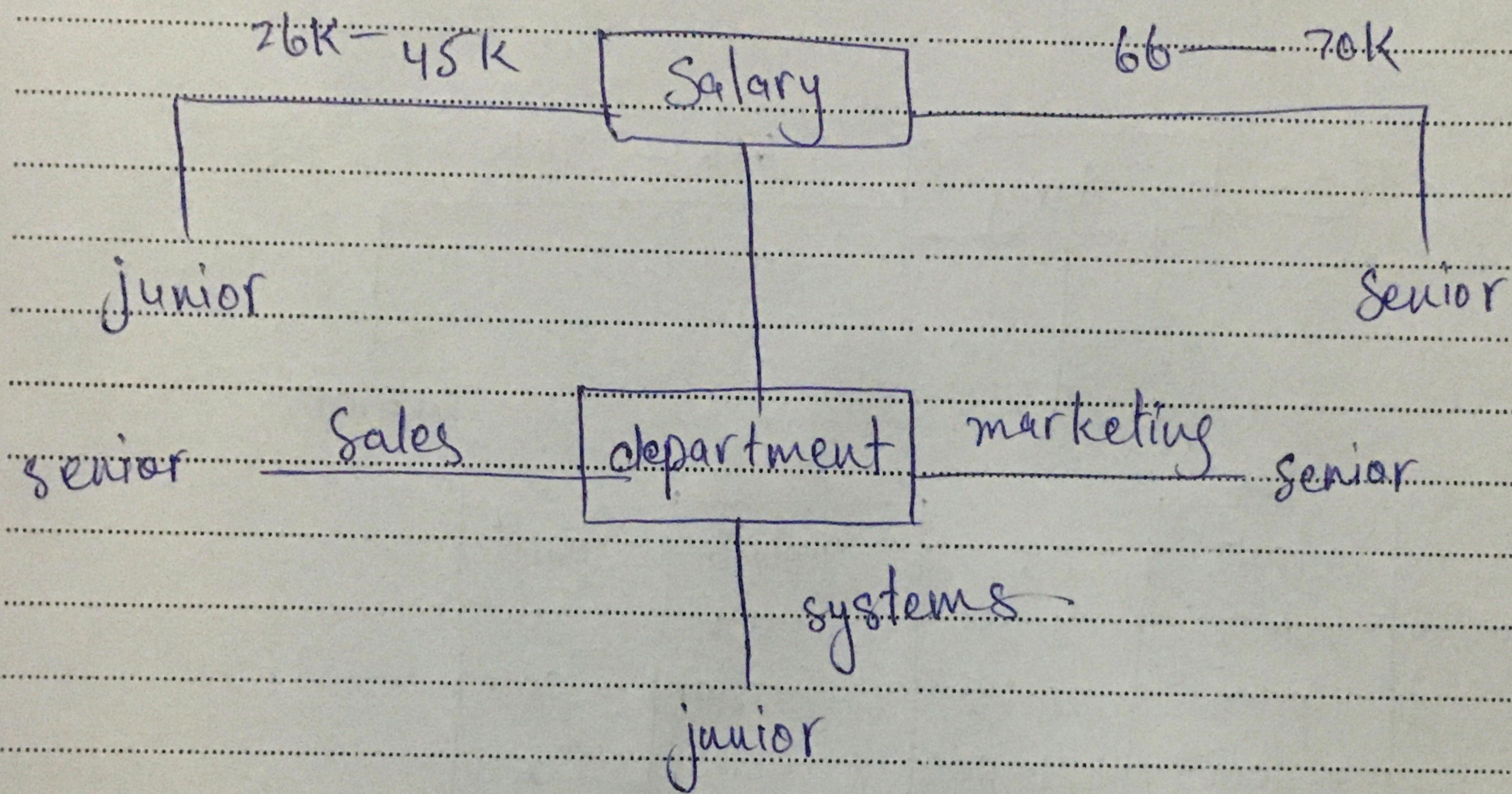
$$\text{Gain ratio} = \frac{1}{1.4611} = 0.68 \rightarrow \text{highest}$$

Age :-

$$\text{Splitting info} = -\frac{20}{63} \log_2\left(\frac{20}{63}\right) - \frac{3}{63} \log_2\left(\frac{3}{63}\right) - \frac{30}{63} \log_2\left(\frac{30}{63}\right)$$
$$= 1.664$$

$$\text{info}(A) = 6$$

$$\text{gain ratio} = \frac{1}{1.664} = 0.60$$



(c)

Classification

"systems", 26, — 30, 45 — 60K

Salary → dept → junior

So, we get a classification of junior for given ~~data~~ record according to dataset.

Question # 2

Class Label = Purchase.

- (A) 150030, Male, 28, 90000, ?  
 (B) 256899, Female, 40, 60000, ?

Using Male = 0, Female = 1 (is not included)

Gender, Age, Salary used in Euclidean distance.

$$= \sqrt{(Gender - gender_0)^2 + (Age - Age_0)^2 + (Salary - Salary_0)^2}$$

Distance

ID	Gender	Age	Salary	Purchased	A	B
15624510	0	19	19K	0	71K	41K
15810944	0	35	20K	0	70K	40K
15668575	1	26	43K	0	47K	17K
15603246	1	27	57K	0	33K	3K
15804002	0	19	76K	0	14K	16K
15728773	0	27	58K	0	32K	2K
15598044	1	27	84K	0	61K	24K
15894829	1	32	180K	1	60K	40K
15600575	0	25	33K	0	57K	27K
15727311	1	35	65K	0	25K	5K

K = 3

FOR A :

	Distance	Label	
①	14000	0	major is 0, so label
②	6000	0	for 150030 is 0
③	25000	0	

FOR B :

	Distance	Label	
①	3000	0	major is 0, so label
②	2000	0	for 256899 is 0
③	5000	0	

Class Label = Gender (Male = 0, Female = 1)

- (C) 150030, ?, 28, \$0000, 1
- (D) 256899, ?, 40, 60000, 0
- (E) 566989, ?, 15, 20000, 0.

Age, Salary, Purchased used in Euclidean Distance.

$$= \sqrt{(Age - Age_0)^2 + (Salary - Salary_0)^2 + (Purchased - Purchased_0)^2}$$

Age	Estimated Salary	Purchased	Gender	C	D	G
19	19K	0	0	71K	41K	1K
35	20K	0	0	70K	40K	20
26	43K	0	1	47K	17K	23K
27	57K	0	1	33K	30K	37K
19	76K	0	0	14K	16K	86K
27	58K	0	0	32K	2K	38K
27	84K	0	1	6K	24K	64K
32	150K	1	1	60K	90K	130K
25	33K	0	0	57K	27K	13K
35	65K	0	1	25K	5K	45K

K = 3

FOR C :	Distance	Label	
①	14000	0	
②	6000	1	
③	28000	1	

} major is 1, So label for 150030 is 1 (Female)

FOR D :	Distance	Label	
①	3000	1	
②	2000	1	
③	5000	1	

} major is 1, So label for 256899 is 1 (Female)

	Distance	Label	
①	1000	0	
②	20	0	
③	13000	0	

} major is 0, So Label for 566989 is 0 (Male)

## DISTANCE WEIGHTED KNN

By using  
Euclidean distance.

TABLE FOR A, B

Purchased	A	B	B
0	$1.4 \times 10^{-5}$		$2.44 \times 10^{-5}$
0	$1.42 \times 10^{-5}$		$2.5 \times 10^{-5}$
0	$2.13 \times 10^{-5}$		$5.88 \times 10^{-5}$
0	$3.03 \times 10^{-5}$		0.00033
0	$7.14 \times 10^{-5}$		$6.28 \times 10^{-5}$
0	$3.13 \times 10^{-5}$		0.0005
0	0.00017		$4.17 \times 10^{-5}$
1	$1.67 \times 10^{-5}$		$1.11 \times 10^{-5}$
0	$1.95 \times 10^{-5}$		$3.73 \times 10^{-5}$
0	$4 \times 10^{-5}$		0.0002

FOR A

$$\begin{aligned} \sum 0 &= 0.000407 \\ \sum 1 &= 0.0000167 \end{aligned} \quad \left. \begin{array}{l} \text{Since sum for 0 is larger, So} \\ \text{Label for 150030 is 0} \end{array} \right\}$$

FOR B

$$\begin{aligned} \sum 0 &= 0.00128 \\ \sum 1 &= 0.000011 \end{aligned} \quad \left. \begin{array}{l} \text{Since sum for 0 is larger, So} \\ \text{Label for 256899 is 0.} \end{array} \right\}$$

TABLE FOR C, D, E

Gender	C	D	E
0	$1.04 \times 10^{-5}$	$2.44 \times 10^{-5}$	0.001
0	$1.42 \times 10^{-5}$	$2.5 \times 10^{-5}$	0.05
1	$2.12 \times 10^{-5}$	$5.88 \times 10^{-5}$	$4.35 \times 10^{-5}$
1	$3.03 \times 10^{-5}$	0.00033	$2.7 \times 10^{-5}$
0	$7.14 \times 10^{-5}$	$6.25 \times 10^{-5}$	$1.8 \times 10^{-5}$
0	$3.125 \times 10^{-5}$	0.0005	$2.6 \times 10^{-5}$
1	0.00017	$4.17 \times 10^{-5}$	$1.6 \times 10^{-5}$
1	$1.67 \times 10^{-5}$	$1.11 \times 10^{-5}$	$7.7 \times 10^{-6}$
0	$1.75 \times 10^{-5}$	$3.7 \times 10^{-5}$	$7.7 \times 10^{-5}$
1	$4 \times 10^{-5}$	0.0002	$2.2 \times 10^{-5}$

FOR C

$$\begin{aligned} \sum 0 &= 0.00015 \\ \sum 1 &= 0.00027 \end{aligned} \quad \left. \begin{array}{l} \text{Since sum for 1 is large, So} \\ \text{Label for 150030 is 1 (Female)} \end{array} \right\}$$

FOR D

$$\begin{aligned} \sum 0 &= 0.0006489 \\ \sum 1 &= 0.0006449 \end{aligned} \quad \left. \begin{array}{l} \text{Since sum for 0 is large, So} \\ \text{Label for 256899 is 0 (Male)} \end{array} \right\}$$

FOR E

$$\begin{aligned} \sum 0 &= 0.05 \\ \sum 1 &= 0.0001 \end{aligned} \quad \left. \begin{array}{l} \text{Since sum for 0 is Large, So} \\ \text{Label for 566989 is 0 (male)} \end{array} \right\}$$