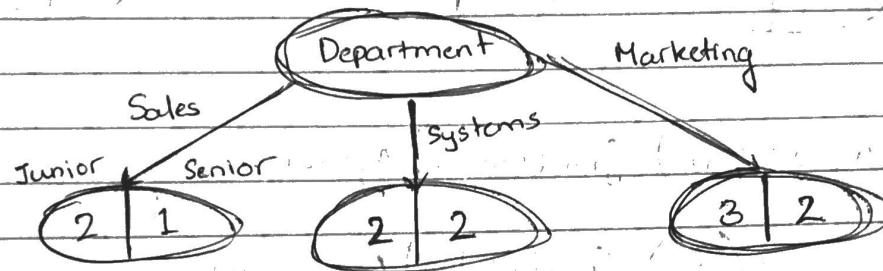


Machine Learning HW2, Anshuman Dikshit

① a Construct Decision Tree w/ multiway split

class: Status {Junior, Senior}

$n=12$



$$\text{Gini}(\text{Dept}) = \left(\frac{3}{12}\right) \cdot \text{Gini}(\text{Status} \mid \text{Dept} = \text{Sales}) +$$

$$\left(\frac{4}{12}\right) \cdot \text{Gini}(\text{Status} \mid \text{Dept} = \text{Systems}) +$$

$$\left(\frac{5}{12}\right) \cdot \text{Gini}(\text{Status} \mid \text{Dept} = \text{Marketing})$$

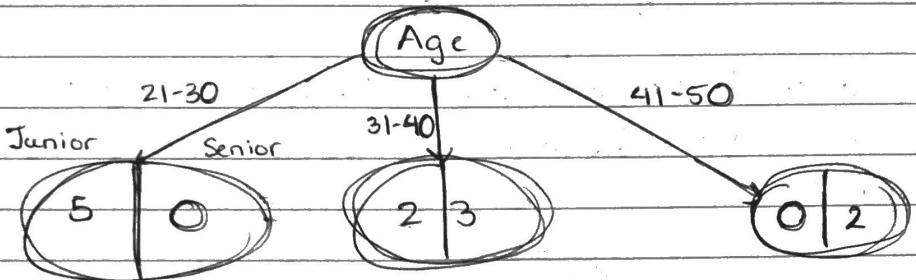
$$\text{Gini}(\text{Status} \mid \text{Dept} = \text{Sales}) = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 1 - 0.4444 - 0.0333 = 0.44$$

$$\text{Gini}(\text{Status} \mid \text{Dept} = \text{Systems}) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5$$

$$\text{Gini}(\text{Status} \mid \text{Dept} = \text{Marketing}) = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0.48$$

$$\text{GINI}(\text{Dept}) = \left(\frac{3}{12}\right) 0.44 + \left(\frac{4}{12}\right) 0.5 + \left(\frac{5}{12}\right) 0.48 = \boxed{0.47667}$$

0.47667



$$Gini(Age) = \left(\frac{5}{12}\right) \cdot Gini(\text{Status} | \text{Age}=21-30)$$

$$+ \left(\frac{5}{12}\right) \cdot Gini(\text{Status} | \text{Age}=31-40)$$

$$+ \left(\frac{2}{12}\right) \cdot Gini(\text{Status} | \text{Age}=41-50)$$

$$Gini(\text{Status} | \text{Age}=21-30) = 0$$

$$Gini(\text{Status} | \text{Age}=31-40) = 1 - \left(\frac{2}{5}\right)^2 = \left(\frac{3}{5}\right)^2 = 0.48$$

$$Gini(\text{Status} | \text{Age}=41-50) = 0$$

$$Gini(Age) = \frac{5}{12} (0.48) = 0.2$$



$$Gini(Salary) = \left(\frac{5}{12}\right) \cdot Gini(\text{Status} | \text{low})$$

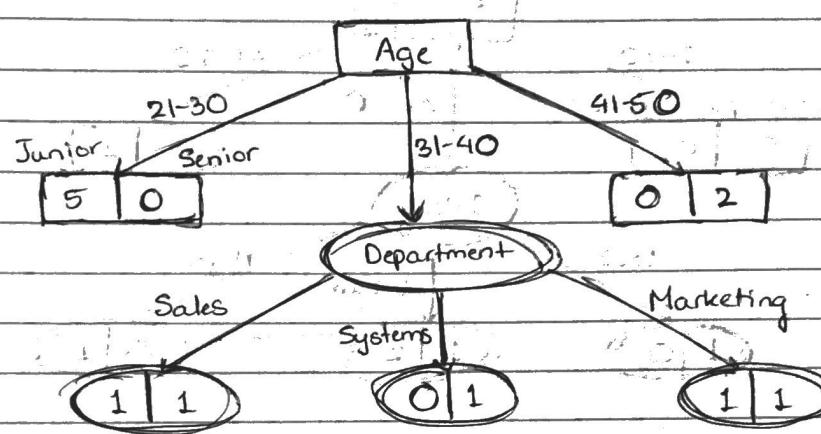
$$+ \left(\frac{1}{2}\right) \cdot Gini(\text{Status} | \text{medium})$$

$$+ \left(\frac{3}{12}\right) \cdot Gini(\text{Status} | \text{high})$$

$$Gini(\text{Status} | \text{medium}) = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 0.44$$

$$Gini(Salary) = \frac{1}{2} (0.44) = 0.22$$

Root Node: Age

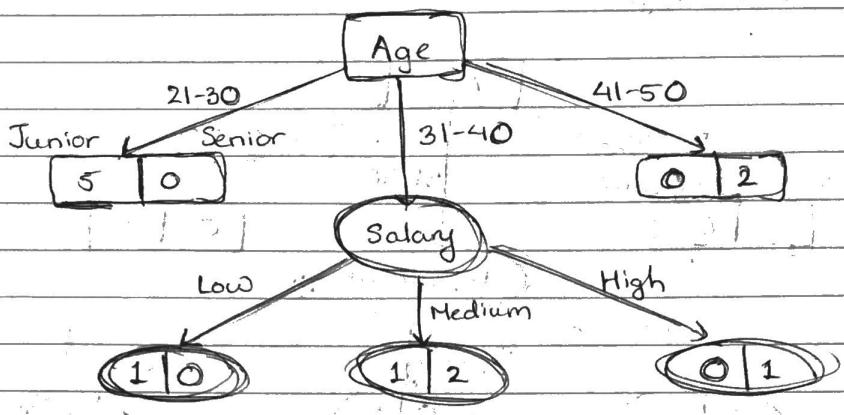


$$\text{Gini}(\text{Department}) = \left(\frac{2}{5}\right) \cdot \text{Gini}(\text{Status} \mid \text{Dept} = \text{Sales}) + \left(\frac{1}{5}\right) \cdot \text{Gini}(\text{Status} \mid \text{Dept} = \text{Systems}) + \left(\frac{2}{5}\right) \cdot \text{Gini}(\text{Status} \mid \text{Dept} = \text{Marketing})$$

$$\text{Gini}(\text{Status} \mid \text{Dept} = \text{Sales}) = 1 - (0.5)^2 - (0.5)^2 = 0.5$$

$$\text{Gini}(\text{Status} \mid \text{Dept} = \text{Marketing}) = 1 - (0.5)^2 - (0.5)^2 = 0.5$$

$$\text{Gini}(\text{Department}) = \left(\frac{2}{5}\right) \cdot \left(\frac{1}{2}\right) + 0 + \left(\frac{2}{5}\right) \cdot \left(\frac{1}{2}\right) = 0.4$$



$$\text{Gini}(\text{Salary}) = \left(\frac{1}{5}\right) \cdot \text{Gini}(\text{Status} \mid \text{Salary} = \text{Low})$$

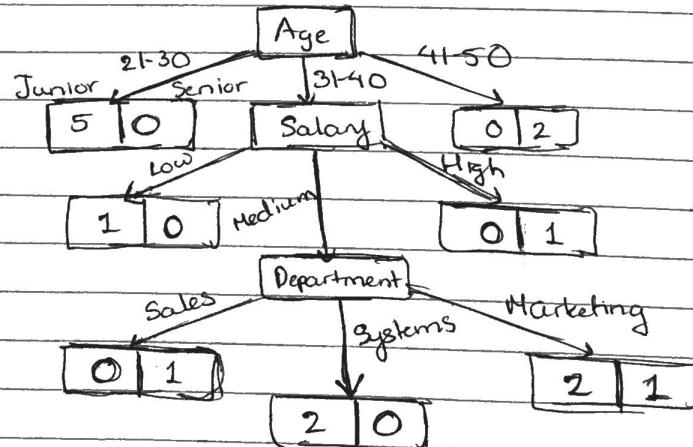
$$+ \left(\frac{3}{5}\right) \cdot \text{Gini}(\text{Salary} \mid \text{Salary} = \text{Medium})$$

$$+ \left(\frac{1}{5}\right) \cdot \text{Gini}(\text{Salary} \mid \text{Salary} = \text{High})$$

$$\text{Gini}(\text{Salary} \mid \text{Salary} = \text{Medium}) = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 0.44$$

$$\text{Gini}(\text{Salary}) = \frac{3}{5} \cdot (0.44) = \underline{0.264}$$

Final Decision Tree



$$\text{Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

$$\text{Precision} = \frac{\text{True Positive}}{\text{T.P.} + \text{F.P.}}$$

b

		Predicted Class	
		Class	Junior Senior
Actual	Junior	7	0
	Senior	1	4

$$\text{recall} = \frac{4}{5}$$

$$\text{precision} = \frac{4}{4}$$

$$F\text{-measure} = 2 \cdot \left( \frac{4}{5} \cdot \frac{4}{4} \right) = \frac{8}{5} = \frac{2}{8} = \frac{2}{\frac{8+20}{8}} = \frac{2}{\frac{28}{8}} = \frac{2}{\frac{7}{2}} = \frac{4}{7}$$

$$\left( \frac{4}{5} + \frac{4}{4} \right) = \frac{16+20}{20} = \frac{36}{20}$$

c	Sales	21-30	High	Junior
	Systems	21-30	Medium	Junior
	Marketing	41-50	High	Senior
	Marketing	31-40	Low	Junior

## ② Naive Bayes Classification

a We have to assume independence between all attributes

$x_i$  where a label is given:

$$P(x_1, x_2, \dots, x_d | Y_j) = P(x_1 | Y_j) \cdot P(x_2 | Y_j) \cdots P(x_d | Y_j)$$

If this stands, then we can estimate  $P(x_i | Y_j), \forall x_i, Y_j$  combinator for all training data

$$b P(A=0 | +) = \frac{2}{5} \quad P(A=0 | -) = \frac{3}{5}$$

$$P(A=1 | +) = \frac{3}{5} \quad P(A=1 | -) = \frac{2}{5}$$

$$P(B=0 | +) = \frac{4}{5} \quad P(B=0 | -) = \frac{3}{5}$$

$$P(B=1 | +) = \frac{1}{5} \quad P(B=1 | -) = \frac{2}{5}$$

$$P(C=0 | +) = \frac{1}{5} \quad P(C=0 | -) = \frac{4}{5}$$

$$P(C=1 | +) = \frac{4}{5} \quad P(C=1 | -) = \frac{1}{5}$$

c

$$P(+ | K) = \frac{(\frac{2}{5})(\frac{1}{5})(\frac{1}{5})(\frac{1}{2})}{K} = 0.008$$

$$P(- | K) = \frac{0}{K} = 0 \Rightarrow [+] \text{ is the class label}$$

d  $P(A=0|+) = \frac{4}{9}$   $P(A=0|-) = \frac{5}{9}$

$$P(A=1|+) = \frac{5}{9} \quad P(A=1|-) = \frac{4}{9}$$

$$P(B=0|+) = \frac{6}{9} \quad P(B=0|-) = \frac{5}{9}$$

$$P(B=1|+) = \frac{3}{9} \quad P(B=1|-) = \frac{9}{9}$$

$$P(C=0|+) = \frac{3}{9} \quad P(C=0|-) = \frac{2}{9}$$

$$P(C=1|+) = \frac{6}{9} \quad P(C=1|-) = \frac{7}{9}$$

e

$$P(+) = \frac{(4/9)(6/9)(6/9)(1/2)}{K} = \frac{0.247}{K}$$

↓

$$P(-) = \frac{(5/9)(4/9)(3/9)(1/2)}{K} = \frac{0.274}{K} \Rightarrow \boxed{\begin{matrix} - \text{ is class} \\ \text{label} \end{matrix}}$$

f The second method, if one of the conditional probabilities is 0, then the calculation of the class probability skews to 0. Via m-estimates we can ensure that doesn't happen