

Assignment - II
Data Mining & Warehouse
 — Abhishek Ray
 Roll-58 Department of CSE (3rd yr)

May 12, 2023

Q.

department	status	age	salary	count
Sales	Senior	31-35	46k-50k	30
Sales	Junior	26-30	26k-30k	40
Systems	Junior	31-35	31k-35k	40
Systems	Junior	21-25	46k-50k	20
Systems	Senior	31-35	66k-70k	5
Systems	Junior	26-30	46k-50k	3
Systems	Senior	41-45	66k-70k	3
Marketing	Senior	36-40	46k-50k	10
Marketing	Junior	31-35	41k-45k	4
Secretary	Senior	46-50	36k-40k	4
Secretary	Junior	26-30	26k-30k	6

a] How would you modify the basic decision tree algorithm to take into consideration the count of each generalized data tuple?

Ans → The basic decision tree algo should be modified as follows to take into consideration the count of each generalized data tuple —

- ① The count of each tuple must be integrated into the calculation of the attribute selection measure (such as information gain)
- ② Take the count into consideration to determine the most common class among the tuples.

b] Use your algo to construct a decision tree from given data.

Ans First we calculate Gini (Gini)-Index for entire dataset

Total-count = 150
 Sales-count = 110
 Systems-count = 28
 Marketing-count = 14
 Secretary-count = 10

$$\begin{aligned} \text{Gini-Total} &= \left[1 - \left(\left(\frac{\text{Sales-count}}{\text{total-count}} \right)^2 + \left(\frac{\text{Systems-count}}{\text{total-count}} \right)^2 + \left(\frac{\text{Marketing-count}}{\text{total-count}} \right)^2 + \left(\frac{\text{Secretary-count}}{\text{total-count}} \right)^2 \right) \right] \\ \text{Gini-total} &= \left[1 - (0.7332 + 0.1872 + 0.0932 + 0.0672) \right] \end{aligned}$$

→ now calculate Gini-Index for each attribute = 0.612

① department

→ Sales: Senior = 0

Junior = $\left[1 - \left(\left(\frac{40}{70} \right)^2 + \left(\frac{30}{70} \right)^2 \right) \right] = \text{0.489}$

→ System: $\text{senior} = \left[1 - \left(\left(\frac{5}{8} \right) \times 2 + \left(\frac{3}{8} \right) \times 2 \right) \right] = \boxed{0.469}$

$\text{junior} = \left[1 - \left(\left(\frac{23}{28} \right) \times 2 + \left(\frac{5}{28} \right) \times 2 \right) \right] = \boxed{0.408}$

→ Marketing: $\text{senior} = 0$ $\text{junior} = 0$

→ Secretary: $\text{senior} = \boxed{0.375}$ $\text{junior} = \boxed{0.5}$

→ Status: $\text{senior} \rightarrow$ $\text{Sales} = 0.469$ $\text{Systems} = 0.375$ $\text{Marketing} = 0.5$ $\text{Secretary} = 0$ | $\text{junior} \rightarrow$ $\text{Sales} = 0.489$ $\text{Systems} = 0.408$ $\text{Marketing} = 0$ $\text{Secretary} = 0.5$

→ Age: $(21-25) = 0$

$(26-30) \rightarrow$ $\text{Sales} = 0.489$ $\text{Systems} = 0.408$ $\text{Marketing} = 0$ $\text{Secretary} = 0.5$

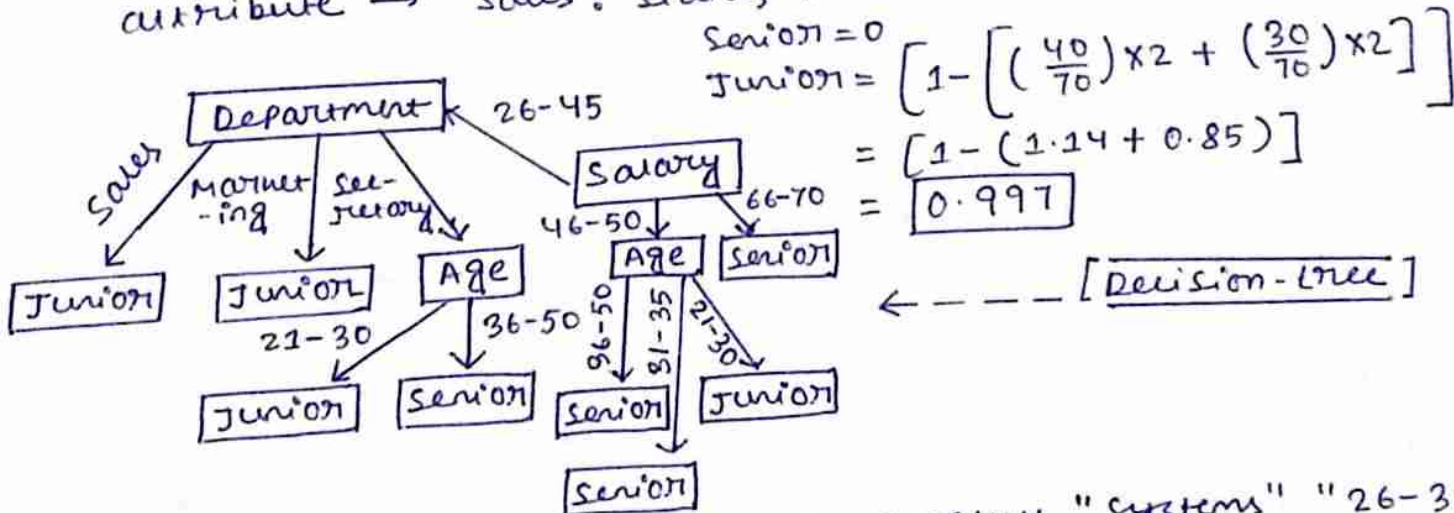
$(31-35) \rightarrow$ $\text{Sales} = 0.489$ $\text{Systems} = 0.469$ $\text{Marketing} = 0$ $\text{Secretary} = 0$ | $(36-40) = 0$ $(41-45) = 0$ $(46-50) \rightarrow$ $\text{Secretary} = 0.375$

→ Salary: $(26k-30k) \rightarrow$ $\text{Sales} = 0.489$ $\text{Systems} = 0$ $\text{Marketing} = 0$ $\text{Secretary} = 0.5$

$(31k-35k) \rightarrow$ $\text{Sales} = 0.489$ $\text{System} = 0.469$ $\text{Marketing} = 0$ $\text{Secretary} = 0$

$(36k-40k) \rightarrow$ $\text{Secretary} = 0.375$ $(41k-45k) \rightarrow$ $\text{Marketing} = 0$ $(46k-50k) \rightarrow$ $\text{Sales} = 0.489$ $\text{System} = 0.469$ $\text{Marketing} = 0$ $\text{Secretary} = 0$

→ Attribute of lowest Gini-Index is Department with value = $\boxed{0.373}$
Split dataset with based on department attribute → $\text{Sales} : \text{Status} \rightarrow$



c) Given a data-tuple having the values "systems", "26-30" & "46-50k" for the attributes department / age / salary respectively what would be the naive Bayesian classification?

→ Given a data tuples with the values — "system" "junior" "26-30" for the attribute department status / age respectively what would a naive Bayesian classification for salary tuple be →

$P(X | \text{senior}) = 0$ | $P(X | \text{junior}) = 0.018$

— Thus a naive Bayesian classification predicts Junior (Am)