

→ DMW Assignment 2

Q1 $P(C_i)$ - Prior Probability of each class

$$P(\text{lenses} = \text{'Non soft contact'}) = \frac{13}{23} \quad \left[\text{Adding 1 to avoid 0 available problem} \right]$$

$$P(\text{lenses} = \text{'soft contact'}) = \frac{5}{23}$$

$$P(\text{lenses} = \text{'hard contact'}) = \frac{5}{23}$$

$$X = \langle \text{Young, Myope, Yes, Reduced} \rangle$$

$$P(\text{Age} = \text{'Young'} \mid \text{lenses} = \text{'Non contact'}) = \frac{5}{13}$$

$$P(\text{Age} = \text{'Young'} \mid \text{lenses} = \text{'Soft contact'}) = \frac{3}{5}$$

$$P(\text{Age} = \text{'Young'} \mid \text{lenses} = \text{'Hard contact'}) = \frac{3}{5}$$

$$P(\text{Spectacle Pres} = \text{'Myope'} \mid \text{lenses} = \text{'Non contact'}) = \frac{8}{13}$$

$$P(\text{Spectacle Pres} = \text{'Myope'} \mid \text{lenses} = \text{'Soft contact'}) = \frac{3}{5}$$

$$P(\text{Spectacle Pres} = \text{'Myope'} \mid \text{lenses} = \text{'Hard contact'}) = \frac{4}{5}$$

$$P(\text{Astigmatism} = \text{'Yes'} \mid \text{lenses} = \text{'Non contact'}) = \frac{7}{13}$$

$$P(\text{Astigmatism} = \text{'Yes'} \mid \text{lenses} = \text{'Soft contact'}) = \frac{1}{5}$$

$$P(\text{Astigmatism} = \text{'Yes'} \mid \text{lenses} = \text{'Hard contact'}) = \frac{5}{5}$$

$$P(\text{Test Production Rate} = \text{"Reduced"} \mid \text{Lenses} = \text{"Non contact"}) = \frac{11}{13}$$

$$P(\text{Test Production Rate} = \text{"Reduced"} \mid \text{Lenses} = \text{"Soft contact"}) = \frac{1}{5}$$

$$P(\text{Test Production Rate} = \text{"Reduced"} \mid \text{Lenses} = \text{"Hard contact"}) = \frac{1}{5}$$

$$P(x \mid \text{Lenses} = \text{"Non contact"}) = \frac{5}{13} \times \frac{2}{13} \times \frac{7}{13} \times \frac{11}{13} = 0.108$$

$$P(x \mid \text{Lenses} = \text{"Soft contact"}) = \frac{3}{5} \times \frac{4}{5} \times \frac{1}{5} \times \frac{1}{5} = 0.014$$

$$P(x \mid \text{Lenses} = \text{"Hard contact"}) = \frac{3}{5} \times \frac{4}{5} \times \frac{5}{5} \times \frac{1}{5} = 0.096$$

$$P(\text{Lenses} = \text{"Non contact"} \mid x) = \frac{13}{23} \times 0.0924 = 0.0522$$

$$P(\text{Lenses} = \text{"Soft contact"} \mid x) = \frac{5}{23} \times 0.0144 = 0.0031$$

$$P(\text{Lenses} = \text{"Hard contact"} \mid x) = \frac{5}{23} \times 0.096 = 0.0209$$

x belongs to the class (Lenses = 'Non contact')

- 2a) The basic decision tree algorithm should be modified as follows to take up into consideration the count of each generalized data tuple
- The count of each tuple must be integrated into the calculation of the attribute selection measure (such as the information gain)
 - Take the count into consideration to determine the most common class among the tuples

b) No of Junior = 113
No of Senior = 52

$$H_0(D) = - \sum_{i=1}^m p_i \log_2(p_i)$$

$$H_0(D) = - \frac{113}{165} \log_2\left(\frac{113}{165}\right) - \frac{52}{165} \log_2\left(\frac{52}{165}\right)$$

$$= 0.8990$$

Department	p_i	N_i	Entropy
Sales	80	30	0.8454
System	23	8	0.8238
Marketing	4	10	0.8631
Secretary	6	4	0.9710

$$\text{Entropy for department} = \frac{110}{165} \times 0.8454 + \frac{31}{165} \times 0.8238$$

$$+ \frac{14}{165} \times 0.8631 + \frac{10}{165} \times 0.9710$$

$$= 0.8565$$

$$\text{Gain} = 0.899 - 0.8565$$

$$= 0.0485$$

Age	P	N	Entropy
21..25	20	0	0
26..30	49	0	0
31..35	44	35	0.9906
36..40	0	0	0
41..45	0	0	0
46..50	0	0	0

$$\text{Age entropy} = \frac{44+35}{165}$$

$$= \frac{79}{165} = 0.4794$$

$$\text{Gain} = 0.8990 - 0.4794$$

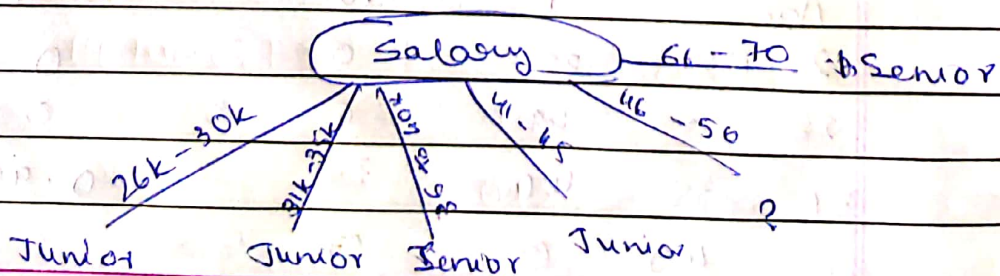
$$= 0.4196$$

Salary	P	N	Entropy
21 - 30K	46	0	0
31 - 35K	40	0	0
36 - 40K	0	4	0
41 - 45K	0	4	0
46 - 50K	23	40	0.9468
66 - 70K	0	8	0

$$H_0(\text{Salary}) \text{ Entropy} = \frac{23+40}{165} \times 0.9468$$

$$= 0.5375$$

Now gain (salary) > gain (age) > Gain (dept)
 \therefore Salary attribute is selected as splitting attribute.



$$Info D_i = \frac{-40}{63} \log\left(\frac{40}{63}\right) - \frac{23}{63} \log\left(\frac{23}{63}\right)$$

$$= 0.9168$$

Consider dept attribute

$$Info = \frac{30}{63} \log_2(sales) + \frac{23}{63} \times Info_{System}$$

$$+ \frac{10}{63} \times Info(marketing)$$

$$= 0$$

$$Gain = 0.9168$$

Consider age attribute

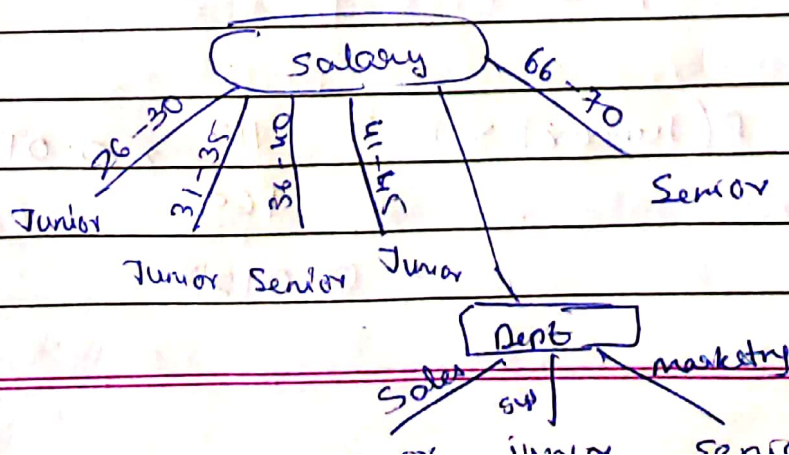
$$Info(D_i) = \frac{20}{63} Info(21-25) + \frac{13}{63} Info(26-30)$$

$$+ \frac{30}{63} Info(31-35) + \frac{10}{63} \times Info(36-40)$$

$$= 0$$

$$Gain = 0.9168$$

∴ Gain for both are same.



3) Prior probabilities

$$P(\text{senior}) = \frac{52}{165}$$

$$P(\text{junior}) = \frac{113}{165}$$

Consider attribute department

$$P(\text{system} | \text{senior}) = \frac{8}{52}$$

$$P(\text{systems} | \text{junior}) = \frac{23}{113}$$

Posterior Probabilities

$$P(x | \text{junior}) = \frac{23 \times 49 \times 23}{113 \times 113 \times 113} = 0.018$$

$$P(x | \text{senior}) = \frac{6 \times 0 \times 40}{52 \times 52 \times 52} = 0$$

$$\therefore \text{Bayes theorem } P\left(\frac{a}{x}\right) = \frac{P(a) \times P(x|a)}{P(x)}$$

$$P(\text{junior} | x) = \frac{113}{165} \times 0.018$$

$$= 0.0123$$

$$P(\text{Senior} | x) = \frac{0 \times 52}{165} = 0$$

$$\therefore 0.43 > 0$$

$x = (\text{System}, 26 \dots 30, 46 \dots 50k)$ will be predicted as status junior

3 To develop a quality classifier to guard against fraudulent credit and transactions with limited fraudulent cases, one can employ various techniques -

① Data collection - Gather a large data set of non-fraudulent transaction, collect a small but diverse set of fraudulent transactions. They must ensure that various types of frauds are represented.

② Data preprocessing
Clean and preprocess data by handling missing values and scaling features. There is imbalance in the dataset due to uneven number of fraudulent transaction, we must try and oversample fraudulent or undersample non fraudulent to balance

③ Classifier Selection

Choose appropriate ML algo for classifier such as random forest or gradient boosting.

④ Model Training.

Train classifier based on preprocessed balanced dataset so that it learns to distinguish between fraudulent and non fraudulent transaction.

⑤ Evaluation metrics.

Deploy suitable metrics such as AUC-ROC and precision, to evaluate the models performance.

⑥ Combine multiple models using bagging and boosting to improve classification accuracy

⑦ Feedback mechanism must be set up so that in case a fraudulent transaction is skipped by the model, feedback can be provided in order to improve and update the model.

4 Accuracy of a classifier on a given test set is the percentage of test set tuples that are correctly classified by the classifier.

$$\Rightarrow \text{accuracy} = \frac{TP + TN}{P + N}$$

TP - True Positive

TN - True Negative.

$$\text{Sensitivity or true positive Rate} = \frac{TP}{P}$$

$$\text{Sensitivity or True negativity Rate} = \frac{TN}{N}$$

$$P \times \text{Sensitivity} = TP - (4)$$

$$N \times \text{Specificity} = TN - (5)$$

Put (4), (5) in (1)

$$\frac{P}{P+N} \times \text{Sensitivity} + \frac{N}{P+N} \times \text{Specificity} = \text{accuracy}$$

Hence, accuracy of a classifier is a function of sensitivity and specificity.

5 Compare & contrast Bagging & Boosting

Bagging	Boosting
① The original dataset is divided into multiple subsets, selecting observations with replacement.	The new subset contains the components misclassified by the previous model.
② This method combines predictions that belong to the same type.	This method combines predictions that belong to different types.
③ Bagging decreases variance	Bagging decreases bias
④ Base classifiers are trained parallelly	Base classifiers are trained sequentially.
⑤ The models are created independently	The model creation is dependent on the previous one