

I(a) Explain the steps involved in the Data mining process with the sketch of the KDD process.

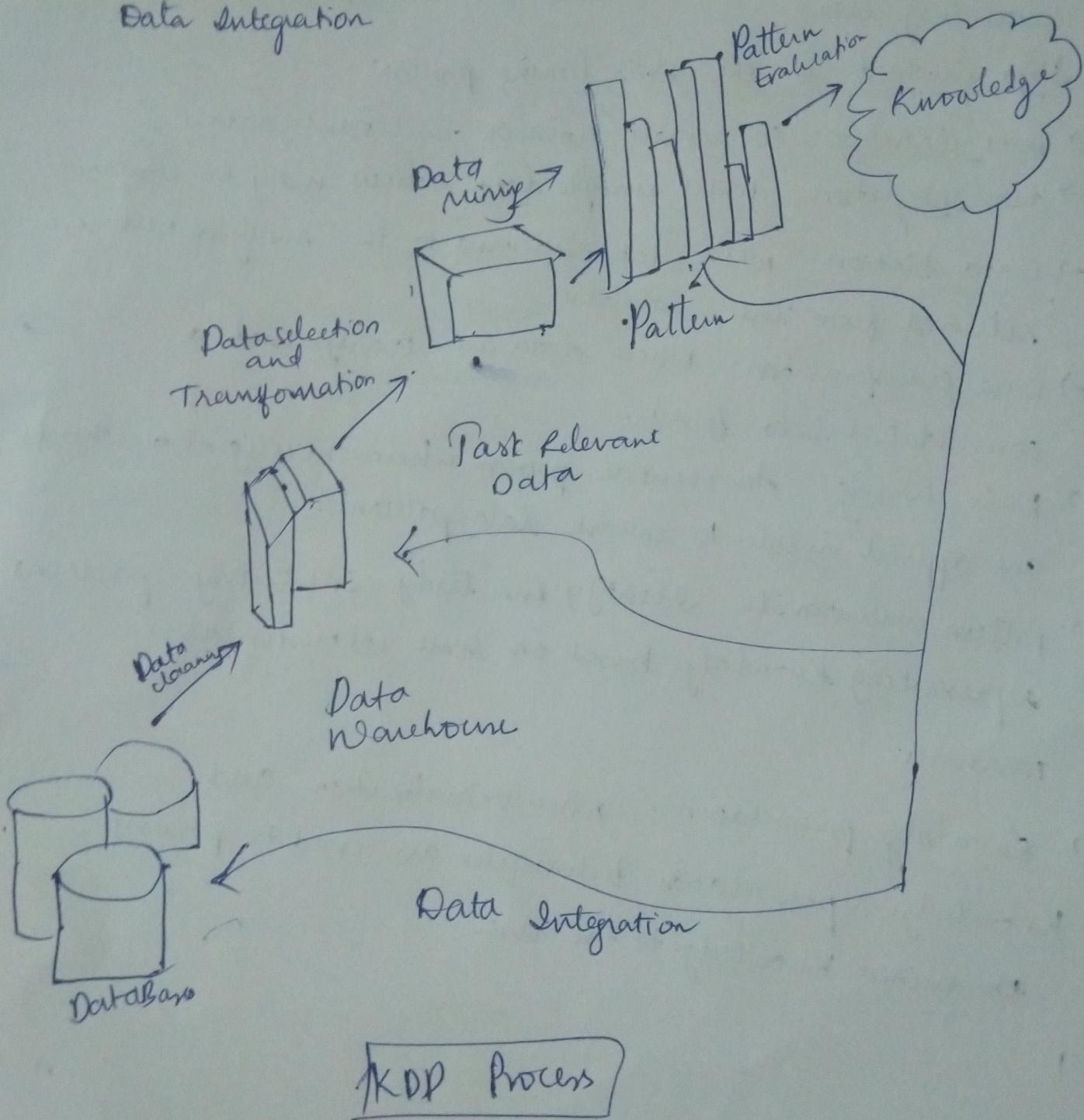
Data mining refers to extracting a "Mining" Knowledge from large amounts of data.

Steps involved in the Data mining process:

- Data cleaning: to remove noise or irrelevant data
- Data integration: where multiple data sources may be combined
- Data selection: where data relevant to the analysis task are retrieved from the database
- Data transformation: where data are transformed & consolidated into forms.
- Data mining: An essential process where intelligent methods are applied in order to extract data patterns.
- Pattern evaluation: to identify the truly interesting patterns representing knowledge based on some interestingness measures.
- Knowledge presentation: where visualization and knowledge representation techniques are used to present the mined knowledge to the user

KDDP: It is an iterative process where evaluation measures can be enhanced, mining can be refined, new data can be integrated and transformed in order to get different and more appropriate results.

→ preprocessing of databases consists of Data cleaning and Data integration



1(b) A sample data containing attributes of Mixed type

Obj Id	Attr-1 (Nominal)	Attr-2 (Ordinal)	Attr-3 (Numeric)
1	AA	Excellent	40
2	BB	Fair	34
3	AA	Excellent	45
4	CC	Good	65
5	BB	Fair	36
6	DD	Bad	22

Construct dissimilarity and similarity matrix:

Ans. Dissimilarity Matrix of Nominal Attr:

$$\text{Dissimilarity between } i \text{ and } j = d(i,j) = \frac{P-m}{P}$$

where P = No of attributes, m = no. of matches

Similarity Matrix:

	1	2	3	4	5	6
1	0					
2	$d(2,1)$	0				
3	$d(3,1)$	$d(3,2)$	0			
4	$d(4,1)$	$d(4,2)$	$d(4,3)$	0		
5	$d(5,1)$	$d(5,2)$	$d(5,3)$	$d(5,4)$	0	
6	$d(6,1)$	$d(6,2)$	$d(6,3)$	$d(6,4)$	$d(6,5)$	0

$$d(2,1) = \frac{1-0}{1} = 1$$

$$d(3,2) = \frac{1-0}{1} = 1$$

$$d(4,3) = \frac{1-0}{1} = 1$$

$$d(3,1) = \frac{1-1}{1} = 0$$

$$d(4,2) = \frac{1-0}{1} = 1$$

$$d(5,3) = \frac{1-0}{1} = 1$$

$$d(4,1) = \frac{1-0}{1} = 1$$

$$d(5,2) = \frac{1-1}{1} = 0$$

$$d(6,3) = 1$$

$$d(5,1) = \frac{1-0}{1} = 1$$

$$d(6,2) = \frac{1-1}{1} = 0$$

$$d(5,4) = 1$$

$$d(6,1) = 1$$

$$d(6,5) = \frac{1-0}{1} = 1$$

	1	2	3	4	5	6
1	0					
2	1	0				
3	0	1	0			
4	1	1	1	0		
5	1	0	1	1	0	
6	1	0	1	1	1	0

Ordinal Attributes:

To create Ordinal Matrix,

Step1: count states, There are 4, (Excellent, good, fair, bad).

Step2: Replace each ordinal data of test 2 by

Rank

Fair=1 Bad=2 Good=3 Excellent=4

Step3: Normalize the ranking $z_{ij} = \frac{x_{ij} - 1}{m_f - 1}$

$$\text{Fair}(1) = \frac{1-1}{4-1} = \frac{0}{3} = 0$$

$$\text{Bad}(2) = \frac{2-1}{4-1} = \frac{1}{3} = 0.3$$

$$\text{Good}(3) = \frac{3-1}{4-1} = \frac{2}{3} = 0.6$$

$$\text{Excellent}(4) = \frac{4-1}{4-1} = 1$$

objid	Test 2 (ordinal)
1	1
2	0
3	1
4	0.6
5	0
6	0.3

We make distance Matrix by using

$$d(x_i, y_j) = |x_{ij} - y_{ij}| + (x_{j2} - y_{j2}) + \dots + (x_{jn} - y_{jn})$$

Distance Matrix :

	1	2	3	4	5	6
1	0					
2	1	0				
3	0	1	0			
4	0.4	0.6	0.4	0		
5	1	0	1	0.6	0	
6	0.7	0.3	0.7	0.3	0.3	0

$$d(2,1) = |0-1| = 1 \quad d(3,1) = |1-1| = 0$$

$$d(3,2) = |1-1| = 0 \quad d(4,3) = |0.6-1| = 0.4$$

$$d(5,1) = |0-1| = 1 \quad d(4,2) = |0.6-0| = 0.6$$

$$d(5,4) = |0-0.6| = 0.6 \quad d(6,4) = |0.3-0| = 0.3$$

$$d(5,3) = |0-0| = 0 \quad d(5,1) = |0-1| = 1$$

$$d(5,2) = |0-0| = 0 \quad d(6,3) = |0.3-1| = 0.7$$

$$d(5,6) = |0-0.3| = 0.3$$

Now Numerical attribute

Test Id	Test 3
1	70
2	34
3	75
4	85
5	36
6	22

formula:

$$d_{ij} = \frac{|x_{if} - x_{jf}|}{\max - \min}$$

Here, $\max = 75$

$\min = 22$

Move to Normalize those rows so

can be Mapped $[0.0, 100]$

$$d(2,1) = \frac{|40-34|}{75-22}$$

$$= 0.67$$

$$d(3,2) = \frac{|45-34|}{75-22}$$

$$= 0.17$$

$$d(3,1) = \frac{|40-35|}{75-22}$$

$$= 0.09$$

$$d(4,2) = \frac{|34-65|}{75-22}$$

$$= 0.58$$

$$d(4,1) = \frac{|40-65|}{75-22} = 0.09$$

$$d(5,1) = \frac{|34-36|}{75-22} = 0.03$$

$$d(5,1) = \frac{|70-36|}{75-22} = 0.64$$

$$d(6,2) = \frac{|34-22|}{75-22} = 0.22$$

$$d(5,4) = \frac{|65-36|}{75-22} = 0.54$$

$$d(6,1) = \frac{|70-22|}{75-22} = 0.90$$

$$d(5,3) = \frac{|75-36|}{75-22} = 0.73$$

$$d(6,3) = \frac{|75-22|}{75-22} = 1$$

$$d(6,5) = \frac{|22-36|}{75-22} = 0.20$$

$$d(6,4) = \frac{|44-65|}{75-22} = 0.81$$

	1	2	3	4	5	6
1	0					
2	0.67	0.17				
3	0.09	0.77	0			
4	0.09	0.58	0.18	0		
5	0.64	0.03	0.73	0.54	0	
6	0.90	0.22	1	0.81	0.26	0

Similarity Matrix:

The Matrix is formed by combining the attributes of dissimilar matrices.

$$d(i,j) = \frac{\sum_{f=1}^p f_{ij} d_{ij}}{\sum_{j=1}^p f_{ij}}$$

where

$f_{ij} = 0$ if x_{if} or y_{jf} missing

or $x_{if} = x_{jf}$ and

d is assortive

$f_{ij} = 1$ otherwise

$$d(2,1) = \frac{(1 \times 1) + (1 \times 1) + (1 \times 0 \times 6)}{1 + 1} \\ = 0.89$$

$$d(3,1) = 0.3$$

$$d(4,1) = 0.49$$

$$d(5,1) = 0.88$$

$$d(6,1) = 0.86$$

$$d(6,2) = 0.97$$

$$d(2,3) = 0.92$$

$$d(4,2) = 0.72$$

$$d(5,2) = 0.01$$

$$d(4,3) = 0.52$$

$$d(5,4) = 0.71$$

$$d(5,3) = 0.91$$

$$d(6,4) = 0.70$$

$$d(6,3) = 0.9$$

$$d(6,5) = 0.52$$

Similarity Matrix:

	1	2	3	4	5	6
1	1					
2	0.89	0				
3	0.3	0.92	0			
4	0.49	0.76	0.52	0		
5	0.68	0.01	0.91	0.71	0	
6	0.86	0.17	0.9	0.70	0.52	0

Therefore for the given data with Mixed-type Attributes

Dissimilarity Matrix

	1	2	3	4	5	6
1	0					
2	0.67	0				
3	0.09	0.77	0			
4	0.09	0.58	0.18	0		
5	0.64	0.03	0.73	0.54	0	
6	0.90	0.22	1	0.81	0.26	0

Similarity Matrix

	1	2	3	4	5	6
1	0					
2	0.89	0				
3	0.3	0.92	0			
4	0.49	0.76	0.52	0		
5	0.68	0.01	0.91	0.71	0	
6	0.86	0.17	0.9	0.70	0.52	0

2a) Data cleaning:

Data cleaning routines attempt to fill missing values, smooth out noise while identifying outliers, and correct inconsistencies in the data.

(i) Approaches to fill missing values:

(a) Ignoring the tuple: This is usually done when the class label is missing (assuming the mining task involves classification or description). This method is not very effective.

(b) Manually filling the missing value: In general, this approach is time consuming and may not be a reasonable task for large datasets with many missing values, especially when the value to be filled in is not easily determined.

(c) Using a global constant to fill in the missing value:

Replace all missing attribute values by the same constant, such as label like "Unknown", or -∞. If missing values are replaced by, let "Unknown" then mining program thinks wrong. Although simple, it is not recommended.

(d) Using attribute mean for Quantitative (numeric) values or attributes mode for Categorical (nominal) values, for all samples belonging to the same class as the given tuple

(e) Using the most probable value to fill in the missing value:

This may be determined with regression, inference-based tools using Bayesian formalism, or decision tree induction.

(iii) Data Discretization:

Discretization techniques can be used to reduce the number of values for a given continuous attribute, by dividing the range of the attribute into intervals.

→ Interval labels can be used to replace actual data values.

Example :

Age: 10, 11, 13, 17, 19, 30, 31, 32, 38, 40, 42, 70, 72, 73, 75

Table Before discretization:

Attribute	Age	Age	Age
	10, 11, 13, 14, 17, 19	30, 31, 32, 38,	70, 72, 73, 75
After discretization	Young	Mature	Old

Concept hierarchy Generation:

A concept hierarchy for a given numeric attribute, defines a discretization of the attribute. Concept hierarchies can be used to reduce the data by collecting and replacing low level concepts (such as numeric values for the attribute age) by higher level concepts (such as, young, middle, or senior).

There are 5 methods for numeric concept hierarchy generation:

- Binning
- Histogram Analysis
- Clustering analysis
- Entropy-based Discretization and
- Data segmentation by "Natural Partitioning".

Q5) Interpret the below methods to normalize the following data

data: 200, 300, 400, 600, 1000

(i) Min-Max Normalization by setting min=0, max=1

min=0 $y = \text{value of Attribute}$
 max=1

$$v_1=200 \quad v_2=300 \quad v_3=400 \quad v_4=600 \quad v_5=1000$$

new max=1, new min=0

$$v_t = \left\{ \frac{y - \text{min } A}{\text{max } A - \text{min } A} \quad (\text{new max}_A - \text{new min}_A) \right\} + \text{new min } A$$

$$\text{for } 200: \text{min max} = \frac{200 - 200(1-0)}{1000 - 200} + 0 = 0.$$

$$\text{for } 300: \text{min max} = \frac{300 - 200(1-0)}{1000 - 200} + 0 = \frac{100}{800} = 0.125$$

$$\text{for } 400: \text{min max} = \frac{400 - 200(1-0)}{800} + 0 = 0.25$$

$$\text{for } 600: \text{min max} = \frac{600 - 200(1-0)}{800} = \frac{400}{800} = 0.5$$

$$\text{for } 1000: \text{min max} = \frac{1000 - 200(1-0)}{1000 - 200} + 0 = 1$$

Original Data:	200	300	400	600	1000
Normalized data (min,max) (0,1)	0	0.125	0.25	0.5	1

(ii) Z-Score Normalization:

Standard deviation:

Now Mean of Data =

$$\frac{\Sigma (\text{every value}) - (\text{mean of data})}{n}$$

∴ Standard deviation =

$$\begin{aligned} &= \sqrt{\frac{(200-500)^2 + (300-500)^2 + (400-500)^2 + (500-500)^2 + (600-500)^2}{5}} \\ &= \sqrt{\frac{10000 + 40000 + 10000 + 10000 + 20000}{5}} \\ &= \sqrt{\frac{100000}{5}} = \sqrt{20000} = 282.8 \end{aligned}$$

$$\begin{aligned} \text{Z-Score of } 200 &= \frac{x - \mu}{\sigma} = \frac{200 - 500}{282.8} = -1.06 \end{aligned}$$

$$\begin{aligned} \text{Z-Score of } 300 &= \frac{300 - 500}{282.8} = -0.7 \end{aligned}$$

$$\begin{aligned} \text{Z-Score of } 400 &= \frac{400 - 500}{282.8} = -0.35 \end{aligned}$$

$$\begin{aligned} \text{Z-Score of } 600 &= \frac{600 - 500}{282.8} = 0.35 \end{aligned}$$

$$\begin{aligned} \text{Z-Score of } 1000 &= \frac{1000 - 500}{282.8} = 1.78 \end{aligned}$$

Original data	200	300	400	600	1000
Z-Score (Normalized data)	-1.06	-0.7	-0.35	0.35	1.78

$$(iii) \text{ Mean} = \frac{200 + 300 + 400 + 600 + 1000}{5} = 500$$

$$\text{Mean Absolute deviation} : \frac{1}{n} \sum_{i=1}^n |x_i - m(x)|$$

where $m(x)$ = Average value of Dataset

n = no. of data values

x_i = data values in set

$$= \frac{|200-500| + |300-500| + |400-500| + |600-500| + |1000-500|}{5}$$

$$MAD = \frac{300 + 200 + 100 + 100 + 500}{5} = \frac{1200}{5} = 240$$

Zscore
with MAD

$$\text{for } 200 : \frac{200-500}{240} \approx -1.25$$

$$\text{for } 300 : \frac{300-500}{240}, \approx -0.83$$

$$\text{for } 400 : \frac{400-500}{240} = -0.832.45$$

$$\text{for } 600 : \frac{600-500}{240}, 2.45 0.41$$

$$\text{for } 1000 : \frac{1000-500}{240} \rightarrow 2.08$$

Original Data	200	300	400	600	1000
Normalized Zscores	-1.25	-0.83	-0.832.45	0.41	2.08

(iv)

Decimal scaling: $V' = \lceil \cdot / 10^j \rceil$

where $j = \text{its smallest integer such that } \max(V_i) \leq 10^j$.

$$V = 200 \rightarrow \lceil 200 / 10^2 \rceil = 0.2$$

$$V = 300 \rightarrow \lceil 300 / 10^3 \rceil = 0.3$$

$$V = 400 \rightarrow \lceil 400 / 10^2 \rceil = 0.4$$

$$V = 600 \rightarrow \lceil 600 / 10^3 \rceil = 0.6$$

$$V = 1000 \rightarrow \lceil 1000 / 10^3 \rceil = 1$$

Original data	200	300	400	600	1000
Normalized data	0.2	0.3	0.4	0.6	1

Ques 3a) List the differences between classification and prediction

3a)

classification

Prediction

1) Classification is the process of identifying to which category, a new observation belongs to on the basis of a training data set containing observations whose category membership is known.

→ In classification, the accuracy depends on finding the label correctly.

→ A model or the classifier is constructed to find the categorical labels.

→ In classification, the model can be known as the classifier.

→ Classification is mostly based on our current or past assumptions.

1) Prediction is the process of identifying the missing or unavailable numerical data for a new observation.

→ In prediction, accuracy depends on how well a given predictor can guess the value of predicted attribute for new data.

→ A model or a predictor will be constructed that predicts a continuous-valued function or ordered value.

→ In prediction, model can be known as the predictor.

→ Prediction is based on our future assumptions.

3b) Make use of the training examples shown in below table for a binary classification problem to compute below operations.

Instance	a_1	a_2	a_3	Target class
1	T	T	1.0	+
2	T	T	6.0	+
3	T	F	5.0	-
4	F	F	4.0	+
5	F	T	7.0	-
6	F	T	3.0	-
7	F	F	8.0	-
8	T	F	7.0	+
9	F	T	5.0	-

a) What is the entropy of this collection of training examples with respect to the class attribute?

$$\begin{aligned}
 \text{Entropy} &= - \sum_{i=0}^{t-1} p(i|t) \log_2 p(i|t) \\
 &= - \left[\left(\frac{4}{9} \right) * \log_2 \left(\frac{4}{9} \right) + \left(\frac{5}{9} \right) * \log_2 \left(\frac{5}{9} \right) \right] \\
 &= - \left[(-0.51997) + (-0.49111) \right] \\
 &= \underline{\underline{0.99107}}
 \end{aligned}$$

b) What are the information gains of a_1 and a_2 relative to these training examples

$$\begin{aligned}
 \text{Entropy}(a_1) &= - \sum_{i=0}^{t-1} p(i|t) \log_2 p(i|t) \\
 &= - \left[\left(\frac{3}{4} \right) * \log_2 \left(\frac{3}{4} \right) + \left(\frac{1}{4} \right) * \log_2 \left(\frac{1}{4} \right) \right] \\
 &= - [(-0.31128) + (-0.5)] \\
 &= \underline{\underline{0.81128}}
 \end{aligned}$$

$$\begin{aligned}\text{Entropy } (a_2) &= - \left[\left(\frac{2}{5}\right) * \log_2 \left(\frac{2}{5}\right) + \left(\frac{3}{5}\right) * \log_2 \left(\frac{3}{5}\right) \right] \\ &\approx - [(-0.52877) + (-0.44218)] \\ &= 0.91095\end{aligned}$$

$$\begin{aligned}\text{Entropy } (a_1) &= - \left[\left(\frac{1}{3}\right) * \log_2 \left(\frac{1}{3}\right) + \left(\frac{4}{3}\right) * \log_2 \left(\frac{4}{3}\right) \right] \\ &\approx - [(-0.46439) + (-0.25754)] \\ &= 0.72193\end{aligned}$$

$$\begin{aligned}\text{Entropy } (a_3) &= - \left[\left(\frac{2}{4}\right) * \log_2 \left(\frac{2}{4}\right) + \left(\frac{2}{4}\right) * \log_2 \left(\frac{2}{4}\right) \right] \\ &= - [(-0.5) + (-0.5)] = \underline{\underline{1}}\end{aligned}$$

Information gain

$$\begin{aligned}a_1: \Delta &= I(\text{parent}) - \sum_{j=1}^k \frac{N(v_j)}{N} I(v_j) \approx \\ &= 0.991 - \left[\left(\frac{4}{9}\right) * (0.81128) + \left(\frac{5}{9}\right) * (0.72193) \right] = \underline{\underline{0.229}} \\ a_2: &0.991 - \left[\left(\frac{5}{9}\right) * (0.971) + \left(\frac{4}{9}\right) * (1) \right] = \underline{\underline{0.007}}\end{aligned}$$

c) for a_3 , which is a continuous attribute, compute information gain for every possible split.

sorted values	1	3	4	5	6	7	8	
split positions	0.5	2	3.5	4.5	5.5	6.5	7.5	8.5
+	0	4	1	3	1	3	2	1
-	0	5	0	5	1	4	1	4
Gain	0	0.143	0.00248	0.0727	0.00113	0.01827	0.1021	0.

split number 1

$$\begin{aligned}\text{Entropy}(t) &= -\left[\left(\frac{4}{9}\right)*\log_2\left(\frac{4}{9}\right)+\left(\frac{5}{9}\right)*\log_2\left(\frac{5}{9}\right)\right] \\ &= -[-0.51997]+(-0.47111) \\ &= 0.99107\end{aligned}$$

Information Gain $\boxed{\Delta = 0}$

split number 2

$$\leq \text{entropy}(t) = -\left[\left(\frac{1}{7}\right)*\log_2\left(\frac{1}{7}\right)+0*\log_2(0)\right] = 0$$

$$\geq \text{entropy}(t) = -\left[\left(\frac{3}{8}\right)*\log_2\left(\frac{3}{8}\right)+\left(\frac{5}{8}\right)*\log_2\left(\frac{5}{8}\right)\right] \\ = -[-0.53064]+(-0.42379) \\ = 0.95443$$

$$\text{Weighted Average} = \left[\left(\frac{1}{9}\right)*0\right]+\left[\left(\frac{8}{9}\right)*0.95443\right] = 0.84899$$

$$\text{Information Gain} = \Delta = 0.991 - 0.84899 = \boxed{0.143} = \Delta$$

split Number 3

$$\leq \text{entropy}(t) = -\left[\left(\frac{1}{2}\right)*\log_2\left(\frac{1}{2}\right)+\left(\frac{1}{2}\right)*\log_2\left(\frac{1}{2}\right)\right] \\ = -[-0.5]+(-0.5) = 1$$

$$\geq \text{entropy}(t) = -\left[\left(\frac{3}{7}\right)*\log_2\left(\frac{3}{7}\right)+\left(\frac{4}{7}\right)*\log_2\left(\frac{4}{7}\right)\right] \\ = -[-0.52388]+(-0.46135) \\ = 0.98523$$

$$\text{Weighted Average} = \left[\left(\frac{2}{9}\right)*1\right]+\left[\left(\frac{7}{9}\right)*(0.98523)\right] \\ = 0.988512$$

Information gain:

$$\therefore \Delta = 0.991 - (0.988512) = \underline{\underline{0.002488}}$$

Split number 4

$$\leq \text{entropy}(t) = -\left[\left(\frac{2}{3}\right) * \log_2\left(\frac{2}{3}\right) + \left(\frac{1}{3}\right) * \log_2\left(\frac{1}{3}\right)\right]$$

$$= -[(-0.38998) + (-0.52832)]$$

$$= 0.9183$$

$$> \text{entropy}(t) = -\left[\left(\frac{2}{6}\right) * \log_2\left(\frac{2}{6}\right) + \left(\frac{4}{6}\right) * \log_2\left(\frac{4}{6}\right)\right]$$

$$= -[(-0.52388) + (-0.38998)]$$

$$= 0.9183$$

Weighted Average: $\left[\left(\frac{3}{9}\right) * 0.9183\right] + \left[\left(\frac{6}{9}\right) * 0.9183\right] = \underline{\underline{0.9183}}$

Information gain:

$$0.991 - (0.9183) = \boxed{\Delta = 0.0127}$$

Split Number 5

$$\leq \text{Entropy}(t) = -\left[\left(\frac{2}{5}\right) * \log_2\left(\frac{2}{5}\right) + \left(\frac{3}{5}\right) * \log_2\left(\frac{3}{5}\right)\right]$$

$$= -[(-0.52871) + (-0.44218)]$$

$$= 0.97095$$

$$> \text{Entropy}(t) = -\left[\left(\frac{2}{9}\right) * \log_2\left(\frac{2}{9}\right) + \left(\frac{2}{9}\right) * \log_2\left(\frac{2}{9}\right)\right]$$

$$= -[(-0.5) + (-0.5)] = 1$$

Weighted Average: $\left[\left(\frac{5}{9}\right) * (0.97095)\right] + \left[\left(\frac{4}{9}\right) * 1\right] = 0.96386$

Information gain:

$$\Delta = 0.991 - (0.96386) = \underline{\underline{0.007139}}$$

Split

Split Number 6

$$\leq \text{Entropy}(t) = - \left[\left(\frac{3}{6}\right) * \log_2\left(\frac{3}{6}\right) + \left(\frac{3}{6}\right) * \log_2\left(\frac{3}{6}\right) \right] = 1$$

$$> \text{Entropy}(t) = - \left[\left(\frac{1}{3}\right) * \log_2\left(\frac{1}{3}\right) + \left(\frac{2}{3}\right) * \log_2\left(\frac{2}{3}\right) \right]$$

$$= \left[(-0.52388) + (-0.38998) \right]$$

$$= 0.9183$$

$$\text{Weighted Average: } \left[\left(\frac{6}{9}\right) * 1 \right] + \left[\left(\frac{3}{9}\right) * 0.9183 \right] = 0.922765$$

$$\text{Information gain: } \Delta = 0.991 - (-0.8889) = 0.10211 \quad \underline{0.018235}$$

Split number 7

$$\leq \text{Entropy}(t) = - \left[\left(\frac{4}{8}\right) * \log_2\left(\frac{4}{8}\right) + \left(\frac{4}{8}\right) * \log_2\left(\frac{4}{8}\right) \right]$$

$$= - [(-0.5) + (-0.5)] = 1$$

$$> \text{Entropy}(t) = - \left[\left(\frac{0}{7}\right) * \log_2\left(\frac{0}{7}\right) + \left(\frac{1}{7}\right) * \log_2\left(\frac{1}{7}\right) \right] = 0$$

$$\text{Weighted Average} = \left[\left(\frac{8}{9}\right) * 1 \right] + \left[\left(\frac{1}{9}\right) * 0 \right] = 0.8889$$

$$\text{Information Gain} = 0.991 + 0.8889 = 0.10211$$

Split Number 8

$$\leq \text{Entropy}(t) = - \left[\left(\frac{4}{9}\right) * \log_2\left(\frac{4}{9}\right) + \left(\frac{5}{9}\right) * \log_2\left(\frac{5}{9}\right) \right]$$

$$= - [(-0.51991) + (-0.47111)] = 0.99108$$

$$> \text{Entropy}(t) = - \left[\left(\frac{0}{0}\right) * \log_2 \dots \right] = 0$$

Information Gain $\Delta = 0$

d) What is the best split (between a_1, a_2, a_3) according to the information gain?

→ According to the information gain, the best split is a_1 due to its higher gain in comparison to a_2 and a_3 .

e) What is the best split (between a_1, a_2) according to the classification error rate?

→ According to the classification error rate, the best split is a_1 due to a lower classification error in comparison to a_2 .

The classification error depicts the accuracy of the sample set; the higher the classification error, the more error the sample contains.

$$a_1: \text{classification error}(t) = 1 - \max [P(C_i|t)] \\ = 1 - \left[\frac{7}{9}, \frac{2}{9} \right] = 1 - \frac{7}{9} = \frac{2}{9} = \underline{\underline{0.222}}$$

$$a_2: \text{classification error}(t) = 1 - \left[\frac{5}{9}, \frac{4}{9} \right] \\ = 1 - \frac{5}{9} = \frac{4}{9} = \underline{\underline{0.444}}$$

f) what is the best split (between a_1 and a_2) according to Gini Index?

According to Gini Index a_1 is the best split because it got lesser Gini Index value

$$a_1: \text{gini}(t) = 1 - \sum_{i=0}^{t-1} [P(C_i|t)]^2 = 1 - \left[\left(\frac{7}{9} \right)^2 + \left(\frac{2}{9} \right)^2 \right] = \underline{\underline{0.3457}}$$

$$a_2: \text{gini}(t) = 1 - \left[\left(\frac{4}{9} \right)^2 + \left(\frac{5}{9} \right)^2 \right] = 1 - \left(\frac{16}{81} + \frac{25}{81} \right) = \underline{\underline{0.4938}}$$