

Bangladesh-Bharat Digital Service and Employment Training

Test-3

Total marks: 100

Date: 14.09.2022

Time: 2 hours

Part-I: Choose the correct option and justify your answer with one or two sentence(s)

(20 X 2 = 40)

1. The unsupervised classification is known as
- Clustering
 - Correlation
 - Classification
 - Logistic regression

Correct Answer: a

Explanation: The unsupervised classification is termed as clustering algorithm.

2. Under which of the following categories does the Bayesian classification method fall?
- Statistical-based methods
 - Distance-based method
 - Error-based method
 - Decision-tree based method

Correct Answer: a

Explanation: Bayesian classifier is a statistical classifier, which performs probabilistic prediction, that is, it predicts class membership probabilities.

3. Which of the following is true about the K-Nearest Neighbour (KNN) classifier?
- It is considered as a lazy learner
 - It is considered as an eager learner
 - The learning strategy of KNN depends upon the dataset
 - The value of k is considered as a constant for all datasets, and always equal to 1

Correct Answer: a

Explanation: KNN is a lazy learner, and not an eager learner. For a discussion on lazy/eager learning of KNN, check. The learning strategy of KNN is fixed, and it does not depend on which dataset upon which KNN is being used. Also, the value of an optimal k is not a constant.

4. Which of the following statement is true about the entropy of an n-dimensional data distributed over k distinct classes?
- The lowest possible value of entropy is -1.
 - The highest possible value of entropy is ∞ .
 - The entropy can be any value between $-\infty$ and $+\infty$.
 - The entropy is always a positive quantity.

Correct Answer: d

Explanation: The entropy is ranges from 0 to 1, hence it is always positive.

5. In the M-estimation approach of Naïve Bayes Classification, the posterior probability is given by?
[Where, n = total number of instances from class, n_{C_i} = number of training examples from class C_i that takes the value $A_j = x$, m = equivalent sample size, p = a user-defined parameter]

- $P(A_j = x|C_i) = \frac{n_{C_i} - mp}{n - m}$
- $P(A_j = x|C_i) = \frac{n_{C_i} + mp}{n - m}$
- $P(A_j = x|C_i) = \frac{n_{C_i} - mp}{n + m}$
- $P(A_j = x|C_i) = \frac{n_{C_i} + mp}{n + m}$

Correct Answer: d

Explanation: The posterior probability in M-estimate approach is given as $P(A_j = x|C_i) = \frac{n_{C_i} + mp}{n + m}$.

6. In Naive Bayesian Classification technique, the M-estimation approach is used when,
- the posterior probability for one of the attributes is infinite
 - the posterior probability for one of the attributes is zero
 - one of the prior probabilities is zero
 - one of the prior probabilities is infinity

Correct Answer: b

Explanation: If the posterior probability for one of the attributes is zero, then the overall class-conditional probability for that class vanishes. In this case the M-estimation technique is used. This generally happens due to the low size of training data.

7. To study the effect of rain on the attendance of a particular school, the data for 10 days are collected. In the table below, the amount of rain of a particular day (Heavy/Medium/Low or H/M/L), and the attendance (High/Low or (H/L)) is provided. What is the value of class conditional probability of rain is High given that attendance is High (P(Rain = High | attendance = High))?

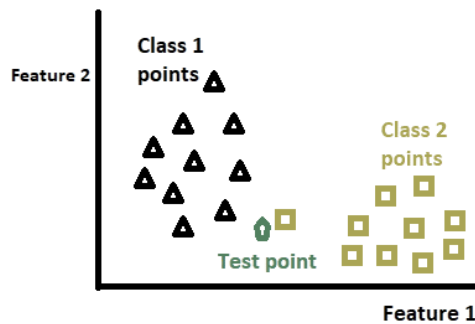
Amount of rain	H	H	L	M	L	M	H	L	L	M
Attendance	L	L	H	L	H	H	H	L	H	H

- $\frac{1}{3}$
- $\frac{1}{5}$
- $\frac{1}{6}$
- $\frac{1}{7}$

Correct Answer: c

Explanation: $P(\text{Rain} = \text{High} | \text{attendance} = \text{High}) = \frac{1}{6}$

8. Consider the following figure. The black triangles are the data points belonging to class 1. The grey rectangles are the data points belonging to class 2. The test point is marked as pentagon. What is the class label of the test point as per kNN classifier with $k = 1$ and $k = 3$?



- Class 1 and Class 1
- Class 1 and Class 2
- Class 2 and Class 1
- Class 2 and Class 2

Correct Answer: c

Explanation: If $k=1$, only one nearest neighbour of the test point is considered, which belongs to Class 2. So, the class label of the test point is Class 2. If $k=3$, three nearest neighbours of the test point is considered, belonging to Class 2, Class 1 and Class 1 respectively. So, the class label of the test point is Class 1, by majority.

9. Calculate Entropy for the data given in the table below-

Sample count	Sample class
9	1
5	2

- 0.94
- 0.94
- 0.36
- 0.36

Part-II: Solve the following problems

(10 X 6 = 60)

1. Consider the dataset shown in the **Table A**. Using the dataset predict the record $X = (\text{Age} = \text{young}, \text{Income} = \text{Medium}, \text{Married} = \text{yes}, \text{Health} = \text{Fair})$ belongs to a class?

Answer:

$$p_i = P(C_i) \times \prod_{j=1}^n P(A_j = a_j | C_i)$$

Calculation of $P(C_i)$

$$P(\text{Select} = \text{'Yes'}) = 9/14 = 0.643$$

$$P(\text{Select} = \text{'No'}) = 5/14 = 0.357$$

Calculation of $P(X | C_i)$ for each class C_i

$$P(\text{Age} = \text{'Young'} | \text{Select} = \text{'Yes'}) = 2/9 = 0.222$$

$$P(\text{Age} = \text{'Young'} | \text{Select} = \text{'No'}) = 3/5 = 0.6$$

$$P(\text{Income} = \text{'Medium'} | \text{Select} = \text{'Yes'}) = 4/9 = 0.444$$

$$P(\text{Income} = \text{'Medium'} | \text{Select} = \text{'No'}) = 2/5 = 0.4$$

$$P(\text{Married} = \text{'Yes'} | \text{Select} = \text{'Yes'}) = 6/9 = 0.667$$

$$P(\text{Married} = \text{'Yes'} | \text{Select} = \text{'No'}) = 1/5 = 0.2$$

$$P(\text{Health} = \text{'Fair'} | \text{Select} = \text{'Yes'}) = 6/9 = 0.667$$

$$P(\text{Health} = \text{'Fair'} | \text{Select} = \text{'No'}) = 2/5 = 0.4$$

Thus,

$$P(X | \text{Select} = \text{'Yes'}) = 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044$$

$$P(X | \text{Select} = \text{'No'}) = 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$$

$P(C_i) \times P(X | C_i)$:

$$P(\text{Select} = \text{'Yes'}) \times P(X | \text{Select} = \text{'Yes'}) = 0.643 \times 0.044 = \mathbf{0.028}$$

$$P(\text{Select} = \text{'No'}) \times P(X | \text{Select} = \text{'No'}) = 0.357 \times 0.019 = 0.007$$

The test data is belonging to yes class.

2. Consider the dataset shown in **Table B** and consider the test data [**Angelina, 5, F**], now find the class of sport for the test data for $k = 1$, and $k = 3$, using the **KNN** classification algorithm.

Answer:

The categorical data [**Male/Female**] converted into numerical data as [**1/2**] and we considered Euclidian distance.

Name	Age	Gender	Distance	Class
Ajay	32	1	27.02	Football
Mark	40	1	35.01	Neither
Sara	16	2	11	Cricket
Zaira	34	2	29	Cricket
Sachin	55	1	50.01	Neither

Rahul	40	1	35.01	Cricket
Pooja	20	2	15	Neither
Smith	15	1	10.05	Cricket
Laxmi	55	2	50	Football
Arun	15	1	10.05	Football

For $k=1$, the class is either **Cricket** or **Football**.

For $k=3$, the class is **Cricket**.

3. Consider a training data set as shown in the **Table C** and answer the following questions.
- Calculate the entropy of the data set.
 - Suppose, you select “Gender” as the splitting attribute. Calculate the following.
 - Information gain
 - Gini index
 - Gain ratio

Answer:

a) **Entropy:**

$$E = \sum_{i=1}^m -p_i \log_2 p_i$$

$$\text{Here, } p_1 = \frac{5}{15} = 0.3333, p_2 = \frac{8}{15} = 0.5333, \text{ and } p_3 = \frac{2}{15} = 0.1333$$

$$\therefore \text{Entropy} = \sum_{i=1}^3 -p_i \log_2 p_i = 0.3333 \times 0.4771 + 0.5333 \times 0.2730 + 0.1333 \times 0.8751 = 1.3996$$

b) For “Gender” as the splitting attribute:

i. **Information gain** = $\alpha(\text{Gender}, D) = E(D) - E_{\text{Gender}}(D)$

Here, $E(D) = 1.3996$ and

$$E_{\text{Gender}}(D) = \frac{9}{15} \times \left(-\frac{4}{9} \log \frac{4}{9} - \frac{5}{9} \log \frac{5}{9} \right) + \frac{6}{15} \times \left(-\frac{1}{6} \log \frac{1}{6} - \frac{3}{6} \log \frac{3}{6} - \frac{2}{6} \log \frac{2}{6} \right) = 1.17829$$

$$\text{Information gain} = \alpha(\text{Gender}, D) = 1.3996 - 1.17829 = \mathbf{0.2213}$$

ii. **Gini index** = $\gamma(A, D) = G(D) - G_A(D)$

$$G(D) = 1 - \left(\frac{5}{15} \right)^2 - \left(\frac{8}{15} \right)^2 - \left(\frac{2}{15} \right)^2 = 0.5867 \text{ and}$$

$$G_{\text{Gender}}(D) = \frac{9}{15} \times \left(1 - \left(\frac{4}{9} \right)^2 - \left(\frac{5}{9} \right)^2 \right) + \frac{6}{15} \times \left(1 - \left(\frac{1}{6} \right)^2 - \left(\frac{3}{6} \right)^2 - \left(\frac{2}{6} \right)^2 \right) = 0.5407$$

$$\therefore \text{Gini index} = 0.5867 - 0.5407 = \mathbf{0.046}$$

iii. **Gain ratio** = $\beta(\text{Gender}, D) = \frac{\alpha(\text{Gender}, D)}{E_{\text{Gender}}^*(D)}$

Now, $\alpha(\text{Gender}, D) = 0.2213$ and

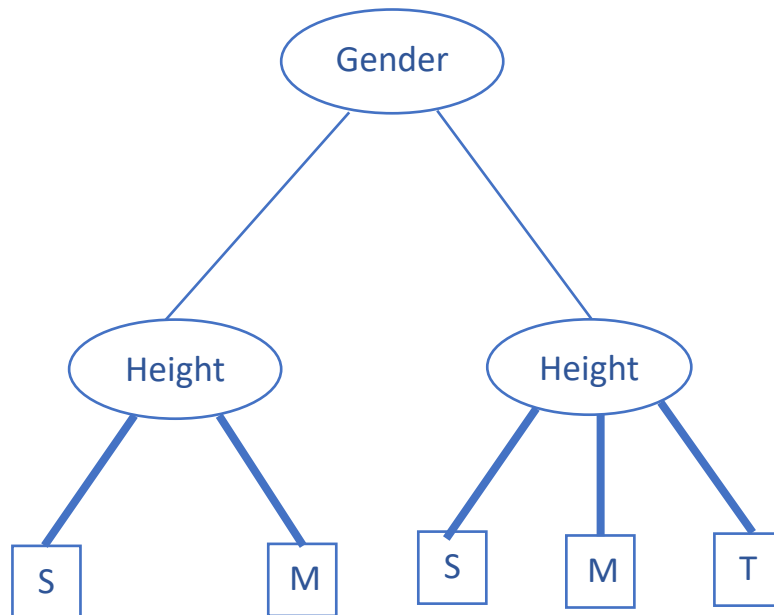
$$E_{\text{Gender}}^*(D) = - \sum_{j=1}^2 \frac{|D_j|}{|D|} \cdot \log \frac{|D_j|}{|D|} = -\frac{9}{15} \log \frac{9}{15} - \frac{6}{15} \log \frac{6}{15} = 0.97$$

$$\therefore \text{Gain Ratio} = \frac{0.2213}{0.97} = \mathbf{0.2281}$$

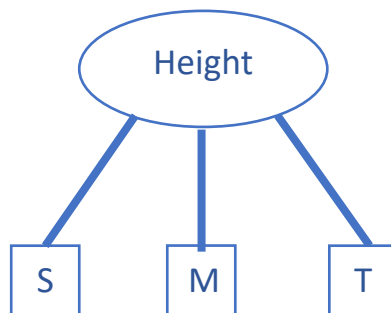
4. Consider the **Table C**, obtain the decision trees with the following splitting order and write the decision rules.
- Gender-Height and
 - Height –Gender.

Answer:

a) Decision tree according to splitting order **Gender-Height**



b) Decision tree according to splitting order **Height-Gender**



5. Consider the training data set shown in the **Table D**. The x and y coordinates of the training points and their corresponding class label are provided. The test point is given as (4,5). Predict the class label of the test point using 3-NN classifier, considering Manhattan Distance as proximity measure. **Note:** The Manhattan distance between the point (x_t, y_t) and any training point (x, y) is given by, $Manhattan\ Distance = |x - x_t| + |y - y_t|$.

Answer:

Datapoint	X	Y	Distance	Class
1	1	3	5	Yes
2	2	2	5	No
3	3	2	4	Yes
4	2	4	3	Yes
5	3	4	2	Yes
6	4	4	1	Yes
7	2	6	3	Yes
8	4	7	2	No
9	5	3	3	No
10	5	7	3	No
11	6	9	6	No
12	7	7	5	No

Test point will be in Yes class

6. A scheme of a training data is stated as below.

Symptom	Duration	Treatment	Class
S1, S2, S3	S (Short), M (Medium), L (Large)	A (Allopathy), H (Homeopath), U (Unani)	Cure (Y), Not Cure (N)

A contingency table is prepared with 400 records of patients, which is shown below.

		Class		Totals
		Y	N	
Symptom	S1	30	10	40
	S2	14	12	26
	S3	24	14	38
Duration	S	22	18	40
	M	36	26	62
	L	32	22	54
Treatment	A	30	24	54
	H	28	18	46
	U	24	16	40
Totals		240	160	400

A test data is given below:

S2	M	H	?
----	---	---	---

You have to classify the test data using the Naïve Bayes' classifier.

- What is the probability that the test data is in class Cure?
- What is the probability that the test data is in class Not Cure?
- In which class the test data will be?
 - The test belongs to class Cure.
 - The test data belongs to class Not Cure.
- What is the entropy of the input data?

Answer:

$$a) P(Y) = \frac{120}{200} \times \frac{7}{120} \times \frac{18}{120} \times \frac{14}{120} = 0.6 \times 0.058 \times 0.15 \times 0.116 = 0.0006$$

$$b) P(N) = \frac{80}{200} \times \frac{6}{80} \times \frac{13}{80} \times \frac{9}{80} = 0.4 \times 0.075 \times 0.162 \times 0.112 = 0.0005$$

c) The predicted class will be Y i.e., it will belong to class Cure.

d) Calculation of the entropy

$$\text{Here, } p_1 = \frac{240}{400} = 0.6 \text{ and } p_2 = \frac{160}{400} = 0.4$$

$$\text{Entropy} = -p_1 \log p_1 - p_2 \log p_2$$

$$= -0.6 \times (-0.7369) - 0.4 \times (-1.3219)$$

$$= 0.44214 + 0.52876$$

$$= \mathbf{0.9709}$$

---END---