# National University of Computer and Emerging Sciences
## Chiniot-Faisalabad Campus

## Data Mining (CS4038)

Date: April 8th 2024

**Course Instructor(s)**

Ms. Ayesha Liaqat

*Solution*

## Sessional-II Exam

| | |
|---|---|
| **Total Time (Hrs):** | 1 |
| **Total Marks:** | 60 |
| **Total Questions:** | 4 |

Roll No      Section      Student Signature

**Attempt all the questions.**

*CLO # 2: Understand the nature of the data and apply data mining techniques to interpret the results.*

**Q1:** Given a dataset where we want to predict whether a customer will purchase a product based on their gender, age group, and preferred payment method. **[10 marks]**

| Gender | Age Group | Payment Method | Purchase |
|---|---|---|---|
| Male | Teen | Credit Card | No |
| Female | Adult | PayPal | Yes |
| Male | Adult | Cash | No |
| Female | Senior | Credit Card | Yes |
| Male | Teen | Cash | No |
| Female | Adult | PayPal | Yes |
| Female | Senior | Credit Card | Yes |
| Male | Teen | Cash | No |
| Male | Senior | PayPal | Yes |
| Female | Adult | Cash | No |

Demonstrate which attribute would you choose as the root node in a decision tree with multi-way splits using entropy impurity measure?

$$Entropy(t) = - \sum_{i=0}^{c-1} P_i(t) \log_2 P_i(t)$$

**① Purchase**

$$\text{Ent}(\text{Purchase}) = -\frac{5}{10}\log_2\left(\frac{5}{10}\right) - \frac{5}{10}\log_2\left(\frac{5}{10}\right)$$

$$= 0.5 + 0.5 = \boxed{1}$$

**Purchase**

| | |
|---|---|
| Yes | 5 |
| No | 5 |

**② Ent(Gender)**

Ent (M/Gender) $= -\frac{1}{5}\log_2\left(\frac{1}{5}\right) - \frac{4}{5}\log_2\frac{4}{5} =$

$= 0.464 + 0.257 = 0.721$

Ent(F/Gender) $= 0.721$

Ent(Gender) $= \frac{5}{10}(0.721) + \frac{5}{10}(0.721) = 0.721$

Gain info $= \text{Ent}(\text{Purchase}) - \text{Ent}(\text{Gender}) = 1 - 0.721 = 0.279$

**Gender**

| | Male | Female |
|---|---|---|
| Yes | 1 | 4 |
| No | 4 | 1 |

**③ Age Group**

Ent(Teen/Age group) $= -\frac{0}{3}\log_2\frac{0}{3} - \frac{3}{3}\log_2\frac{3}{3} = 0$

Ent (Adult/Age group) $= 1$

Ent (Senior/Age group) $= 0$

Ent (Age Group) $= \frac{3}{10}(0) + \frac{4}{10}(1) + \frac{3}{10}(0) = 0.4$

Gain info $= 1 - 0.4 = 0.6$

**Age group**

| | Teen | Adult | Senior |
|---|---|---|---|
| Yes | 0 | 2 | 3 |
| No | 3 | 2 | 0 |

**④ Payment Method**

Ent(CC/PM) $= -\frac{2}{3}\log_2\left(\frac{2}{3}\right) - \frac{1}{3}\log_2\frac{1}{3} = 0.385 + 0.533$
$= 0.9181$

Ent(Paypal/PM) $= 0$

Ent(Cash/PM) $= 0$

Ent (Payment Method) $= \frac{3}{10}(0.9181) + \frac{3}{10}(0) + \frac{4}{10}(0) = 0.2754$

Gain info $= 1 - 0.2754 = \boxed{0.7246}$

**Payment Method**

| | CC | Paypal | Cash |
|---|---|---|---|
| Yes | 2 | 3 | 0 |
| No | 1 | 0 | 4 |

So, As per the max Entropy and Gain info, Payment Method serve as the Root node for decision tree

Construction

**CLO # 2: Understand the nature of the data and apply data mining techniques to interpret the results.**

**Q2: Consider the following dataset,** [10+5+10=25 marks]

a) Construct a kdtree with y-x split order using median method. Store datapoints at leaf nodes only and attributes at non-leaf nodes. Keep splitting if there are more than 2 points in any region.

| Points | X | Y | Class |
|--------|----|----|-------|
| P1 | 8 | 17 | - |
| P2 | 6 | 20 | + |
| P3 | 4 | 21 | - |
| P4 | 9 | 13 | + |
| P5 | 12 | 12 | + |
| P6 | 4 | 26 | + |
| P7 | 3 | 22 | - |
| P8 | 7 | 23 | - |

(i)

Sorted on y

| 12 | 12 |
|----|----|
| 9 | 13 |
| 8 | 17 |
| 6 | 20 ] Median |
| 4 | 21 ] = 20.5 |
| 3 | 22 |
| 7 | 23 |
| 4 | 26 |

(ii)

Sorted on x

$M=8.5$ 
$\begin{cases} 6 & 20 \\ 8 & 17 \\ 9 & 13 \\ 12 & 12 \end{cases}$ $\Rightarrow$

sorted on x

$M=4$ 
$\begin{cases} 3 & 22 \\ 4 & 21 \\ 4 & 26 \\ 7 & 23 \end{cases}$

**KD Tree**



(iii)

Sorted on y

| 8 | 17 | 12 | 12 |
|---|----|----|----|
| 6 | 20 | 9 | 13 |

$M=18.5$

Sorted on y

| 4 | 21 | 7 | 23 |
|---|----|----|----|
| 3 | 22 | 4 | 26 |

$M=22$

| ≤ | | > | |
|---|---|---|---|
| 4 | 21 | 4 | 26 |
| 3 | 22 | | |

**b)** Find the nearest neighbor for points P9(4,8) and P10(11,10) using the above constructed tree. Explicitly mention which branch you moved at each level.

① $P9(4,8) \Rightarrow l_1 \rightarrow l_2 - \text{then} \rightarrow P1$ $\overset{+2.5}{\cancel{/}}$ so nearest

neighbor is P1. $\quad$ P2

② $P10(11,10) \Rightarrow l_1 - l_2 \text{ then} - P5$ $\overset{+2.5}{\cancel{/}}$ so nearest

neighbor is P5. $\quad$ P4

**c)** Classify the point P9(4,8) using 5 nearest neighbors using.
    i)      Majority voting approach
    ii)     Distance weighted majority weighting

(i)

~~P1~~ Majority Voting $\qquad\qquad (+5)$

(a) $P1(8,17), P9(4,8) = \sqrt{(8-4)^2 + (17-8)^2} = 9.84 - - \cdot 5$

(b) $P2(6,20), P9(4,8) = \sqrt{(6-4)^2 + (20-8)^2} = 12.16 + - \cdot 5$

(c) $P3(4,21), P9(4,8) = \sqrt{(4-4)^2 + (21-8)^2} = 13 - - \cdot 5$

(d) $P4(9,13), P9(4,8) = \sqrt{(9-4)^2 + (13-8)^2} = 7.07 + - \cdot 5 \quad +4$

(e) $P5(12,12), P9(4,8) = \sqrt{(12-4)^2 + (12-8)^2} = 8.94 + - \cdot 5$

(f) $P6(4,26), P9(4,8) = \sqrt{(4-4)^2 + (26-8)^2} = 18 \quad \cdot 5$

(g) $P7(3,22), P9(4,8) = \sqrt{(3-4)^2 + (22-8)^2} = 14.03 \quad \cdot 5$

(h) $P8(7,23), P9(4,8) = \sqrt{(7-4)^2 + (23-8)^2} = 15.29 \quad \cdot 5$

So 5 nearest neighbors are P1, P2, P3, P4, P5 and according to majority voting we can assign the class label to our test data P9(4,8) as " + ".

(ii) Distance weighted majority voting

distance of $P_9(4,8)$ with below datapoints and weighted distance

(a) $P_1(8,17)$ = 9.84 ; $\frac{1}{d^2} = \frac{1}{(9.84)^2} = 0.0103$ weighted distance

(b) $P_2(6,20)$ = 12.16 = $1/(12.16)^2 = 0.0067$

(c) $P_3(4,21)$ = 13 = $1/(13)^2 = 0.00591$

(d) $P_4(9,13)$ = 7.07 = $1/(7.07)^2 = 0.02000$

(e) $P_5(12,12)$ = 8.94 = $1/(8.94)^2 = 0.01251$

(f) $P_6(4,26)$ = 18 = $1/(18)^2 = 0.00308$

(g) $P_7(3,22)$ = 14.03 = $1/(14.03)^2 = 0.00508$

(h) $P_8(7,23)$ = 15.29 = $1/(15.29)^2 = 0.00427$

As 5 nearest neighbors are $P_1, P_2, P_3, P_4, P_5$
and $P_1, P_3$ belong to class "−" and $P_2, P_4, P_5$ to "+".

So,

d) for "−"

$$[(9.84 * 0.0103) + (13 * 0.00591)] = 0.101352 + 0.07683$$
$$= \boxed{0.178182}$$

For "+"

$$[(12.16 * 0.0067) + (7.07 * 0.02000) + (8.94 * 0.01251)]$$
$$= 0.081472 + 0.1414 + 0.1118394.$$
$$= \boxed{0.3347114}$$

As the value for class label "+" is larger than "−"

So we assign label "+" to our test datapoint $P_9(4,8)$.

**CLO # 2: Understand the nature of the data and apply data mining techniques to interpret the results.**

**Q3:** Consider the following ratings matrix with three users and six items. Ratings are on a 1-5 star scale. Compute the following from the data of this matrix.　　**[5+10=15 marks]**

|  | Item 1 | Item 2 | Item 3 | Item 4 | Item 5 | Item 6 |
|---|---|---|---|---|---|---|
| $U_1$ User 1 | 4 | 5 | $\circ$ | 5 | 1 | $\circ$ |
| $U_2$ User 2 | $\circ$ | 3 | 4 | 3 | 1 | $\circ$ |
| $U_3$ User 3 | 2 | $\circ$ | 1 | 3 | $\circ$ | 4 |

a) Treat missing values as 0. Compute the cosine similarity between each pair of users.

(i) user 1 and user 2

$\|U1\| = 16+25+0+25+1+0 = \sqrt{67} = 8.18$
$\|U2\| = 0+9+16+9+1+0 \sqrt{35} = 5.91$

$U1.U2 = 4*0 + 5*3 + 0*4 + 5*3 + 1*1$ to c
$= 31$

$cos(U1, U2) = \dfrac{U1.U2}{\|U1\|\|U2\|} = \dfrac{31}{8.18 * 5.91} = \boxed{0.6412}$

(ii) user 2 and user 3

$\|U2\| = \sqrt{35} = 5.91$
$\|U3\| = 4+0+1+9+0+16 \sqrt{30} = 5.47$

$U2.U3 = 0*2 + 3*0 + 4*1 + 3*3 + 1*0 + 0*4 = 13$

$cos(U2, U3) = \dfrac{U2.U3}{\|U2\|\|U3\|} = \dfrac{13}{5.9} = \boxed{0.402}$

(iii) user 1 and user 3

$U1.U3 = 4*2 + 5*0 + 0*1 + 5*3 + 1*0 + 0*4 = 23$　　$cos(U1, U3) = \dfrac{23}{8.18 (5.9)} = \boxed{0.514}$

b) Treat the above given data for user1 and user2 as vector1 and vector2 with missing values treated as "0", apply a scaling factor of 2 and translation factor of 5 to vector2 and demonstrate and prove that cosine similarity is either invariant to scaling and translation or not.

$vector1 = (4, 5, 0, 5, 1, 0)$
$vector2 = (0, 3, 4, 3, 1, 0)$

→ Scaling factor of 2 to vector 2

$vector2_s * 2 = (0, 6, 8, 6, 2, 0)$ $+2$

→ Translation factor of 5 to vector 2.

$vector2_t + 5 = (5, 8, 9, 8, 6, 5)$ $+2$

$(V_1, V_2)$

(i) $Cosine\left(\dfrac{V1 \cdot V2}{\|V1\| \|V2\|}\right) = \boxed{0.64 \$ 2}$ (Already driven in part(a))

$(V_1, V_{s2})$

(ii) $\|V2\| = 0+36+64+36+4+0 = \sqrt{140} = 11.8321$ $+2$

$V_i \cdot V_{s2} = 4*0 + 5*6 + 0*8 + 5*6 + 1*2 + 0*0 = 62$

$$Cos\left(\frac{V1 \cdot V_{ss}}{\|V1\| \|V_{ss}\|}\right) = \frac{62}{(67)(\|+83\cdot21)} = \boxed{0.6405}$$
$$(8.18)$$

(iii) $\|V_2\| = 25 + 64 + 81 + 64 + 36 + 25 = \sqrt{295} = 17.1755$

$V_1 \cdot V_{t2} = 4*5 + 5*8 + 0*9 + 5*8 + 1*6 + 0*5 = 106$

$$Cos\left(\frac{V1 \cdot V_{t2}}{\|V1\| \|V_{t2}\|}\right) = \frac{106}{(67)(\;-4\;)} = 0.00536$$
$$= \frac{106}{8.18 \times 17.1755} = \boxed{0.7544}$$

So, proved that scaling factor Cosine similarity is invariant to scaling and not invariant to translation.

**CLO # 3: Evaluate the performance and effectiveness of different data mining models using appropriate matrices and validation techniques.**

Q4: Consider a binary classification problem where a model is trained to predict whether a patient has a certain medical condition (positive class) or not (negative class). The model is evaluated using a test dataset containing 100 patient records. The following confusion matrix is obtained from the model's predictions: **[1*4 + 2+2+2 = 10 marks]**

| Actual | Predicted | |
|---|---|---|
| | Condition | No condition |
| Condition | 30  TP | 10  FN |
| No condition | 5  FP | 55  TN |

Write appropriate formula and calculate the following:

a) No. of individuals classifier fails to identify who have the medical condition?

$FN = 10$

b) No. of individuals classifier incorrectly identifies as having the medical condition?

$FP = 5$

c) No. of individuals the classifier correctly identifies who do not have the medical condition?

$TN = 55$

d) No. of individuals classifier correctly predicts that medical condition when they actually do have it.

$TP = 30$

e) Accuracy

$$TP+TN \mid TP+FP+FN+TN = \frac{30+55}{30+57+10+55} = 0.85$$

f) Specificity

$$TN \mid TN+FP = \frac{55}{55+5} = 0.91$$

g) Precision

$$TP \mid TP+FP = \frac{30}{5+30} = 0.85$$