

Homework no -1

Course ID : CSCE 5380

Name : Sai Harsha Reddy

ID : 11589323

Data Mining

1b) Data mining is the process of discovering patterns or relationships in large data sets, with the goals of extracting meaningful information and insights from the data.

Data Mining is a complex process that combines elements from several fields, including databases, statistics and machine learning. Data mining is more than just the simple application of technology from databases, statistics and machine learning. It involves a holistic approach to data analysis that integrates these fields and leverages their strengths to extract valuable insights from large and complex datasets.

2a) Let's assume that the data set contains the recordings of patients weight, height, age, temperature, blood pressure, glucose level and heart rate for every visit they come.

Classification :

Row : Records all patients

Column : Various measures like temperature, blood pressure, blood glucose and heart rate.

Clustering

Row: Patient ID

Column: Various measures like temperature, blood pressure, blood glucose and heart rate.

Association rule mining

Row: Patient ID

Column: various measure like temperature, blood pressure and heart rate.

Anomaly detection

Row: Patient ID

Column: various measures like temperature, blood pressure, blood glucose and heart rate.

4)

5) Continuous: Can be any value within the given range without any restriction

Discrete: Data can have only a specific value and it may be the case the number of values is infinite but each will be unique and non-overlapping

a) Continuous, quantitative, interval

For coverage use of the internet is continuous and its a multi quantity

and this if find with in a ~~certain~~ time
of interval.

b) Discrete, qualitative, ordinal

The mark will be considered a
specific one so its discrete then it
will be a quality of the performance
of the study.

c) Discrete as it will be distinct and
qualitative ordinal as needs to be
assigned in order - for the credit card
numbers.

d) Discrete, qualitative, ordinal

It implied a salary of all employees
in an organisation.

5)

a) The limitation of the approach to compute
similarity if we replace the null values
in sales commission by 0 is that we
can't distinguish between sales and
non sales person occupation.

If user wants to calculate the similarity
between all the users based on their
sales commission and replace null

but if we have a problem null stands for not applicable and hence means that the user does not need to fill in that value of column for that table. If we replace null with 0 will be unable to distinguish sales ~~and~~ non sales peoples as their sales peoples with 0 sales commission. Hence we now have any method to distinguish between sales and non sales person occupation. This is the short coming of replacing null ~~as~~ 0 in the sales commission attribute.

6) a) Distance b/w given objects (22, 1, 42, 10)
 $\times (20, 0, 36, 8)$

$$\begin{aligned}
 &= \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2 + (p_2 - p_1)^2} \\
 &= \sqrt{(20-22)^2 + (0-1)^2 + (36-42)^2 + (8-10)^2} \\
 &= \sqrt{4+1+36+4} \\
 &= \sqrt{45} \\
 &= 6.7082
 \end{aligned}$$

7) a) let consider the following dataset we want to decide whether the Patients have lung cancer or not.

where

class P = s \in lung cancer = "Yes"

class N = s lung cancer = "No"

Step 1: Computer the gini index of the overall collection of training set.

formula $Gini(D) = 1 - \sum_{i=1}^m p_i^2$

$$\begin{aligned}
 gini(D) &= 1 - (S/10)^2 - (S/10)^2 \\
 &= 0.5
 \end{aligned}$$

Step 2:

Select the best split among tobacco smoking
Random exposure and chronic cough) For
the root dist all the split you consider
together their corresponding values.

formula

$$gini_A(D) = \frac{|D_1|}{|D|}$$

$$= gini(D_1) + \frac{|D_2|}{|D|}$$

Let's now consider Tobacco smoking: {yes, No}

It is a Binary attribute

gini

Tobacco Smoking ∈ {yes} (D) = $(\frac{4}{10})$

$$(1 - (\frac{4}{10}))^2 - (\frac{1}{10})^2 + (\frac{3}{10})(1 - (\frac{1}{10}))^2 (\frac{4}{10})$$

$$= 0.32 \Rightarrow gini \text{ Tobacco Smoking}$$

Let Now consider Random exposure {yes, no}

It is binary attribute

gini random exposure $\in \{ \text{yes} \}$ (o)

$$(1 - (2/2)^2 - 0) + (8/10) \left(1 - \frac{3}{8}\right)^2 - \left(\frac{5}{8}\right)^2$$

$= 0.375 \Rightarrow$ gini random exposure $\in \{ \text{No} \}$ (o)

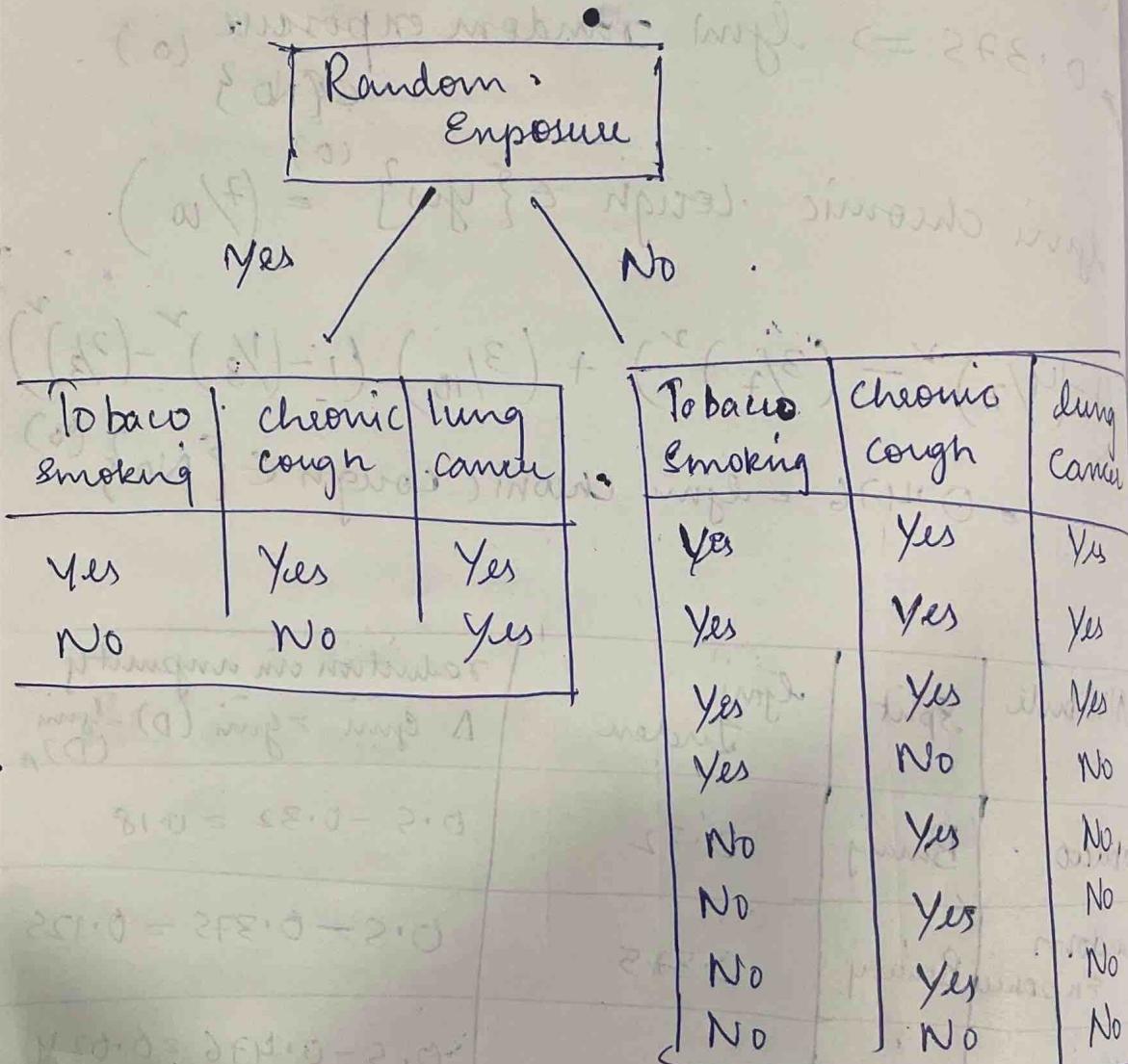
gini chronic cough $\in \{ \text{yes} \}$ (o) $= (7/10)$

$$(1 - (4/7)^2 - (3/7)^2) + (3/10) \left(1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2\right)$$
 $= 0.476 =$ gini chronic cough $\in \{ \text{No} \}$ (o)

Attribute	Split	gini Index	reduction in impurity $\Delta \text{gini} = \text{gini}(D) - \text{gini}(D_s)$
Tobacco	Binary	0.32	$0.5 - 0.32 = 0.18$
Random Exposure	Binary	0.375	$0.5 - 0.375 = 0.125$
Chronic cough	Binary	0.476	$0.5 - 0.476 = 0.024$

The attribute that minimizes the reduction in impurity or equivalently has the minimum gini index gini split(n) is selected as the splitting attribute

Tobacco smoking is selected with minimum index under sand highest reduction in impurity.



$$\text{Step 1: } \text{min}(n) = 1 - \left(\frac{3}{8}\right)^2 - \left(\frac{5}{8}\right)^2$$

$$= 0.468$$

Step 2: Let now ~~Tobacco smoking~~
= {Yes, No}

It is a binary attribute

gini Tobacco smoking $\in \{ \text{Yes} \}^{(0)}$

$$= \left(\frac{4}{8} \right) \left(1 - \left(\frac{3}{4} \right)^2 - \left(\frac{1}{4} \right)^2 \right) + \left(\frac{4}{8} \right) \left(1 - 0 - \left(\frac{4}{4} \right)^2 \right)$$

$$= 0.1875 + 0 \quad \begin{array}{l} \text{add up} \\ \text{answers} \end{array}$$

$$= 0.187$$

gini tobacco smoking $\in \{ \text{No} \}^{(0)}$

let now consider chronic cough $\{ \text{yes}, \text{no} \}$

It is a binary attribute

gini chronic cough $\in \{ \text{yes} \}^{(0)}$

$$= \left(\frac{6}{8} \right) \left(1 - \left(\frac{3}{6} \right)^2 - \left(\frac{3}{6} \right)^2 \right) + \left(\frac{2}{8} \right) \left(1 - 0 - \left(\frac{2}{2} \right)^2 \right)$$

$$= 0.375 \quad \Rightarrow \text{gini chronic cough } \in \{ \text{no} \}^{(0)}$$

Attribute	split	gini index	reduction in entropy
Tobacco smoking	binary	0.187	0.281
chronic cough	binary	0.375	0.093

Random exposure

lung = yes

Tobacco smoking

Yes

No.

Chronic cough	lung cancer	Chronic cough	lung cancer
Yes	Yes	Yes	No
Yes	Yes	Yes	No
Yes	Yes	Yes	No
No	No	No	No

8)
a)

We can modify the basic decision tree algorithm by integrating the count of each type into the calculation of the attribute selection measure taking the count into consideration to determine the most common class among the types.