

Sai Mahesh Muditavat

11527206

Data Mining

Section - 002

(Q) Data mining

- Data mining is process of automatically discovering useful information in large data repositories. It provides tools to discover knowledge from data.
- In nutshell data mining turns a large collection of data into knowledge.
- So it is ~~more than simple transformation technology~~
- (b) Data mining is definitely ~~a~~ ^{more than} simple transformation of technology which evolved into Machine learning, databases, Statistics because huge data when started to accumulate to better manage those it first evolved into data management systems, to retrieve the require data we started to use SQL. As Centuries decades passed data started to accumulate into advance database & mining Complex data we later used statistics and Machine learning to get right kind of data to solve a problem.

② Ambulance Medical care dataset

Features: gender, age, duration of visit, diagnosis, demography, symptoms, medication.

Task: Classification

Row: Patient (or) Patient name
ID

Column: Symptoms, medication, visit, age

Target: To find whether patient has disease
(a) yes
(b) no

Clustering

Test Row: Patient demography (Patient visits, marital status, age)
Column: Symptoms, age, demography etc

Target: Clustering can be done by having age group of people in with no symptoms and other set of age, with symptoms and predict whether this patient has chance of getting disease.

Patients also be clustered based on demography
Conditions and predict whether there
is a possibility to predict this
Set of people of certain demography have
chance of getting disease (not interested)

Associate Rule mining

Row: Patient visit (Patient id)

Column:

Symptoms, age, diagnosis of medical condition.

Target: We identify if Symptoms and medical conditions coexist

We can also see demograph age and medical conditions Correlation between them

Anomaly detection

labeled 16 Mar 2013

Row: patient visit

Column: Symptoms, demography, age, medication

Target: Young age people with sudden & illnes Certain demography showing certain common disease symptoms (e.g.: diabetes)

④ A Average number of hours were spent on Internet in week

Continuous, quantitative, ratio

Continuous because it is specific has range of values and real numbers

Quantitative

Ratio: As it average hours it can also mean of hours of usage across the week

(b) GPA of Student

→ Continuous, Qualitative, Ordinal
GPA Values are 4.1, 4.2, 7.0, 8.3,

So they are continuous
ordinal because values are in certain order

(c)

Credit Card number

Discrete, Qualitative, Nominal

Discrete due its distinctness

Qualitative Nominal: There are distinguishable

and formed with certain permutation

(d) Salary above median Salary of all employees

Continuous, Quantitative, Interval

Continuous, Salary figures are always in growing order

Interval: Salary are within this interval range
Hence it is also Quantitative

5① Null means not applicable, if we replace with zero will be unable to distinguish non sales people, as there is possibility of sales person with 0. This is a disadvantage of replacing null with 0 in particular in sales commission attribute → Mary will be more similar to Bob and Lisa

$$P = (22, 1, 42, 10)$$

$$Q = (20, 0, 36, 8)$$

Distance $d(P, Q) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$

$$\begin{aligned} d &= \sqrt{(22-20)^2 + (1-0)^2 + (42-36)^2 + (10-8)^2} \\ &= \sqrt{4 + 1 + 36 + 4} \\ &= 6.708 \end{aligned}$$

7A

~~Entropy related~~

$$P+ = 5 \xrightarrow{0.5 \text{ good}} 0.5 \text{ (good)}$$

$$P- = 5 \xrightarrow{0.5 \text{ - normal}} 0.5 \text{ (normal)}$$

$$\text{Entropy} = -0.5 \log_2 (0.5) - (0.5)^{\log_2 (0.5)} \\ = 1 \text{ spol}^2 - (2^1)^{\text{spol}^2} \\ (\text{NH}_3)$$

Care
Tobacco Smoking Yes - Lung Cancer Yes
No - Lung Cancer No

14

Tobacco Smoking No - Lung Cancer No
1 (20.)

EC_{4,1})

$$= -\left(\frac{4}{5} \log_2 \frac{4}{5}\right) - \left(\frac{1}{5} \log_2 \frac{1}{5}\right) \\ = 0.72$$

EC_{1,4}) \rightarrow (Tobacco Smoking Yes - Lung Cancer No
Tobacco Smoking No - Lung Cancer No

$$= 0.72$$

Information gain ($\text{Entropy}_{\text{parent}} - \sum \text{Entropy}_{\text{child}}$)
 $1 - 0.72$
 $\underline{= 0.278}$

Radon Exposure

$$\begin{array}{ll}
 \text{Radon Yes} & \text{lung cancer yes} \\
 \text{Radon No} & \text{lung cancer - Yes} \\
 \end{array}$$

$(2, 3) - (2-0)$ ≈ 0.9
 $= \frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right)$
 $= 0.52 - (-0.44)$

0.96

Radon Yes

Radon No

L.C - Yes

L.C - No

$(0, 5)$

$\equiv 0$

Information gain

$$0.76 \pm 0.23$$

Chronic Cough $\xrightarrow{\text{Yes}}$ Lung Cancer Yes - lead to predict
 $\xrightarrow{\text{No}}$ Lung Cancer No (S.P.)

Chronic Cough
No

Lung Cancer Yes
(S.P.)

$$E(Y_0) = (1, 1)$$

$$- \frac{4}{15} \log_2 \frac{4}{15} - \frac{1}{15} \log_2 \frac{1}{5}$$

$$\approx 0.72$$

Yes, Yes
Cough No, No

L.C.

$(3, 2)$

$$- \left(\frac{3}{15} \log_2 \frac{3}{15} \right) - \frac{2}{15} \log_2 \frac{2}{15} = 0.96$$

Information gain:

$$1 - 0.96 = \underline{0.04}$$

informational content and predictive power

Weight loss - yes L.C: Yes
Weight loss NO L.C: Yes
(3, 2)

$$- \frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = (0.97) \text{ J}$$

$$= 0.97$$

Weight loss Yes

L.C - NO

Weight loss: NO

L.C: No

(2, 3)

OK, OK

OH, OH

$$- \frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = (0.97) \text{ J}$$

$$= 0.97$$

$$= (2 \text{ J}, 2 \text{ J})$$

Information gain

$$1 - 0.971 = 0.029$$

Tobacco smoking has highest Information gain

Q

Tobacco smoking Yes

Radon exposure Yes

Radon exposure No

(1, 3)

$$\left(-\frac{1}{5} \log_2 \frac{1}{5} \right) - \left(\frac{3}{5} \log_2 \frac{3}{5} \right)$$

$$(-0.46) - (-0.44) \rightarrow 0$$

R.Exposure (Yes, No) → 0

Entropy total:

$$0.649$$

Info gain: 0.072

Chronic Cough

(Yes, Yes)

No, Yes

(4, 0)

(Yes, No)
No, No

(0, 1)

1
0

→ Info gain: 0.72

Weight loss $T \cdot S \rightarrow$ Yes

lost promotion

Yes Yes

No, Yes

(1, 0) = 0

lost savings

on savings?

(E, 1)

T.S: Yes

(Yes, No)
No, No

(2, 2) = 0.91%

(E, E) - (E, N)

(N, E) - (N, N)

Information gain: 0.17%

We choose Chronic Cough for Tobacco

Smoking = Yes

(Yes, Yes)
No, No

(No, No)

lost work

(Yes, Yes)

lost ch

(0, 0)

Tobacco Smoking

= No

2201 Bf p64
01.2.7

Radon: Yes, L.C = Yes \rightarrow 0.5
 \searrow No, L.C = Yes \rightarrow 0.4
 $(1, 0)$ $\quad \quad \quad 0 = (0, 1)$
 \vdots
 $= 0$

Radon: Yes, L.C = Yes \rightarrow 0.5

Radon No : L.C = No \rightarrow 0.5

$(0, 4)$

Information gain: $\rightarrow 0.72$

Chronic Cough, Lung Cancer: T-S = No

(Yes, Yes)
No, Yes

$(0, 1) = 0$

$(3, 1) = 0.40$

Information gain: 0.32

Weight loss

T.S = No

Yes

L.C = Yes

Weight loss

No

Yes

$$(1, 0) = 0$$

$$(0, 1)$$

Weight loss

$$(2, 2)$$

$$= 0.91$$

$$\text{Entropy} = 0.55$$

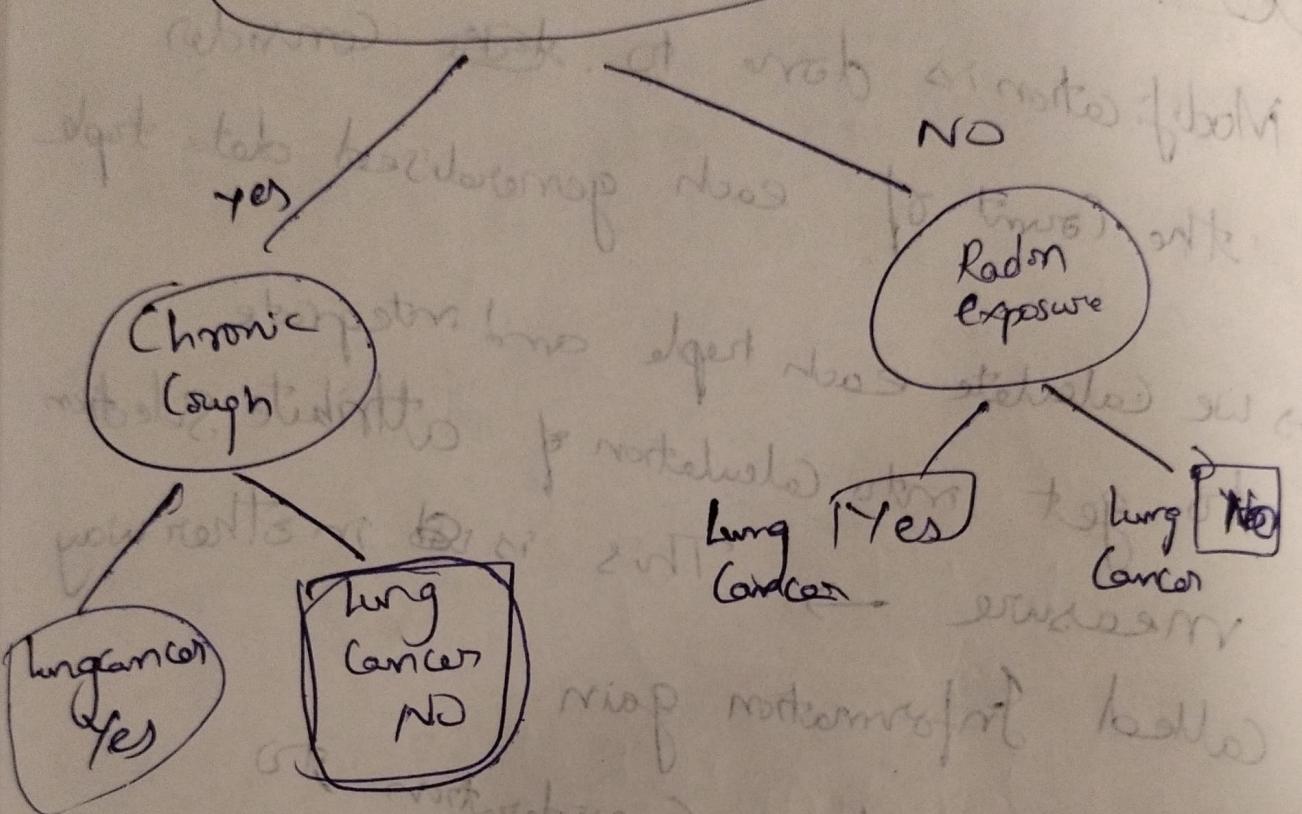
Information: 0.170

Here best Split for Tobacco Smoking = Radon exposure

If Radon exposure = Yes, L.C = Yes
i.e., " " = no L.C = no

$$0.170 : (1, 0) = (1, \epsilon)$$

Tobacco Smoking



Modification is done to consider
the count of each generalised data tuple

→ We calculate each tuple and integrate
to get into calculation of attribute selection
measure → This is in other way

called Information gain

→ We also count in consideration
determine common class among the tuples.