Name: M. ASLIN REDDY

Roll No: 211121

Data Warehouse & Data Mining

Assignment - 1

Implementation of Decision Tree Classifier

Data Set:

| Age | Income | Student | Credit-Rating | Buys Computer |
|-----|--------|---------|---------------|---------------|
| ≤ 30 | high | no | fair | no |
| ≤ 30 | high | no | excellent | no |
| 31-40 | high | no | fair | yes |
| >40 | medium | no | fair | yes |
| >40 | low | yes | fair | yes |
| >40 | low | yes | excellent | no |
| 31-40 | low | yes | excellent | yes |
| ≤ 30 | medium | no | fair | no |
| ≤ 30 | low medium | yes | fair | yes |
| >40 | medium | yes | fair | yes |
| ≤ 30 | medium | yes | excellent | yes |
| 31-40 | medium | no | excellent | yes |
| 31-40 | high | yes | fair | yes |
| >40 | medium | no | excellent | no |

First we need to calculate Information gain for every attribute other than target attribute which is Buys-Computer

From the attribute Which has highest ~~for info~~ information gain we can partition the dataset from that attribute and continue the Process.

Partition is stopped when the following Conditions met:

i. All Samples belong to the same class

ii. No Samples left

iii. No Remaining Attributes for further Partioning

~~Info~~ Entropy of Target $= -\sum\limits^{m} P_i \log_2(P_i)$

Class P: buys_Computer = "yes" = 9

Class N: buys_Computer = "no" - 5

$$E(Buys\_Computer) = -\frac{9}{14}\log_2(9/14) - \frac{5}{14}\log_2(5/14)$$

$$E(Target) = 0.940$$

$$E(Buys\_computer, Age) = P(\leq30)\cdot E(target, \leq30) +$$
$$P(31-40)\cdot E(target, 31-40) +$$
$$P(>40)\cdot E(target, >40)$$

| Age | | Buys-Computer yes | No |
|---|---|---|---|
| | ≤30 | 2 | 3 |
| | 31-40 | 4 | 0 |
| | >40 | 3 | 2 |

$$= P(\leq 30) \cdot E(2,3) + P(31-40) \cdot E(4,0) + P(>40) \cdot E(3,2)$$

$$= \frac{5}{14}\left[-\frac{2}{5} \cdot \log_2(2/5) - \frac{3}{5}\log_2(3/5)\right] + \frac{4}{14}\left[-\frac{4}{4}\log(4/4) - \frac{0}{4}\log(0/4)\right]$$

$$+ \frac{5}{14}\left[-\frac{3}{5}\log_2(3/5) - \frac{2}{5}\log(2/5)\right]$$

$$= \frac{5}{14}(0.970) + \frac{4}{14} \cdot 0 + \frac{5}{14}(0.970)$$

E(Buys Comp, age)

$$= 0.6935$$

Information Gain (Target) = E(Target) - E(Target, age)
Age

$$= 0.940 - 0.6935$$

Gain (Target) = 0.246
Age

E(Target, Income) = $\quad$ P(high) · E(Target, high) +

$$P(medium) \cdot E(Target, medium) +$$

$$P(Low) \cdot E(Target, Low)$$

|  |  | Buys_computer | |
|---|---|---|---|
|  |  | yes | No |
| Income | High | 2 | 2 |
|  | medium | 4 | 2 |
|  | Low | 3 | 1 |

$$= \frac{4}{14} \cdot E(2,2) + \frac{6}{14} E(4,2) + \frac{4}{14} E(3,1)$$

$$= \frac{4}{14}\left[-\frac{2}{4}\log_2(2/4) - \frac{2}{4}\log_2(2/4)\right] +$$

$$\frac{6}{14}\left[-\frac{4}{6}\log_2(4/6) - \frac{4}{6}\log_2(4/6)\right] +$$

$$\frac{4}{14}\left[-\frac{3}{4}\log_2(3/4) - \frac{1}{4}\log_2(1/4)\right]$$

$$= \frac{4}{14}(1) + \frac{6}{14}(0.918) + \frac{4}{14}(0.811)$$

$$E\!\left(\begin{array}{c}\text{target}\\ \text{Income}\end{array}\right) = 0.910$$

$$\text{Gain}\!\left(\begin{array}{c}\text{Target}\\ \text{Income}\end{array}\right) = E(\text{Target}) - E(\text{Target, Income})$$

$$= 0.940 - 0.910$$

$$\text{Gain}\,(\text{Target}) = 0.03$$
$$\text{income}$$

$$E(\text{target, Student}) = P(\text{yes}) \cdot E(\text{target, yes}) + P(\text{No}) E(\text{target, no})$$

| | Buys − Computer | |
|---|---|---|
| | yes | No |
| Student. yes | 26 | 01 |
| no | 03 | 14 |

$$= \frac{7}{14} E(6,1) + \frac{7}{14} E(3,4)$$

$$= \frac{7}{14}\left[-\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 (\frac{1}{7})\right] +$$
$$\frac{7}{14}\left[-\frac{3}{7} \log_2 (\frac{3}{7}) - \frac{4}{7} \log_2 (\frac{4}{7})\right]$$

$$= \frac{7}{14}\,(0.591) + \frac{7}{14}(0.985)$$

$$E(\text{target, Student}) = 0.788$$

$$\text{Gain}\,(\text{Target}) = E(\text{Target}) - E(\text{Target, Student})$$
$$\text{Student}$$

$$= 0.940 - 0.788$$

$$= 0.152$$

$$E(\text{target, Credit Rating}) = P(\text{fair}) \cdot E(\text{target, fair}) +$$
$$P(\text{excellent}) \cdot E(\text{target, Excellent})$$

| Credit Rating | | Buys_Computer | |
|---|---|---|---|
| | | yes | No |
| | fair | 3 | 3 |
| | excellent | 6 | 2 |

$$= \frac{6}{14} \cdot E(3,3) + \frac{8}{14} E(6,2)$$

$$= \frac{6}{14}(1) + \frac{8}{14}(0.811)$$

$$E(\text{target, credit}) = 0.892$$

$$\text{Gain (Target)} = E(\text{Target}) - E(\text{target, Credit})$$
$$\text{Credit}$$

$$= 0.940 - 0.892$$

$$= 0.048$$

By Comparing Information Gain from 4 attributes

| Age | Income | Student | Credit_Rating |
|---|---|---|---|
| 0.246 | 0.03 | 0.152 | 0.048 |

Age has highest Information Gain.

Root node: **Age**

Branches: **≤30**, **31–40**, **>40**

### ≤30

| Income | Student | Credit-rating | Class |
|---|---|---|---|
| high | no | fair | no |
| high | no | excellent | no |
| medium | no | fair | no |
| low | yes | fair | yes |
| medium | yes | excellent | yes |

### 31–40

| Income | Student | Credit-Rating | Class |
|---|---|---|---|
| high | no | fair | yes |
| low | yes | excellent | yes |
| medium | no | excellent | yes |
| high | yes | fair | yes |

### >40

| Income | Student | Credit-Rating | Class |
|---|---|---|---|
| medium | no | fair | yes |
| low | yes | fair | yes |
| low | yes | excellent | no |
| medium | yes | fair | yes |
| medium | no | excellent | no |

→ In the samples of 31-40, all belong to same class of the target. So, 31-40 group buys computer.

Now consider ≤30 group/node.

Entropy of ≤30

$$E(target)_{≤30} = E(2,3)$$    2-yes's
                                3- No's

$$= -\frac{2}{5} Log_2(2/5) - \frac{3}{5} Log_2(3/5)$$

$$E(target)_{≤30} = 0.970$$

$$E(target_{≤30}, income) = P(high)E(y, high) + P(medium)E(y, med) + P(Low)E(y, Low)$$

y

| | yes | no |
|---|---|---|
| high | 0 | 2 |
| medium | 1 | 1 |
| Low | 1 | 0 |

$$= \frac{2}{5} \frac{2}{5} E(0,2) + \frac{2}{5} E(1,1)$$
$$+ \frac{1}{5} E(1,0)$$

$$= \frac{2}{5}(0) + \frac{2}{5} + \frac{1}{5}(0)$$

$$= 0.4$$

Gain (target, ≤30) = E(target, ≤30) - E(target, ≤30, income)

income

$$= \boxed{0.970 - 0.4}$$

$$= 0.570$$

$$E(target, ≤30, Student) = P(no) \cdot E(y, no) + P(yes) \cdot E(y, no)$$

y

$$= \frac{3}{5} E(3,3) + \frac{2}{5} E(2,0)$$

$$= \frac{3}{5}(0) + \frac{2}{5}(0)$$

$$= 0$$

$$\text{Gain (target, } \leq 30) = 0.970 - 0$$
Student

$$= 0.970$$

$$E(\text{target, } \leq 30, \text{ Credit. Rating}) = P(\text{fair}) \cdot E(y, \text{fair}) + P(\text{excellent})$$
$$\underbrace{\quad}_{y} \qquad \qquad \qquad \cdot E(y, \text{excell}_{ent})$$

$$= \tfrac{3}{5} E(1,2) + \tfrac{2}{5} (0 \, E(1,1))$$

$$= \boxed{0.6591}$$

$$\text{Gain} \quad (\text{target, } \leq 30) = 0.970 - 0.6591$$
Credit Rating

$$= 0.310$$

By Comparing Information Gain from 3 Attributes

in $\leq 30$ group

| Income | Student | Credit - Rating |
|--------|---------|-----------------|
| 0.570  | ~~0.6591~~ 0.970 | 0.310 |

Student has highest Information Gain.

$\leq 30$

Student ?

No         yes

| Income | Credit. | Class |
|--------|---------|-------|
| high | fair | no |
| high | excellent | no |
| ~~high~~ · med | fair | no |

| Income | Credit. | class |
|--------|---------|-------|
| low | fair | yes |
| med. | excellent | yes |
|  |  | ~~yes~~ |

As samples of both yes & no have ~~classes~~ Pure Class

~~#~~ ~~their~~ ~~own~~ , further Partioning is not needed.

$\leq 30$

Student?

no      yes

no      yes

Now Consider >40 group

$$E \text{ (target, } >40) = E(3,2)$$
$$\underbrace{\qquad}_{k}$$

3 - yes's
2 - no's

$$= 0.970$$

$$E(\text{k, income}) = P(\text{medium}) \cdot E(k, \text{med}) + P(\text{low}) \cdot E(k, \text{low})$$

$$= \frac{3}{5} E(2,1) + \frac{2}{5} E(1,1)$$

$$= \cancel{0.950} \ 0.951$$

$$E(k, \text{Student}) = P(\text{yes}) \cdot E(k, \text{yes}) + P(\text{no}) \cdot E(k, \text{no})$$

$$= \frac{3}{5} \cdot E(2,1) + \frac{2}{5} E(1,1)$$

$$= 0.951$$

$$E(k, \text{Credit Rating}) = P(\text{fair}) \ E(k, \text{fair}) + P(\text{excellent}) \cdot E(k, \text{excellent})$$
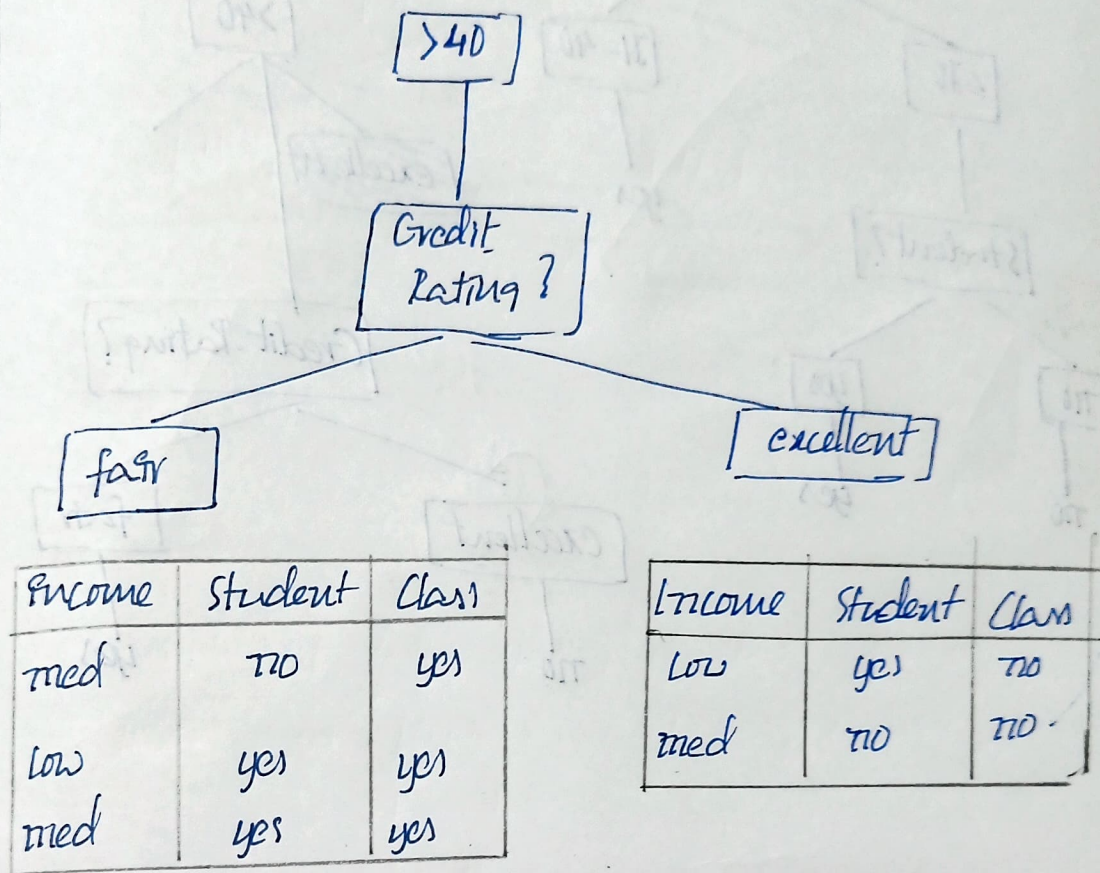
$$= \frac{3}{5} E(3,0) + \frac{2}{5} E(0,2)$$

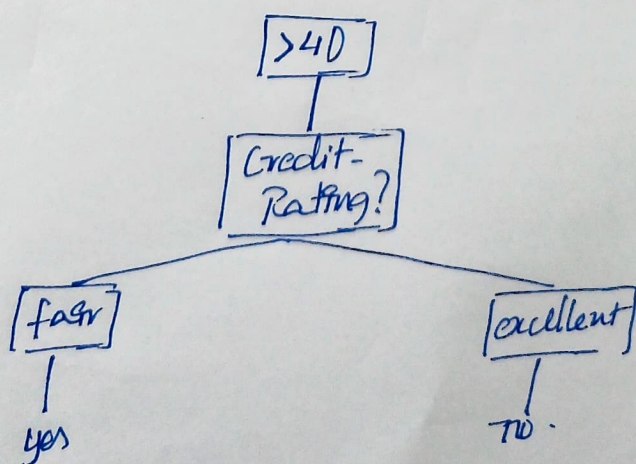$$= 0$$

$$\text{Gain (target, } >40) = 0.970 - 0.951 = 0.02$$
income

$$\text{Gain (target, } >40) = 0.970 - 0.951 = 0.02$$
Student

$$\text{Gain (target, } >40) = 0.970 - 0 = 0.970.$$
credit

Credit_Rating has highest information gain in >40 group.

>40

Credit Rating ?

fair

excellent

| Income | Student | Class |
|--------|---------|-------|
| med | no | yes |
| low | yes | yes |
| med | yes | yes |

| Income | Student | Class |
|--------|---------|-------|
| low | yes | no |
| med | no | no |

As samples of both fair & excellent have Pure Classes, further Partitioning is not needed and assign the labels respectively.

>40

Credit-Rating?

fair

yes

excellent

no

```
                    ( Root )
                       |
                    ( Age ? )
          ≤30  /        |  31-40     \  >40
             /          |              \
      ( Student? )    ( yes )      ( Credit Rating? )
      no /    \ yes              fair /      \ excellent
        /      \                    /          \
     ( No )  ( yes )            ( yes )        ( No )
```