# University of Victoria

# CSC - 503

# Report (Assignment - 1)

**Submitted by :**

Sannath Reddy Vemula - V00949217

Shaima Patel - V00949940

**Submitted To : Alex Thomo**

thomo@cs.uvic.ca

**1) (4 points) Construct the root and the first level of a decision tree for the titanic dataset. Use the ID3 algorithm. Show the details of your construction (entropies calculated for each step). You can use a spreadsheet or SQL database to compute the counts.**

**Then, check your solution with Weka and submit a text file of your classifier output window.**

Assignment – 1   Data Mining

શ્રી સ્વામિનારાયણ ડિવાઇન મિશન

Dt.: / /    Pg.no.:

| Age | Sex | Pclass | Survived |
|---|---|---|---|
| A\|C | M/F | 1/2/3/ crew | Y/N |

Total rows = 2201
Yes = 711
No = 1490

Age :   Adult  N→ 1438    = 2092
              Y→ 654
        Child  N→ 52      = 109
              Y→ 57

Sex :   Male   N→ 1364    = 1731
              Y→ 367
        Female N→ 126     = 470
              Y→ 344

Pclass :  1st  N→ 122     = 325
              Y→ 203
         2nd   N→ 167     = 285
              Y→ 118
         3rd   N→ 528     = 706
              Y→ 178
         Crew  N→ 673     = 885
              Y→ 212

$\rightarrow$ Entropy $\left(\text{Survived}/P=1st\right)=$

$-P\left(\dfrac{No}{\underset{1st}{Total}}\right) * \log_2 P_3\left(\dfrac{No}{\underset{1st}{Total}}\right) - P\left(\dfrac{Yes}{\underset{1st}{Total}}\right) *$

$\log_2 P\left(\dfrac{Yes}{Total\ 1st}\right)$

$\boxed{\underline{P\text{-}class}}$

1st $= -\left(\dfrac{203}{325}\right)\log_2\left(\dfrac{203}{325}\right) - \left(\dfrac{122}{325}\right)\log_2\left(\dfrac{122}{325}\right)$

$= +0.4240 + 0.5306$

$= \boxed{0.9546}$

2nd $= -\left(\dfrac{118}{285}\right)\log_2\left(\dfrac{118}{285}\right) - \left(\dfrac{167}{285}\right)\log_2\left(\dfrac{167}{285}\right)$

$= +0.4518 + 0.5267$

$= \boxed{\cdot 0.9785}$

3rd $= -\left(\dfrac{528}{706}\right)\log_2\left(\dfrac{528}{706}\right) - \left(\dfrac{178}{706}\right)\log_2\left(\dfrac{178}{706}\right)$

$= +0.3134 + 0.5011$

$= \boxed{\cdot 0.8145}$

New $= -\left(\dfrac{673}{885}\right)\log_2\left(\dfrac{673}{885}\right) - \left(\dfrac{212}{885}\right)\log_2\left(\dfrac{212}{885}\right)$

$$= + 0.3004 + 0.4938$$
$$= \boxed{0.7942}$$

selected $\downarrow$     $\boxed{Sex}$

$$Male = -\left(\frac{1364}{1731}\right) \log_2\left(\frac{1364}{1731}\right) - \left(\frac{367}{1731}\right) \log_2\left(\frac{367}{1731}\right)$$
$$= + 0.2708 + 0.4744$$
$$= \boxed{0.7452}$$

$$Female = -\left(\frac{126}{470}\right) \log_2\left(\frac{126}{470}\right) - \left(\frac{344}{470}\right) \log_2\left(\frac{344}{470}\right)$$
$$= 0.5091 + 0.3295$$
$$= \boxed{0.8386}$$

$\boxed{Age}$

$$child = -\left(\frac{52}{109}\right) \log_2\left(\frac{52}{109}\right) - \left(\frac{57}{109}\right) \log_2\left(\frac{57}{109}\right)$$
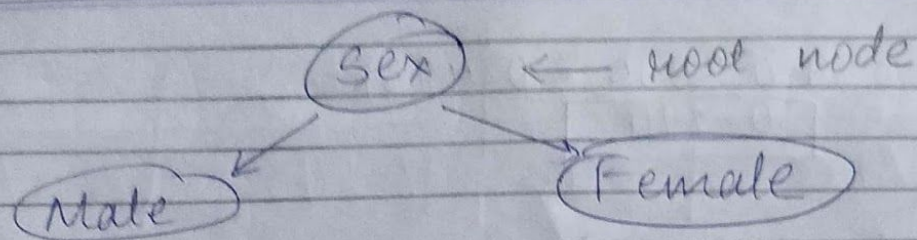$$= \boxed{0.9983}$$

$$Adult = -\left(\frac{1438}{2092}\right) \log_2\left(\frac{1438}{2092}\right) - \left(\frac{654}{2092}\right) \log_2\left(\frac{654}{2092}\right)$$
$$= 0.3713 + 0.5244$$
$$= \boxed{0.8967}$$

Sex ← root node

Male                    Female

Male                        Female

[Pclass]

1st $\xrightarrow{Y}$ 62      = 180      1st $\xrightarrow{Y}$ 141   = 145
    $\xrightarrow{N}$ 118                   $\xrightarrow{N}$ 4

2nd $\xrightarrow{Y}$ 25   = 179      2nd $\xrightarrow{Y}$ 93   = 106
    $\xrightarrow{N}$ 154                  $\xrightarrow{N}$ 13

3rd $\xrightarrow{Y}$ 88   = 510      3rd $\xrightarrow{Y}$ 90   = 196
    $\xrightarrow{N}$ 422                  $\xrightarrow{N}$ 106

crew $\xrightarrow{Y}$ 192 = 862     crew $\xrightarrow{Y}$ 20   = 23
     $\xrightarrow{N}$ 670                 $\xrightarrow{N}$ 3

[Age]

Male $\xrightarrow{Y}$ 338 = 1667    Male $\xrightarrow{Y}$ 316   = 425
     $\xrightarrow{N}$ 1329               $\xrightarrow{N}$ 109

Female $\xrightarrow{Y}$ 29 = 64     Female $\xrightarrow{Y}$ 28   = 45
       $\xrightarrow{N}$ 35                 $\xrightarrow{N}$ 17

(Male)

## Pclass

1st $= -\left(\dfrac{62}{180}\right) \log_2 \left(\dfrac{62}{180}\right) - \left(\dfrac{118}{180}\right) \log_2 \left(\dfrac{118}{180}\right)$

$= +0.5296 + 0.3993$

Selected $= \boxed{0.9289}$

(2nd) $= -\left(\dfrac{25}{179}\right) \log_2 \left(\dfrac{25}{179}\right) - \left(\dfrac{154}{179}\right) \log_2 \left(\dfrac{154}{179}\right)$

$= 0.3966 + 0.1867$

$= \boxed{0.5833}$

3rd $= -\left(\dfrac{88}{510}\right) \log_2 \left(\dfrac{88}{510}\right) - \left(\dfrac{422}{510}\right) \log_2 \left(\dfrac{422}{510}\right)$

$= 0.4373 + 0.2261$

$= \boxed{0.6634}$

Crew $= -\left(\dfrac{192}{862}\right) \log_2 \left(\dfrac{192}{862}\right) - \left(\dfrac{670}{862}\right) \log_2 \left(\dfrac{670}{862}\right)$

$= 0.4825 + 0.2825$

$= \boxed{0.7650}$

## Age

Adult $= -\left(\dfrac{338}{1667}\right) \log_2 \left(\dfrac{338}{1667}\right) - \left(\dfrac{1329}{1667}\right) \log_2 \left(\dfrac{1329}{1667}\right)$

$= 0.4667 + 0.2606$

$= \boxed{0.7273}$

Follow the river and will

child $= -\left(\dfrac{29}{64}\right) \log_2 \left(\dfrac{29}{64}\right) - \left(\dfrac{35}{64}\right) \log_2 \left(\dfrac{35}{64}\right)$

$= 0.5174 + 0.4761$

$= \boxed{0.9935}$

(Female)

selected

$\boxed{Pclass}$

①st $= -\left(\dfrac{141}{145}\right) \log_2 \left(\dfrac{141}{145}\right) - \left(\dfrac{4}{145}\right) \log_2 \left(\dfrac{4}{145}\right)$

$= 0.039 + 0.142$

$= \boxed{0.181}$

2nd $= -\left(\dfrac{93}{106}\right) \log_2 \left(\dfrac{93}{106}\right) - \left(\dfrac{13}{106}\right) \log_2 \left(\dfrac{13}{106}\right)$
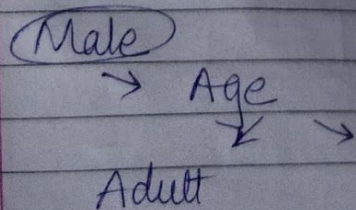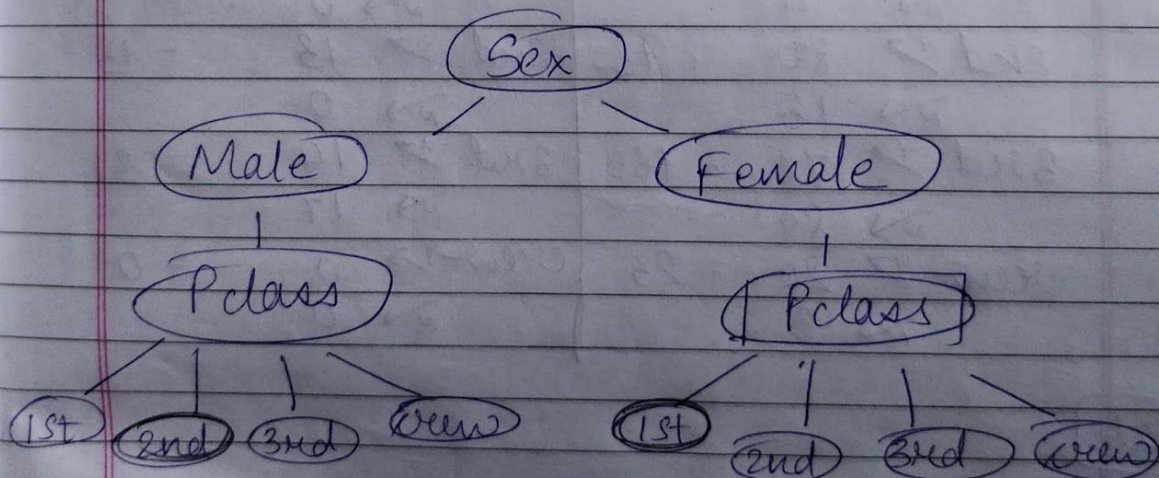
$= 0.1656 + 0.3712$

$= \boxed{0.5368}$

3rd $= -\left(\dfrac{90}{196}\right) \log_2 \left(\dfrac{90}{196}\right) - \left(\dfrac{106}{196}\right) \log_2 \left(\dfrac{106}{196}\right)$

$= 0.5155 + 0.4795$

$= \boxed{0.995}$

new $= -\left(\dfrac{20}{33}\right) \log_2 \left(\dfrac{20}{33}\right) - \left(\dfrac{3}{33}\right) \log_2 \left(\dfrac{3}{33}\right)$

$= 0.1753 + 0.3832$

$= \boxed{0.5585}$

$$\boxed{\text{Adult}} \quad \boxed{\text{Age}}$$

$$\text{Adult} = -\left(\frac{316}{425}\right) \log_2\left(\frac{316}{425}\right) - \left(\frac{109}{425}\right) \log_2\left(\frac{109}{425}\right)$$

$$= +0.3178 + 0.5034$$

$$= \boxed{0.8212}$$

$$\text{child} = -\left(\frac{28}{45}\right) \log_2\left(\frac{28}{45}\right) - \left(\frac{17}{45}\right) \log_2\left(\frac{17}{45}\right)$$

$$= 0.4259 + 0.5305$$

$$= \boxed{0.9569}$$

Sex

Male          Female

Pclass          Pclass

1st  2nd  3rd  crew          1st  2nd  3rd  crew

Male

→ Age

Adult          child

**2) (4 points) Construct two rules using PRISM for the weather dataset. Show the details of your construction. Then, check your solution with Weka and submit a text file of your classifier output window.**

Q2  (1)  We will make 2 rules using
          Prism Algorithm

            Yes ✓                    No

Rule ① : If Outlook = ⟨ Sunny →
                        Overcast
                        Rainy

          play should = [Yes]

          P(Sunny) = 2/5
          P(Overcast) = 4/4
          P(Rainy) = 3/5

          As Overcast = 4/4 = 1 → which
            means regardless of Temperature
          Humidity & Windy attributes
            Play is always    Yes

          Hence;
            If (Outlook = Overcast) = (Play = Yes)

Rule ② : Now, lets see the
            condition when Play = No
              always

| Outlook | Temp | Hum | Windy | Play |
|---------|------|-----|-------|------|
| S | | | | |
| Sunny | hot | high | F | No |
| sunny | hot | high | F | No |
| rainy | cool | normal | T | No |
| Sunny | mild | high | F | No |
| rainy | mild | high | T | No |

If (Outlook = sunny) ∧ (Humidity = High)

= (~~survived~~ = No)

play

**3)(4 points) Classify using Naïve Bayes method on the titanic dataset the data items:**
**2nd child male ?**
**2nd adult female ?**
**Then, check your solution with Weka (the dataset is included with Weka).**



Q3 (1) 2nd | child | male    Survived ?

• P( Survived = No | E) =
 P( Pclass = 2nd | Survived = No) *
 P(Age = child | Survived = No) * P(Sex =
 male | Survived = No) / P(E) =

$$= \left(\frac{52}{1490}\right)\left(\frac{167}{1490}\right)\left(\frac{1364}{1490}\right)\left(\frac{1490}{2201}\right)$$
 _____
 P(E)

$$= \frac{11844976}{48866440100}$$
 P(E)

$$= 0.00242 / P(E)$$

• P (Survived = Yes | E) =
 P( Pclass = 2nd | Survived = Yes) *
 P(Age = child | survived = Yes) * P(Sex =
 male | survived = Yes) / P(E)

$$= \left(\frac{57}{711}\right)\left(\frac{118}{711}\right)\left(\frac{367}{711}\right)\left(\frac{711}{2201}\right)$$
 _____
 P(E)

$$= \frac{2468442}{1112651721}$$
 P(E)

$$= 0.00221 / P(E)$$

P(Survived = Yes|E) + P(Survived = No|

0.00221 /P(E) + 0.00242 /P(E) = 1

P(E) = 0.00221 + 0.00242

So,

P(survived = Yes|E) = $\dfrac{0.00221}{(0.0022 + 0.00242)}$

= 0.477   = 47.7%

P(Survived = No|E) = $\dfrac{0.00242}{0.00221 + 0.00242}$

= 0.522   = 52.2%

Thus, | Survived = No |

② 2nd | adult | female            Survived ?

• P(Survived = No| E) =

P(Pclass = 2nd | survived = No) *

P(Age = Adult | Survived = No) * P(Sex = female| survived = No) /P(E)

= $\dfrac{\left(\dfrac{167}{1490}\right)\left(\dfrac{438}{1490}\right)\left(\dfrac{126}{1490}\right)\left(\dfrac{1490}{2201}\right)}{P(E)}$

= $\dfrac{\dfrac{30258396}{4886440100}}{P(E)}$

$= 0.00619 / P(E)$

○ $P(\text{survived} = \text{Yes} | E) =$
$P(\text{class} = 2\text{nd} | \text{survived} = \text{Yes}) *$
$P(\text{Age} = \text{adult} | \text{survived} ) * P(\text{sex} = \text{Female} | \text{survived} = \text{Yes}) / P(E)$

$= \dfrac{\left(\dfrac{118}{711}\right)\left(\dfrac{654}{711}\right)\left(\dfrac{344}{711}\right)\left(\dfrac{711}{2201}\right)}{P(E)}$

$= \dfrac{\dfrac{26547168}{111265172 1}}{P(E)}$

$= 0.0238 / P(E)$

$P(\text{Survived} = \text{Yes} | E) + P(\text{survived} = \text{No} | E) = 1$
$0.00619 / P(E) + 0.0238 / P(E) = 1$
$P(E) = 0.00619 + 0.0238$

So,
$P(\text{survived} = \text{Yes} | E) = \dfrac{0.0238}{0.00619 + 0.0238}$
$= 0.7959 \quad \boxed{= 79.59\%}$

$P(\text{survived} = \text{No} | E) = \dfrac{0.00619}{0.00619 + 0.0238}$
$= 0.2064 \quad \boxed{= 20.64\%}$

Thus, $\boxed{\text{survived} = \text{Yes}}$

**4.1. Run your classifier by training on traindata.txt and trainlabels.txt then testing on traindata.txt and trainlabels.txt. Report the accuracy in results.txt (along with a comment saying what files you used for the training and testing data). In this situation, you are training and testing on the same data. This is a sanity check: your accuracy should be very high i.e. > 90%**

**(.py and .jpynb code are in separate file which is attached along with results.txt )**

```
Python 3.7.4 Shell                                                    —   □   X
File Edit Shell Debug Options Window Help
Python 3.7.4 (tags/v3.7.4:e09359112e, Jul  8 2019, 19:29:22) [MSC v.1916 32 bit (Intel)] on win32
Type "help", "copyright", "credits" or "license()" for more information.
>>>
 RESTART: C:\Users\Sannath\OneDrive\Documents\Visual Studio 2015\Desktop\python\Programs\Complete code.py
count of accurate labels for traindata.txt :  311
count of inaccurate labels for traindata.txt : 11
Accuracy % : 96.58385093167702
>>> |
```

**4.2. Run your classifier by training on traindata.txt and trainlabels.txt then testing on testdata.txt and testlabels.txt. Report the accuracy in results.txt (along with a comment saying what files you used for the training and testing data). We will not be letting you know beforehand what your performance on the test set should be.**

**(determined test labels of testdata.txt are saved into the  file - 'd_test_labels.txt'.)**