

Homework 4 UG

October 28, 2021

Jose Carlos Munoz

3.10) To find the cost function of both of the Decision tree we use this formula.

$$F_x(n) = \text{Nodes} * \log_2 m + \text{Leafs} * [\log_2 k] + \text{Errors} * \log_2 n \quad (1)$$

Where m is the number of Attributes, k is the number of classes and n is the unknown Sample size. Which are 16, 3 and n respectively

Now we find the cost function for decision Tree A, where Nodes = 2, Leafs = 3, 7 errors

$$\begin{aligned} F_a(n) &= 2 * \log_2 16 + 3 * [\log_2 3] + 7 * \log_2 n \\ F_a(n) &= 2 * 4 + 3 * 2 + 7 * \log_2 n \\ F_a(n) &= 8 + 6 + 7 * \log_2 n \\ F_a(n) &= 14 + 7 * \log_2 n \end{aligned} \quad (a)$$

Now we find the cost function for decision Tree b, where Nodes = 2, Leafs = 3, and 4 errors.

$$\begin{aligned} F_b(n) &= 4 * \log_2 16 + 5 * [\log_2 3] + 4 * \log_2 n \\ F_b(n) &= 4 * 4 + 5 * 2 + 4 * \log_2 n \\ F_b(n) &= 16 + 10 + 4 * \log_2 n \\ F_b(n) &= 26 + 4 * \log_2 n \end{aligned} \quad (b)$$

Using the MDL paradigm we need to find a $L_S(h) + \sqrt{\frac{\log(\frac{2}{\delta}) + |h|}{2m}}$. The better decision tree is the one that gives us the lowest value. m is found as the sample size which is 200, δ is given as .99. and the $L_S(h)$ is the cost function of the Decision tree that we found.

solving for both, when n is 200, Decision tree B is much lower in value. Therefore it is the best Decision Tree of the two.

4.6)a

$$\begin{aligned} P(S|UG) &= .15 \\ P(S|G) &= .23 \\ P(G) &= .2 \\ P(UG) &= .8 \end{aligned} \quad (1)$$

These are the known probabilities.

From this we can find $P(G|S)$.

Because of Bayes Theorem $P(G|S)$ is the same as the following

$$P(G|S) = \frac{P(S|G) * P(G)}{P(S)} \quad (2)$$

$P(S)$ can be found as

$$\begin{aligned} P(S) &= P(S|G) * P(G) + P(S|UG) * P(UG) \\ P(S) &= .23 * .2 + .15 * .8 \\ P(S) &= .166 \end{aligned} \quad (3)$$

Therefore

$$\begin{aligned} P(G|S) &= \frac{.23 * .2}{.166} \\ P(G|S) &= .277 \end{aligned} \tag{4}$$

So the probability that a smoker is a graduate student is .277

4.6)c

The probability that a smoker is a graduated student can be written as $P(UG|S)$.

$$\begin{aligned} P(UG|S) &= \frac{P(S|UG) * P(UG)}{P(S)} \\ P(UG|S) &= \frac{.23 * .8}{.277} \\ P(UG|S) &= .857 \end{aligned} \tag{5}$$

So the probability that a smoker is an undergrad is .857.

Since $P(UG|S) > P(G|S)$ we can conclude we have a higher chance of finding an undergrad that is a smoker

4.6)d

$$\begin{aligned} P(D|UG) &= .1 \\ P(D|G) &= .3 \\ P(D) &= P(D|UG) * P(UG) + P(D|G) * P(G) \\ P(D) &= 0.1 * .8 + .2 * .3 \\ P(D) &= .14 \\ P(D,S|G) &= P(D|G) * P(S|G) \\ P(D,S|G) &= .3 * .23 \\ P(D,S|G) &= .069 \\ P(D,S|UG) &= P(D|UG) * P(S|UG) \\ P(D,S|UG) &= .1 * .15 \\ P(D,S|UG) &= 0.015 \\ P(D,S) &= Q \end{aligned} \tag{6}$$

These are the known probabilities. Since we don't know what $P(D,S)$ is, we set it as a constant Q

Now we can find the values for $P(G|D,S)$ and $P(UG|D,S)$

$$\begin{aligned}
P(\text{UG}|\text{D},\text{S}) &= \frac{P(\text{D},\text{S}|\text{UG}) * P(\text{UG})}{P(\text{D},\text{S})} \\
P(\text{UG}|\text{D},\text{S}) &= \frac{.015 * .8}{Q} \\
P(\text{UG}|\text{D},\text{S}) &= \frac{.012}{Q} \\
P(\text{G}|\text{D},\text{S}) &= \frac{P(\text{D},\text{S}|\text{UG}) * P(\text{UG})}{P(\text{D},\text{S})} \\
P(\text{G}|\text{D},\text{S}) &= \frac{.069 * .2}{Q} \\
P(\text{G}|\text{D},\text{S}) &= \frac{.0139}{Q}
\end{aligned} \tag{7}$$

From these results we can conclude that the chance that we find a graduate that lives in a dorm and is a smoker is higher than the chance that we find an undergraduate that lives in a dorm and is a smoker.

4.7)a

$$\begin{aligned}
P(\text{A}=0|+) &= \frac{2}{5} = .4 \\
P(\text{A}=0|-) &= \frac{3}{5} = .6 \\
P(\text{A}=1|+) &= \frac{3}{5} = .6 \\
P(\text{A}=1|-) &= \frac{2}{5} = .4 \\
P(\text{B}=0|+) &= \frac{4}{5} = .8 \\
P(\text{B}=0|-) &= \frac{3}{5} = .6 \\
P(\text{B}=1|+) &= \frac{1}{5} = .2 \\
P(\text{B}=1|-) &= \frac{2}{5} = .4 \\
P(\text{C}=0|+) &= \frac{3}{5} = .6 \\
P(\text{C}=0|-) &= \frac{0}{5} = 0 \\
P(\text{C}=1|+) &= \frac{2}{5} = .4 \\
P(\text{C}=1|-) &= \frac{5}{5} = .1
\end{aligned} \tag{8}$$

4.7)b

we are task to find $P(A=0,B=1,C=0|+)$. Using the Bayes Therm we canfind the value as

$$\begin{aligned}
P(+|A=0,B=1,C=0) &= \frac{P(A=0,B=1,C=0|+) * P(+)}{P(A=0,B=1,C=0)} \\
P(+|A=0,B=1,C=0) &= \frac{P(A=0|+) * P(B=1|+) * P(C=0|+) * P(+)}{P(A=0,B=1,C=0)} \\
P(+|A=0,B=1,C=0) &= \frac{.4 * .2 * .6 * .5}{P(A=0,B=1,C=0)} \\
P(+|A=0,B=1,C=0) &= \frac{0.024}{P(A=0,B=1,C=0)}
\end{aligned} \tag{9}$$

$$\begin{aligned}
P(-|A=0,B=1,C=0) &= \frac{P(A=0,B=1,C=0|-) * P(-)}{P(A=0,B=1,C=0)} \\
P(-|A=0,B=1,C=0) &= \frac{P(A=0|-) * P(B=1|-) * P(C=0|-) * P(-)}{P(A=0,B=1,C=0)} \\
P(-|A=0,B=1,C=0) &= \frac{.6 * .4 * 0 * .5}{P(A=0,B=1,C=0)} \\
P(-|A=0,B=1,C=0) &= 0
\end{aligned} \tag{10}$$

From these results we canconlude that the class label for $(A=0, B=1, C=0)$ will be Class +.

4.7)c

We will be looking at the conditional probabillites for the them all over again with the m-estimate. When $m=4$ and $p = 1/2$; to find the new Conditonal probabillites we use this equation

$$\frac{n_c + m * p}{n + m} \tag{11}$$

so now the The conditional probablities will be

$$\begin{aligned}
P(A=0|+) &= \frac{2+2}{5+4} = \frac{4}{9} \\
P(A=0|-) &= \frac{3+2}{5+4} = \frac{5}{9} \\
P(A=1|+) &= \frac{3+2}{5+4} = \frac{5}{9} \\
P(A=1|-) &= \frac{2+2}{5+4} = \frac{4}{9} \\
P(B=0|+) &= \frac{4+2}{5+4} = \frac{6}{9} \\
P(B=0|-) &= \frac{3+2}{5+4} = \frac{5}{9} \\
P(B=1|+) &= \frac{1+2}{5+4} = \frac{3}{9} \\
P(B=1|-) &= \frac{2+2}{5+4} = \frac{4}{9} \\
P(C=0|+) &= \frac{3+2}{5+4} = \frac{5}{9} \\
P(C=0|-) &= \frac{0+2}{5+4} = \frac{2}{9} \\
P(C=1|+) &= \frac{2+2}{5+4} = \frac{4}{9} \\
P(C=1|-) &= \frac{5+2}{5+4} = \frac{7}{9}
\end{aligned} \tag{12}$$

4.7)d

we repeat b) but with the m-estimate conditional probabilities

$$\begin{aligned}
P(+|A=0,B=1,C=0) &= \frac{P(A=0,B=1,C=0|+) * P(+)}{P(A=0,B=1,C=0)} \\
P(+|A=0,B=1,C=0) &= \frac{P(A=0|+) * P(B=1|+) * P(C=0|+) * P(+)}{P(A=0,B=1,C=0)}
\end{aligned} \tag{13}$$

$$\begin{aligned}
P(+|A=0,B=1,C=0) &= \frac{\frac{4}{9} * \frac{3}{9} * \frac{5}{9} * .5}{P(A=0,B=1,C=0)} \\
P(+|A=0,B=1,C=0) &= \frac{0.0142}{P(A=0,B=1,C=0)}
\end{aligned}$$

$$\begin{aligned}
P(-|A=0,B=1,C=0) &= \frac{P(A=0,B=1,C=0|-) * P(-)}{P(A=0,B=1,C=0)} \\
P(-|A=0,B=1,C=0) &= \frac{P(A=0|-) * P(B=1|-) * P(C=0|-) * P(-)}{P(A=0,B=1,C=0)}
\end{aligned} \tag{14}$$

$$\begin{aligned}
P(-|A=0,B=1,C=0) &= \frac{\frac{5}{9} * \frac{4}{9} * \frac{2}{9} * .5}{P(A=0,B=1,C=0)} \\
P(-|A=0,B=1,C=0) &= \frac{0.0274}{P(A=0,B=1,C=0)}
\end{aligned}$$

From these results we can conclude that the class label for $(A=0, B=1, C=0)$ is class +

4.7)e

The better method would be the m-estimate because we do not want our entire expression to be zero

Problem B)

A) A benefit of having only 2 dimensions is that it will simplify the computation and its complexity is reduced. A disadvantage is that it will be missing out on other relationships with other attributes. This can greatly decrease the accuracy of the model.

A benefit of having many dimensions is that the accuracy of it is much higher. However, a problem is that computation is slow and the complexity of it is much higher

b) If we have more unlabeled samples of A, P, and I, then the first algorithm will be more accurate. However, if the reverse occurs, then the second algorithm will be much more accurate.