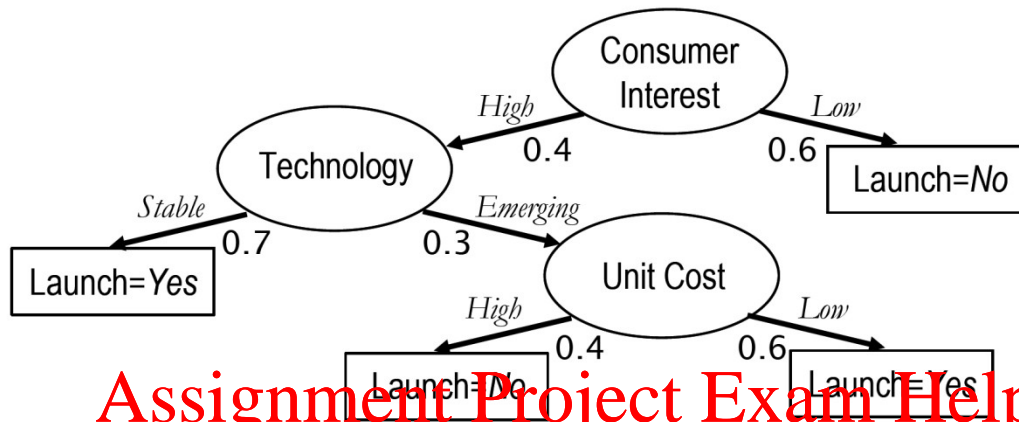


Assignment 3

273 Business Intelligence for Analytical Decisions

This assignment must be completed individually. Submit Word file to online drop box on Canvas. Write your name in the Word file.

Q.1. Consider a decision tree (as shown below) for launching new technology products:

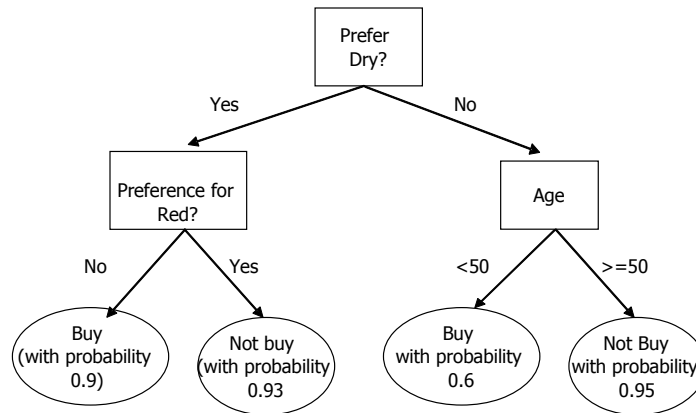


The branching probabilities are provided. Given this decision tree, find the probability $\Pr[\text{Launch} = \text{Yes} \mid \text{evidence}]$ in each of the following cases.

- (a) Consumer Interest = *High*, Technology = *Emerging*
- (b) Consumer Interest = *High*, Unit cost = *Low*
- (c) Technology = *Emerging*

Q. 2. A winery maintains a dataset containing information about customers who subscribe to its tasting events and special offers for wine cases. The winery occasionally mails tasting samples of new wines in an effort to increase sales. The chief marketing officer is aiming to send samples of a newly produced wine to customers who are NOT likely to place an order for the new wine (when probability of purchase is less than 0.5). Based on a survey conducted amongst its customers and their willingness to buy a case of the new wine before tasting it, the firm collected tree attributes – whether customers *prefer dry*, or have *preference for red wine* and their *ages*. The classification problem is

to predict whether customers will buy or not buy wine, with certain confidence/probability. The following classification tree was induced:



Answer the following questions:

- 1) After running a data mining software, suppose the above tree is generated. If you have to choose a single attribute to predict which customer is likely to place an order or not, which attribute would you use (check one and briefly justify your choice)?

- A. Preference for red wine
- B. Age (whether the customer is older or younger than 50 years old)
- C. Preference for dry wine
- D. Impossible to determine given the information provided.

- 2) The Winery manager, Barbara, wants to consult with you if she should send a sample to a new customer named George. She tells you that George prefers dry wine and strongly prefers red wine. What's your recommendation (on whether to send a sample to George)? Justify. (Note: The chief marketing officer is aiming to send samples of a newly produced wine to customers who are NOT likely to place an order for the new wine (when probability of purchase is less than 0.5).)
- 3) Assume that the cost of producing and shipping a wine sample to a customer is \$8 and the revenue from each order is \$100. Assume the company ascertained that the probability of purchase would become 17% for a customer who receives a sample of a new wine. What is the expected incremental revenue and the incremental cost of shipping a sample to customers preferring dry and red wines? Would you suggest shipping a sample to customers preferring dry and red wines? Explain your answer.

Q. 3. An NBA specialist is trying to predict each team's likelihood of winning the championship (so this is the dependent variable). She decides to use two predictors (i.e. independent variables) – 1) whether a team won more than 55 games during the past season and 2) whether the team as a whole is healthy. The dataset below contains information about the top eight teams and whether the specialist thinks they have a chance to win.

Team	Win less than 55 games?	Team Healthy?	Contender to win the championship?
Phoenix Suns	No	Fair	No
Detroit Pistons	No	Excellent	No
San Antonio Spurs	No	Fair	Yes
Miami Heat	No	Fair	Yes
Denver Nuggets	Yes	Fair	Yes
Seattle Supersonics	Yes	Excellent	No
Houston Rockets	Yes	Excellent	Yes
Dallas Mavericks	No	Excellent	No

Use the examples in the above database to determine which attribute you should split on **first**, in order to build a decision tree to predict whether a team is a contender to win. Explain each step and show **all** relevant computations. To simplify your computation, you are given that the information gain if splitting on “team healthy” is 0.189.

You may need the following logarithm values for answering the question:

$$\log_2 \frac{2}{3} = -0.585, \quad \log_2 \frac{1}{3} = -1.585, \quad \log_2 \frac{1}{2} = -1, \quad \log_2 \frac{3}{5} = -0.737, \quad \log_2 \frac{2}{5} = -1.322$$

$$\log_2 \frac{1}{4} = -2, \quad \log_2 \frac{3}{4} = -0.415$$

Q. 4. A bank selected 12 inputs to predict whether each applicant defaults. The output variable (BAD) indicates whether an applicant defaulted on the home equity line of credit. The following table describes the variables used.

Name	Type	Description
BAD	Binary	1=applicant defaulted on loan or seriously delinquent 0=applicant paid loan
CLAGE	Numeric	Age of oldest credit line in months
CLNO	Numeric	Number of credit lines
DEBTINC	Numeric	Debt-to-income ratio
DELINQ	Numeric	Number of delinquent credit lines
DEROG	Numeric	Number of major derogatory reports
JOB	Nominal	Occupational categories
LOAN	Numeric	Amount of the loan request
MORTDUE	Numeric	Amount due on existing mortgage
NINQ	Numeric	Number of recent credit inquiries
REASON	Binary	DebtCon=debt consolidation HomeImp=home improvement
VALUE	Numeric	Value of current property
YOJ	Numeric	Years at present job

Question 4.a: Download the file HMEQ.arff from Canvas. Load HMEQ.arff into WEKA. What's the class distribution for BAD (i.e. percentage of 1s vs. percentage of 0s)? Build a classification model to predict whether a customer will default using decision tree model J48 (classifiers -> trees -> J48). Please use the default parameter setting for J48, and choose Percentage split (66%) for Test options. Please report the classification accuracy rate of the model built, as well as the confusion matrix. Use confusion matrix to comment on how good the model is.

4.b Now, go back to the Preprocess window where you have the HMEQ.arff file open. Under Filter, choose filters-> supervised->instance->Resample. In the parameter window for Resample, change biasToUniformClass to 1 while leaving other parameters unchanged (i.e. invertSelection=False, noReplacement=False, sampleSizePercent=100). Click on Apply and now report the class distribution of BAD. Use the same model set up in 4.a to build the model. Please report the classification accuracy rate of the model built, as well as the confusion matrix. Please compare with the prior class distribution and comment on how good the model is, and also use confusion matrix to comment on how good the model is, especially compared to the model in 4.a.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder