

CS 584: Machine Learning

Pradyot Mayank(A20405826)

Spring 2019 Assignment 3

Question 1 (50 points)

You will use the CART algorithm to build profiles of credit card holders. The data is the CustomerSurveyData.csv. The analysis specifications are:

Target Variable

- **CarOwnership.** The type of car ownership. This variable has three non-missing categories which are *Leased*, *None*, and *Own*.
- Drop all missing values in the target variable.

Nominal Predictors

- **CreditCard.** The type of credit card held. This variable has five categories which are *American Express*, *Discover*, *MasterCard*, *Others*, and *Visa*.
- **JobCategory.** The category of the job held. This variable has six non-missing categories which are *Agriculture*, *Crafts*, *Labor*, *Professional*, *Sales*, and *Service*.
- Recode all the missing values into the *Missing* category.

You will use the Entropy metric as the splitting criterion. You may want to write a Python program to assist you in answering the questions.

- a) (5 points). What is the Entropy metric for the root node?

Root node Entropy=1.0744900777690842

- b) (5 points). How many possible binary-splits that you can generate from the CreditCard predictor?

K=5

No of binary splits= $2^{k-1}-1=15$

- c) (10 points). Calculate the Entropy metric for each possibly binary split that you can generate from the CreditCard predictor. List your answers in a table. The table should have three columns: an index of the split, the contents of the two branches, and the split entropy metric.

SplitIndex	Left_Child	Right_Child	SplitEntropy	Counter
[1.0468094317797334, 1.1479670231237966]	['Visa']	['American Express', 'Discover', 'MasterCard', 'Others']	1.072038135	1
[1.0725070489912254, 1.075008859357403]	['American Express']	['Discover', 'MasterCard', 'Others', 'Visa']	1.073000406	1
[1.1085058090287196, 0.973200661032307]	['Discover']	['American Express', 'MasterCard', 'Others', 'Visa']	1.072135785	1
[1.0689643982639212, 1.174264078521314]	['Others']	['American Express', 'Discover', 'MasterCard', 'Visa']	1.073660764	1
[1.0760277584181488, 1.069359232936541]	['MasterCard']	['American Express', 'Discover', 'Others', 'Visa']	1.074427312	1
[1.0348662734428489, 1.1211075887224382]	['American Express', 'Visa']	['Discover', 'MasterCard', 'Others']	1.073381645	2
[1.0852881500406055, 1.0638912768034012]	['Discover', 'Visa']	['American Express', 'MasterCard', 'Others']	1.07420029	2
[1.0369375321912813, 1.1522460232150684]	['Others', 'Visa']	['American Express', 'Discover', 'MasterCard']	1.070838229	2
[1.0351668498221591, 1.1109476274999246]	['MasterCard', 'Visa']	['American Express', 'Discover', 'Others']	1.072253962	2
[1.1167776104815863, 1.018286062470057]	['American Express', 'Discover']	['MasterCard', 'Others', 'Visa']	1.070880549	2
[1.065626118599882, 1.0992153713223232]	['American Express', 'Others']	['Discover', 'MasterCard', 'Visa']	1.073748	2
[1.073826633163597, 1.0739768609820293]	['American Express', 'MasterCard']	['Discover', 'Others', 'Visa']	1.073892313	2
[1.1034337362362638, 1.007246961317521]	['Discover', 'Others']	['American Express', 'MasterCard', 'Visa']	1.073288801	2
[1.1269534892368755, 1.0203280170186693]	['Discover', 'MasterCard']	['American Express', 'Others', 'Visa']	1.072702449	2
[1.068678222097913, 1.087853923191551]	['MasterCard', 'Others']	['American Express', 'Discover', 'Visa']	1.074135627	2
[1.087853923191551, 1.068678222097913]	['American Express', 'Discover', 'Visa']	['MasterCard', 'Others']	1.074135627	3
[1.0203280170186693, 1.1269534892368755]	['American Express', 'Others', 'Visa']	['Discover', 'MasterCard']	1.072702449	3
[1.007246961317521, 1.1034337362362638]	['American Express', 'MasterCard', 'Visa']	['Discover', 'Others']	1.073288801	3
[1.0739768609820293, 1.073826633163597]	['Discover', 'Others', 'Visa']	['American Express', 'MasterCard']	1.073892313	3
[1.0992153713223232, 1.065626118599882]	['Discover', 'MasterCard', 'Visa']	['American Express', 'Others']	1.073748	3
[1.018286062470057, 1.1167776104815863]	['MasterCard', 'Others', 'Visa']	['American Express', 'Discover']	1.070880549	3
[1.1109476274999246, 1.0351668498221591]	['American Express', 'Discover', 'Others']	['MasterCard', 'Visa']	1.072253962	3
[1.1522460232150684, 1.0369375321912813]	['American Express', 'Discover', 'MasterCard']	['Others', 'Visa']	1.070838229	3
[1.0638912768034012, 1.0852881500406055]	['American Express', 'MasterCard', 'Others']	['Discover', 'Visa']	1.07420029	3
[1.1211075887224382, 1.0348662734428489]	['Discover', 'MasterCard', 'Others']	['American Express', 'Visa']	1.073381645	3
[1.069359232936541, 1.0760277584181488]	['American Express', 'Discover', 'Others', 'Visa']	['MasterCard']	1.074427312	4
[1.174264078521314, 1.0689643982639212]	['American Express', 'Discover', 'MasterCard', 'Visa']	['Others']	1.073660764	4
[0.973200661032307, 1.1085058090287196]	['American Express', 'MasterCard', 'Others', 'Visa']	['Discover']	1.072135785	4
[1.075008859357403, 1.0725070489912254]	['Discover', 'MasterCard', 'Others', 'Visa']	['American Express']	1.073000406	4
[1.1479670231237966, 1.0468094317797334]	['American Express', 'Discover', 'MasterCard', 'Others']	['Visa']	1.072038135	4

- d) (5 points). What is the optimal split for the CreditCard predictor?

Optimal Split entropy for card: 1.0708382285522746

['Others', 'Visa'] and ['American Express', 'Discover', 'MasterCard']

- e) (5 points). How many possible binary-splits that you can generate from the JobCategory predictor?

K=7

No of binary splits= $2^{k-1}-1=63$

- f) (10 points). Calculate the Entropy metric for each possibly binary split that you can generate from the JobCategory predictor. List your answers in a table. The table should have three columns: an index of the split, the contents of the two branches, and the split Entropy metric.

[illegible]

- g) (5 points). What is the optimal split for the JobCategory predictor?

Optimal Split entropy for Job: 1.0720111150297396 since it has the lowest entropy.

['Agriculture', 'Sales'] and ['Crafts', 'Labor', 'Missing', 'Professional', 'Service']

- h) (5 points). Between the CreditCard and the JobCategory predictors, which predictor will you choose for producing the second layer (i.e., depth 1) of your decision tree?

Split Entropy=1.0708382285522746

We will choose the predictor Creditcard or producing the second layer of the tree ['Others', 'Visa'] and ['American Express', 'Discover', 'MasterCard'] as the entropy is lowest from both the predictors.

Question 2 (50 points)

In 2014, Allstate provided the data on Kaggle.com for the Allstate Purchase Prediction Challenge which is open. The data contain transaction history for customers that ended up purchasing a policy. For each Customer ID, you are given their quote history and the coverage options they purchased.

The data is available on the Blackboard as Purchase_Likelihood.csv. It contains 665,249 observations on 97,009 unique Customer ID. We are going to use the MNLogit function to build a multinomial logistic model to predict purchase likelihood of coverage A using three predictors. The target variable is **A** which have these categories 0, 1, and 2. The nominal predictors are (categories are inside the parentheses):

1. **group_size**. How many people will be covered under the policy (1, 2, 3 or 4)?
2. **homeowner**. Whether the customer owns a home or not (0=no, 1=yes)
3. **married_couple**. Does the customer group contain a married couple (0=no, 1=yes)

Please build a multinomial logistic model using and answer the following questions.

- a) (2 points) Suppose you start with a model with only the Intercept term (i.e., without any predictors). How many parameters are in this model?

For intercept only model we only have one parameter is involved

- b) (3 points) What are the marginal counts of the categories of the target variable A?

ToatlCount_0	ToatlCount_1	ToatlCount_2	ToatlCount
143691	426067	95491	665249

- c) (5 points) Without calling the MNLogit function, what are the maximum likelihood estimates of the predicted probabilities $\pi_{ij}, j = 1, 2, 3$ of this Intercept-only model? Show all the necessary steps and the estimates for the $\pi_{ij}, j = 1, 2, 3$. (Hint: equate the first derivatives of the log-likelihood function to zeros for this Intercept-only model)

Total Observation	MLE	Odds	MLE Intercept
143691.0	0.215996	1.000000	0.000000
426067.0	0.640462	2.965161	1.086931
95491.0	0.143542	0.664558	-0.408633

- d) (3 points) What is the log-likelihood value of this Intercept-only model? (Hint: the log-likelihood function is $l = \sum_{i=1}^m \sum_{j=1}^3 n_{ij} \log_e(\pi_{ij})$)

(tototalCount_0)log_e(prob_0) + (tototalCount_1)log_e(prob_1) + (tototalCount_2)log_e(prob_2)
143691 log_e(0.215996) + 426067 log_e(0.640462) + 95491 log_e(0.143542)
Log Likelihood: -595406.7618844224

- e) (5 points) Next, you are asked to mathematically calculate the maximum likelihood estimates of the Intercept terms $\beta_{j0}, j = 1, 2, 3$. The convention is to set the Intercept term to zero for the target category A = 0, i.e., $\beta_{10} = 0$. (Hint: use the mathematical formula of the logit of π_{ij} (i.e., $\log_e(\pi_{ij}/\pi_{i1})$ for this Intercept only model, then solve for the betas)?

J = 1, 2, 3

$\beta = [\beta_{10}, \beta_{20}, \beta_{30}]$

π_{i1} = Probability of 0= 0.21599581510081187

π_{i2} = Probability of 1= 0.64046244338586

π_{i3} = Probability of 2= 0.1435417415133281

$\beta_j = \log_e(\pi_{ij}/\pi_{i1})$

$\beta_1 = \log_e(\pi_{i0}/\pi_{i1}) = \log_e(1) = 0$

$\beta_2 = \log_e(\pi_{i1}/\pi_{i1}) = \log_e(0.6404624433858/0.2159958151008) = 1.08693$

$\beta_3 = \log_e(\pi_{i3}/\pi_{i1}) = \log_e(0.143541741513/0.2159958151008) = -0.40863$

- f) (5 points) Create and display a contingency table where group_size, homeowner, and married_couple are on the row dimension, and A is on the column dimension. The cell contents are the row percentages of the categories of A per each level combination of group_size, homeowner, and married_couple.

group_size	homeowner	married_couple	0	1	2	All
1	0	0	25.77842	59.14609	15.07549	100
1	0	1	32.14196	51.69897	16.15907	100
1	1	0	18.02423	68.62947	13.34631	100
1	1	1	22.17153	62.39966	15.42881	100
2	0	0	27.27079	55.59455	17.13465	100
2	0	1	20.55084	64.5793	14.86986	100
2	1	0	25.33176	59.48539	15.18285	100
2	1	1	16.06721	70.20728	13.72551	100
3	0	0	27.55906	66.92913	5.511811	100
3	0	1	23.30023	59.51108	17.18869	100
3	1	0	25.99532	59.48478	14.51991	100
3	1	1	26.15788	56.30644	17.53568	100
4	0	0	100	0	0	100
4	0	1	18.75	67.85714	13.39286	100
4	1	0	48.48485	42.42424	9.090909	100
4	1	1	33.33333	53.00546	13.6612	100
All			21.59958	64.04624	14.35417	100

- g) (2 points) Based on the contingency table in e), do you expect the separation or the quasi-separation phenomenon to occur when we build the multinomial logistic model which has group_size, homeowner, and married_couple as the predictors.

Based on the contingency table in (f), we expect the Quasi-complete separation because none of the predictor completely define particular category of target variable.

- h) (5 points) Now, you will use the MNLogit function to build the multinomial logistic model which has group_size, homeowner, and married_couple as the predictors. What value of the target variable A is used by the MNLogit function as the reference category? Next, what is the log-likelihood value of this model? Finally, how many parameters (including the redundant ones) are in the model?

Reference Category Value of A=0

Log-Likelihood=-591936.7938327907

Since there are 9 parameters and each parameter has its occurrence and nonoccurrence i.e.

So,Parameters Involved=9*2=18

- i) (10 points) What are the values of group_size, homeowner, and married_couple such that the odd $\text{Prob}(A=1)/\text{Prob}(A=0)$ will attain its maximum? What is the maximum odd $\text{Prob}(A=1)/\text{Prob}(A=0)$ value?

Predicted probability of 0: 0.1606721019540383

Predicted probability of 1: 0.7020727585496748

Maximum odd Probability: 4.36959963809093

- j) (5 points) According to the multinomial logistic model, what is the odds ratio for group_size = 3 versus group_size = 1, and A = 2 versus A = 0? Mathematically, the odds ratio is $(\text{Prob}(A=2)/\text{Prob}(A=0) \mid \text{group_size} = 3) / ((\text{Prob}(A=2)/\text{Prob}(A=0) \mid \text{group_size} = 1))$.

group_size_3A2=-0.107280

group_size_1A0=0.411100

group_size_3A2-group_size_1A0

Required Odds Ratio: 0.5954844518092098

- k) (5 points) According to the multinomial logistic model, what is the odds ratio for group_size = 1 versus group_size = 3, and A = 2 versus A = 1? Mathematically, the odds ratio is $(\text{Prob}(A=2)/\text{Prob}(A=1) \mid \text{group_size} = 1) / ((\text{Prob}(A=2)/\text{Prob}(A=1) \mid \text{group_size} = 3))$.

group_size_3A2=-0.026011

group_size_1A1=0.108286

group_size_1A1-group_size_3A2

Required Odds Ratio: 1.1437324577458428