CS 584-04: Machine Learning

Fall 2018 Assignment 3

Suhas Sreenivas - A20423132

Question 1

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

The Gini metric for the root node is 0.7670692.

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

$$2^{k-1}-1$$

$$=2^{3-1}-1=3$$

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

Contents of the 2 branches	Sequence index of the split	Split GINI metric
(Lease, none) and own	0,1 and 2	0.7660055295576353
(None, own) and Lease	1,2 and 0	0.7657091011467262
(Own, Lease) and None	0,2 and 1	0.7667762687160046

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

(Own, None) and Lease with a GINI index of 0.7657091011467262

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

$$2^{7-1} - 1 = 64 - 1 = 63$$

f) (10 points). Calculate the Gini 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: the sequence index of the split, the contents of the two branches, the split Gini metric.

	Contents of the 2 branches	Sequence index	Split GINI metric
1	Agriculture and (craft,labor,missing,professional,sales,service)	[(0), (1,2,3,4,5,6)]	0.7670516826657 839
2	craft and (Agriculture,labor,missing,professional,sale s,service)	[(1,), (0, 2, 3, 4, 5, 6)]	0.7669942084821 879
3	labor and (craft, Agriculture,missing,professional,sales,servi ce)	[(2,), (0, 1, 3, 4, 5, 6)]	0.7670315160755 341
4	missing and (craft,labor, Agriculture,professional,sales,service)	[(3,), (0, 1, 2, 4, 5, 6)]	0.7668979204279 504
5	professional and (craft,labor,missing, Agriculture,sales,service)	[(4,), (0, 1, 2, 3, 5, 6)]	0.7667678180799 103
6	sales and (craft,labor,missing,professional, Agriculture,service)	[(5,), (0, 1, 2, 3, 4, 6)]	0.7666854311563 086
7	service and (craft,labor,missing,professional,sales, Agriculture)	[(6,), (0, 1, 2, 3, 4, 5)]	0.7670150066283 693
8	Agriculture, Crafts and (labor,missing,professional,sales,service)	[(0, 1), (2, 3, 4, 5, 6)]	0.7670091884141 732
9	Agriculture Labor and (craft,missing,professional,sales,service)	[(0, 2), (1, 3, 4, 5, 6)]	0.7670390199142 361
10	Agriculture, Missing and (craft,labor,professional,sales,service)	[(0, 3), (1, 2, 4, 5, 6)]	0.7670229036125 564
11	Agriculture, Professional and (craft,labor,missing,sales,service)	[(0, 4), (1, 2, 3, 5, 6)]	0.7667752040720 033
12	Agriculture, Sales and (craft,labor,missing,professional,service)	[(0, 5), (1, 2, 3, 4, 6)]	0.7667187854366 833
13	Agriculture, Service and (craft,labor,missing,professional,sales)	[(0, 6), (1, 2, 3, 4, 5)]	0.7670277019046 212
14	craft Labor and (Agriculture, missing, professional, sales, service)	[(1, 2), (0, 3, 4, 5, 6)]	0.7669835211091 989
15	craft, Missing and (Agriculture, labor, professional, sales, service)	[(1, 3), (0, 2, 4, 5, 6)]	0.7669933691118 804
16	craft, Professional and (Agriculture, labor, missing, sales, service)	[(1, 4), (0, 2, 3, 5, 6)]	0.7667013453310 396

17	craft, Sales and (Agriculture, labor, missing, professional, service)	[(1, 5), (0, 2, 3, 4, 6)]	0.8059158490275
18	craft, Service and (Agriculture, labor, missing, professional, sale s)	[(1, 6), (0, 2, 3, 4, 5)]	0.7670440336504 847
	,		
19	Labor, Missing and (Agriculture,	[(2, 3), (0, 1, 4, 5,	0.7670487631566
	Crafts,professional,sales,service)	6)]	11
20	Labor, Professional and (Agriculture,	[(2, 4), (0, 1, 3, 5,	0.7668199919361
	Crafts,missing,sales,service)	6)]	81
21	Labor, Sales and (Agriculture,	[(2, 5), (0, 1, 3, 4,	0.7667483725125
	Crafts,missing,professional,service)	6)]	881
22	Labor, Service and (Agriculture,	[(2, 6), (0, 1, 3, 4,	0.7669986871637
	Crafts,missing,professional,sales)	5)]	38
	Crares, missing, professional, sales	3/1	
23	Professional, Missing and (Agriculture,	[(3, 4), (0, 1, 2, 5,	0.7667445819476
25	Crafts, Labor, sales, service)	6)]	135
24	Missing, Sales and (Agriculture,	[(3, 5), (0, 1, 2, 4,	0.7666929425599
24	Crafts, Professional, sales, service)	6)]	277
25	Missing, Service and (Agriculture,	[(3, 6), (0, 1, 2, 4,	0.7670279564538
23	Crafts, Labor, professional, service)	5)]	968
	Crarts, Labor, professionar, service)	3/]	308
26	Professional, Sales and (Agriculture, Crafts,	[(4, 5), (0, 1, 2, 3,	0.7670179293290
20	Labor, Missing, Service)	6)]	057
27	Professional, Service and (Agriculture,	[(4, 6), (0, 1, 2, 3,	0.7668363333333
-′	Crafts, Labor, Missing, Sales	5)]	333
	Oranto, Labor, missing, cares	7,1	
28	Sales , Services and (Agriculture, Crafts,	[(5, 6), (0, 1, 2, 3,	0.7666814182609
	Labor, Missing, Professional)	4)]	785
		-71	
29	[(Agriculture, Craft, Labor), (Missing,	[(0, 1, 2), (3, 4, 5,	0.7669926818873
	Professional, Sales, Services)]	6)]	669
30	[(Agriculture, Craft, Missing), (Labor,	[(0, 1, 3), (2, 4, 5,	0.7670014692093
	Professional, Sales, Services)]	6)]	515
31	[(Agriculture, Craft, Professional), (Labor,	[(0, 1, 4), (2, 3, 5,	0.7667009008842
	Missing, Sales, Services)]	6)]	012
32	[(Agriculture, Craft, Sales), (Labor, Missing,	[(0, 1, 5), (2, 3, 4,	0.7668272791206
32	Professional, Services)]	6)]	003
33	[(Agriculture, Craft, Services), (Labor,	[(0, 1, 6), (2, 3, 4,	0.7670474978957
	Missing, Professional, Sales)]	5)]	71
		~/1	, · -
34	[(Agriculture, Labor, Missing), (Craft,	[(0, 2, 3), (1, 4, 5,	0.7670475434223
"	Professional, Sales, Services)]	6)]	493
35	[(Agriculture, Labor, Professional), (Craft,	[(0, 2, 4), (1, 3, 5,	0.7668140712782
	Missing, Sales, Services)]	6)]	201
	IVIIISSIIIB, Saics, Scivices/		201

36	[(Agriculture, Labor, Sales), (Craft, Missing,	[(0, 2, 5), (1, 3, 4,	0.7667641833913
	Professional, Services)]	6)]	141
37	[(Agriculture, Labor, Services), (Craft,	[(0, 2, 6), (1, 3, 4,	0.7670088886265
	Missing, Professional, Sales)]	5)]	447
38	[(Agriculture, Missing, Professional), (Craft,	[(0, 3, 4), (1, 2, 5,	0.7667497953159
	Labor, Sales, Services)]	6)]	906
39	[(Agriculture, Missing, Sales), (Craft, Labor,	[(0, 3, 5), (1, 2, 4,	0.7667220133097
	Professional, Services)]	6)]	015
40	[(Agriculture, Missing, Services), (Craft,	[(0, 3, 6), (1, 2, 4,	0.7670320641598
	Labor, Professional, Sales)]	5)]	185
		,,	
41	[(Agriculture, Professional, Sales), (Craft,	[(0, 4, 5), (1, 2, 3,	0.7670153981729
-	Labor, Missing, Services)]	6)]	524
42	[(Agriculture, Professional, Services), (Craft,	[(0, 4, 6), (1, 2, 3,	0.7668315404951
	Labor, Missing, Sales)]	5)]	73
	20001, 1411351116, 30103/1	3/1	7.5
43	[(Agriculture, Sales, Services), (Craft, Labor,	[(0, 5, 6), (1, 2, 3,	0.7666997631427
45	Missing, Professional)]		3
	ivilssing, Professionary	4)]	3
11	[(Croft Labor Missing) (Agriculture	[/1 2 2) /0 / 5	0.7660070153605
44	[(Craft, Labor, Missing), (Agriculture,	[(1, 2, 3), (0, 4, 5,	0.7669970153695
45	Professional, Sales, Services)]	6)]	415
45	[(Craft, Labor, Professional), (Agriculture,	[(1, 2, 4), (0, 3, 5,	0.7667090459089
4.6	Missing, Sales, Services)]	6)]	398
46	[(Craft, Labor, Sales), (Agriculture, Missing,	[(1, 2, 5), (0, 3, 4,	0.7668155372116
	Professional, Services)]	6)]	556
47	[(Craft, Labor, Services), (Agriculture,	[(1, 2, 6), (0, 3, 4,	0.7669972882930
	Missing, Professional, Sales)]	5)]	872
	1/2 6 20 2 2 6 2 10 10 10 10	<u> </u>	
48	[(Craft, Missing, Professional), (Agriculture,	[(1, 3, 4), (0, 2, 5,	0.7666819099761
	Labor, Sales, Services)]	6)]	641
49	[(Craft, Missing, Sales), (Agriculture, Labor,	[(1, 3, 5), (0, 2, 4,	0.7668173702917
	Professional, Services)]	6)]	922
50	[(Craft, Missing, Services), (Agriculture,	[(1, 3, 6), (0, 2, 4,	0.7670535229559
	Labor, Professional, Sales)]	5)]	383
51	[(Craft, Professional, Sales), (Agriculture,	[(1, 4, 5), (0, 2, 3,	0.7670235232160
	Labor, Missing, Services)]	6)]	99
52	[(Craft, Professional, Services), (Agriculture,	[(1, 4, 6), (0, 2, 3,	0.7667760989651
	Labor, Missing, Sales)]	5)]	224
53	[(Craft, Sales, Services), (Agriculture, Labor,	[(1, 5, 6), (0, 2, 3,	0.7668008442520
	Missing, Professional)]	4)]	571
	73	,,	
54	[(Labor, Missing, Professional), (Agriculture,	[(2, 3, 4), (0, 1, 5,	0.7668099526557
	Craft, Sales, Services)]	6)]	102
L	J. J. J. Gales, Sci 11003/j	~/1	

55	[(Labor, Missing, Sales), (Agriculture, Craft,	[(2, 3, 5), (0, 1, 4,	0.7667638050721
	Professional, Services)]	6)]	957
56	[(Labor, Missing, Services), (Agriculture,	[(2, 3, 6), (0, 1, 4,	0.7670188726283
	Craft, Professional, Sales)]	5)]	493
57	[(Labor, Professional, Sales), (Agriculture,	[(2, 4, 5), (0, 1, 3,	0.7670517628446
	Craft, Missing, Services)]	6)]	133
58	[(Labor, Professional, Services),	[(2, 4, 6), (0, 1, 3,	0.7668304838752
	(Agriculture, Craft, Missing, Sales)]	5)]	236
59	[(Labor, Sales, Services), (Agriculture, Craft,	[(2, 5, 6), (0, 1, 3,	0.7666785924047
	Missing, Professional)]	4)]	138
60	[(Missing, Professional, Sales), (Agriculture,	[(3, 4, 5), (0, 1, 2,	0.7670039804995
	Craft, Labor, Services)]	6)]	727
61	[(Missing, Professional, Services),	[(3, 4, 6), (0, 1, 2,	0.7668233960538
	(Agriculture, Craft, Labor, Sales)]	5)]	839
62	[(Missing, Sales, Services), (Agriculture,	[(3, 5, 6), (0, 1, 2,	0.7666943383195
	Craft, Labor, Professional)]	4)]	367
63	[(Professional, Sales, Services), (Agriculture,	[(4, 5, 6), (0, 1, 2,	0.7670006273964
	Craft, Labor, Missing)]	3)]	458

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

Optimal Split		GINI for split
[(Labor, Sales, Services), (Agriculture, Craft,	[(2, 5, 6), (0, 1, 3,	0.7666785924047
Missing, Professional)]	4)]	138

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

I will choose CarOwnership as the predictor for the second layer as the GINI for CarOwnership is lesser than the GINI for JobCategory.

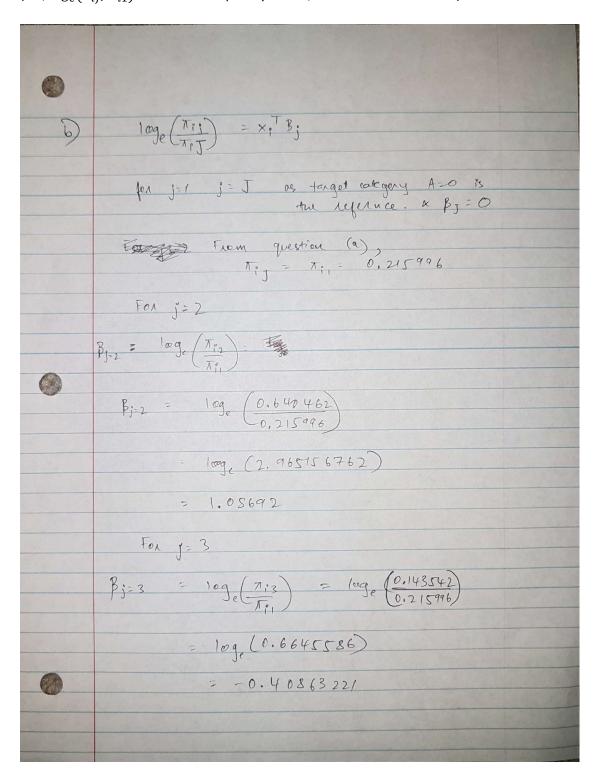
Question 2

a) (10 points) Suppose you start with a model with only the Intercept term (i.e., without any predictors). You are asked to mathematically calculate the maximum likelihood estimates of the predicted probabilities π_{ij} , j=1,2,3 without calling the MNLogit function. Show all the necessary steps and the estimates for the π_{ij} , j=1,2,3. (Hint: equate the first derivatives of the log-likelihood function to zeros for this Intercept-only model)

	Figureting the first derivative of the log-likehood function to zero.
	$\leq x_{is}(n_{ij} - n_{i}\pi_{i}) = 0$
	Since there are no subpopulations, m-1
	xis(ni, - ni Tij) = 0
	uij - uitij = O
9	min
	No. of observations A 143 691 (N1) p 4 26067 (N1) 1 9 5 491 (N1s) 2 3 abservations: 665249
	Total
	T: 143691 = 0.215996
0	Ti2 = 426067 = 0.640462
	Ti3 = 95 491 = 0.1435 42 665249

	0
lace likelihered K	
21- 8 & N. 1000 (X.)	
log likelihood, m k -> l= 2 2 N; lege (Tig)	
= 143891 x log (0.215996) +	
426067 x log (0.640462)+	
= 143891 x leq (0.215996) + 426067 x leq (0.640462) + 95 491 x leoq. (0.143542)	
= -595406.7619	
log L-t	
log L= l	
log L = 1 log L = -595406, 7619	N. H.
Je 50 C406, 7319	
L= 2 = 1000 Moximum	
Like hood.	
The second secon	
The town of the to	
The state of the s	
A THE CAST OF THE PARTY OF THE	
CAPACES CAPACES	
2200 C200	

b) (10 points) Next, you are asked to mathematically calculate the maximum likelihood estimates of the Intercept terms β_{j0} , j=1,...,K. 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)?



- c) (4 points) Now, you will use the MNLogit function to build the multinomial logistic model. What value of the target variable A is used by the MNLogit function as the reference category? The statsmodels.api.MNLogit function conventionally takes the lexically first target category as the reference. In this case, A = 0 is the reference category.
- d) (2 points) How many iterations are performed before convergence is achieved?
 57 Iterations were performed before convergence.
 With default settings, the algorithm did not converge. Hence the maximum number of iterations was set to 100.

```
thisFit = logit.fit(method='newton', full_output = True, maxiter = 100, tol = 1e-8)
```

- e) (4 points) How many parameters (including the redundant ones) are in the model? There are 7 parameters.
- f) (5 points) When group_size = 2, homeowner = 1, and married_couple = 1, what are the predicted probabilities: Prob(A = 0), Prob(A = 1), and Prob(A = 2)?

[[0.16697404 0.69415127 0.13887469]]

Prob(A = 0) = 0.16697404Prob(A = 1) = 0.69415127

Prob(A = 2) = 0.13887469

g) (10 points) What are the values of the predictors group_size, homeowner, and married_couple such that Prob(A = 0) will attain its maximum? What is the maximum Prob(A = 0) value? [1,0,0,0,1,0,0]

The values of the predictors are:

group_size = 4 homeowner = 0 married_couple = 0 Maximum Prob(A = 0) value = 0.42117104

h) (5 points) According to the logistic model, what is the odds ratio for group_size = 4 versus group_size = 1, and A = 1 versus A = 0? Mathematically, the odds ratio is (Prob(A=1)/Prob(A=0) | group_size = 4) / ((Prob(A=1)/Prob(A=0) | group_size = 1).

```
For group_size , 4 = 0.408236 - 0.359197 = 0.049039
For group_size, 1 = 0.408236 + 0.424367 = 0.832603
```

Group_size 4 - group_size 1 = -0.783564

 $log_e (Prob(A=1)/Prob(A=0) \mid group_size = 4) - log_e (Prob(A=1)/Prob(A=0) \mid group_size = 1) = -0.783564$

 $log_e \{ (Prob(A=1)/Prob(A=0) \mid group_size = 4) / (Prob(A=1)/Prob(A=0) \mid group_size = 1) \} = 1$

-0.783564

 $(Prob(A=1)/Prob(A=0) \mid group_size = 4) / (Prob(A=1)/Prob(A=0) \mid group_size = 1) = e^{0.783564} = 0.45678$