

# CS 584-04: Machine Learning

## Autumn 2019 Assignment 4

---

### Question 1 (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. You will build a multinomial logistic model with the following specifications.

1. The nominal target variable is **A** which have these categories 0, 1, and 2
2. The nominal features are (categories are inside the parentheses):
  - a. **group\_size**. How many people will be covered under the policy (1, 2, 3 or 4)?
  - b. **homeowner**. Whether the customer owns a home or not (0 = No, 1 = Yes)?
  - c. **married\_couple**. Does the customer group contain a married couple (0 = No, 1 = Yes)?
3. Include the Intercept term in the model
4. Enter the five model effects in this order: `group_size`, `homeowner`, `married_couple`, `group_size * homeowner`, and `homeowner * married_couple` (No forward or backward selection)
5. The optimization method is Newton
6. The maximum number of iterations is 100
7. The tolerance level is  $1e-8$ .
8. Use the `sympy.Matrix().rref()` method to identify the non-aliased parameters

Please answer the following questions based on your model.

- a) (5 points) List the aliased parameters that you found in your model.

**Solution:-**

group_size_4
homeowner_1
married_couple_1
group_size_1 * home_owner_1
group_size_2 * home_owner_1
group_size_3 * homeowner_1
group_size_4 * homeowner_0
group_size_4 * homeowner_1
homeowner_0 * married_couple_1

homeowner_1 * married_couple_0
homeowner_1 * married_couple_1

- b) (5 points) How many degrees of freedom do you have in your model?

**Solution:-** Degrees of freedom = **20**

- c) (10 points) After entering a model effect, calculate the Deviance test statistic, its degrees of freedom, and its significance value between the current model and the previous model. List your Deviance test results by the model effects in a table.

**Solution:-**

Model effect	Chi-square statistic	Degree of freedom	Significance
group_size	987.5766005262267	6	4.347870389027117e-210
homeowner	5867.781500353478	2	0.0
married_couple	84.57800238393247	2	4.3064572180356084e-19
group_size * homeowner	254.07812536344863	6	5.5121059685664295e-52
homeowner * married_couple	70.84227676945738	2	4.13804354793157e-16

- d) (5 points) Calculate the Feature Importance Index as the negative base-10 logarithm of the significance value. List your indices by the model effects.

**Solution:-**

Model Effect	Feature Importance Index
group_size	209.36172341080683
homeowner	Infinity
married_couple	18.365879862870976
group_size * homeowner	51.2586824418404
homeowner * married_couple	15.383204943219134

- e) (10 points) For each of the sixteen possible value combinations of the three features, calculate the predicted probabilities for  $A = 0, 1, 2$  based on the multinomial logistic model. List your answers in a table with proper labelling.

**Solution:-**

Possibilities	Pr(A=0)	Pr(A=1)	Pr(A=2)
group_size_1_homeowner_0_married_couple_0	0.2596505294674 657	0.58917499922 24112	0.1511744713101 232
group_size_1_homeowner_0_married_couple_1	0.2600917257573 4996	0.59210579575 04566	0.1478024784921 9338
group_size_1_homeowner_1_married_couple_0	0.1836024871259 8187	0.68202954795 90752	0.1343679649149 4285
group_size_1_homeowner_1_married_couple_1	0.1540230016390 0653	0.70991797504 80099	0.1360590233129 8358
group_size_2_homeowner_0_married_couple_0	0.2219361377687 4224	0.62110513967 97336	0.1569587225515 2415
group_size_2_homeowner_0_married_couple_1	0.2223208689711 648	0.62421616161 4875	0.1534629694139 6015
group_size_2_homeowner_1_married_couple_0	0.2025095494428 9226	0.65977268413 50491	0.1377177664220 5874
group_size_2_homeowner_1_married_couple_1	0.1705515551369 7823	0.68944950933 34505	0.1399989355295 7146
group_size_3_homeowner_0_married_couple_0	0.2395700821998 8554	0.60461592070 4135	0.1558139970959 7948
group_size_3_homeowner_0_married_couple_1	0.2399917487557 1322	0.60766047069 60655	0.1523477805482 2107
group_size_3_homeowner_1_married_couple_0	0.3011398506574 032	0.53129677962 25498	0.1675633697200 4698
group_size_3_homeowner_1_married_couple_1	0.2590173731094 2163	0.56701664362 09104	0.1739659832696 6792
group_size_4_homeowner_0_married_couple_0	0.1944846874453 793	0.66968590537 77223	0.1358294071768 984
group_size_4_homeowner_0_married_couple_1	0.1946921010052 6804	0.67259209013 01082	0.1327158088646 2372
group_size_4_homeowner_1_married_couple_0	0.3877190926982 077	0.48497444674 97125	0.1273064605520 7973
group_size_4_homeowner_1_married_couple_1	0.3391717120963 9435	0.52640408607 88882	0.1344242018247 1734

- f) (5 points) Based on your model, what values of group\_size, homeowner, and married\_couple will maximize the odds value  $\text{Prob}(A=1) / \text{Prob}(A=0)$ ? What is that maximum odd value?

**Solution:-** group\_size = 1, homeowner = 1, married\_couple = 1

Maximum odd value = **4.609168549460486**

- g) (5 points) Based on your 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  $(\text{Prob}(A=2)/\text{Prob}(A=0) \mid \text{group\_size} = 3) / ((\text{Prob}(A=2)/\text{Prob}(A=0) \mid \text{group\_size} = 1))$ .

**Solution:-**

$$\begin{aligned} & (\text{Prob}(A=2)/\text{Prob}(A=0) \mid \text{group\_size} = 3) / ((\text{Prob}(A=2)/\text{Prob}(A=0) \mid \text{group\_size} = 1)) \Rightarrow \\ & \log((\text{Prob}(A=2)/\text{Prob}(A=0) \mid \text{group\_size} = 3) / ((\text{Prob}(A=2)/\text{Prob}(A=0) \mid \text{group\_size} = 1))) = \\ & \log(\text{Prob}(A=2)/\text{Prob}(A=0) \mid \text{group\_size} = 3) - \log(\text{Prob}(A=2)/\text{Prob}(A=0) \mid \text{group\_size} = 1) \end{aligned}$$

$$\rightarrow (\text{group\_size}_3 - \text{group\_size}_1) + (\text{group\_size}_3 * \text{homeowner}_0 - \text{group\_size}_1 * \text{homeowner}_0)(1-h)$$

$$\rightarrow (-0.274022 + 0.384741) * (1-h)$$

$$(\text{Prob}(A=2)/\text{Prob}(A=0) \mid \text{group\_size} = 3) / ((\text{Prob}(A=2)/\text{Prob}(A=0) \mid \text{group\_size} = 1))$$

$$\Rightarrow \exp(0.527471 - 0.801493) + (-0.5987 + 0.983441)(1-h)$$

$$\Rightarrow \exp((-0.274022 + 0.384741) * (1-h))$$

Here the odds ratio depends on the values of group\_size and homeowner. So h take values of ( 0 or 1) where h is homeowner.

- h) (5 points) Based on your model, what is the odds ratio for homeowner = 1 versus homeowner = 0, and A = 0 versus A = 1? Mathematically, the odds ratio is  $(\text{Prob}(A=0)/\text{Prob}(A=1) \mid \text{homeowner} = 1) / ((\text{Prob}(A=0)/\text{Prob}(A=1) \mid \text{homeowner} = 0))$ .

**Solution:-**  $(\text{Prob}(A=0)/\text{Prob}(A=1) \mid \text{homeowner} = 1) / ((\text{Prob}(A=0)/\text{Prob}(A=1) \mid \text{homeowner} = 0)) \Rightarrow$   
 $\log(\text{Prob}(A=0)/\text{Prob}(A=1) \mid \text{homeowner} = 1) - \log(\text{Prob}(A=0)/\text{Prob}(A=1) \mid \text{homeowner} = 0) =$   
 $(0.800157 - 1.505554 * g_1 - 1.164638 * g_2 - 0.654639 * g_3 + 0.212483 * (1-m)).$

$$(\text{Prob}(A=0)/\text{Prob}(A=1) \mid \text{homeowner} = 1) / ((\text{Prob}(A=0)/\text{Prob}(A=1) \mid \text{homeowner} = 0))$$

$$\rightarrow \exp((0.800157 - 1.505554 * g_1 - 1.164638 * g_2 - 0.654639 * g_3 + 0.212483 * (1-m))).$$

Here the odds ratio depends on the values of group\_size and married\_couple.

So  $g_1, g_2, g_3, m$  take values of ( 0 or 1) where  $g_1, g_2, g_3$  is group\_size\_1, group\_size\_2, group\_size\_3 respectively. And m is married\_couple.

## Question 2 (50 points)

You are asked to build a Naïve Bayes model using the same Purchase\_Likelihood.csv. The model specifications are:

1. No smoothing is needed. Therefore, the Laplace/Lidstone alpha is zero
2. The nominal target variable is **A** which have these categories 0, 1, and 2
3. The nominal features are (categories are inside the parentheses):
  - a. **group\_size**. How many people will be covered under the policy (1, 2, 3 or 4)?
  - b. **homeowner**. Whether the customer owns a home or not (0 = No, 1 = Yes)?
  - c. **married\_couple**. Does the customer group contain a married couple (0 = No, 1 = Yes)?

Please answer the following questions based on your model.

- a) (5 points) Show in a table the frequency counts and the Class Probabilities of the target variable.

**Solution:-**

	A	Count	Proportion
0	0	143691	0.215996
1	1	426067	0.640462
2	2	95491	0.143542

- b) (5 points) Show the crosstabulation table of the target variable by the feature group\_size. The table contains the frequency counts.

**Solution:-**

group_size				
A	1	2	3	4
0	115460	25728	2282	221
1	329552	91065	5069	381
2	74293	19600	1505	93

- c) (5 points) Show the crosstabulation table of the target variable by the feature homeowner. The table contains the frequency counts.

**Solution:-**

homeowner		
A	0	1
0	78659	65032
1	183130	242937
2	46734	48757

- d) (5 points) Show the crosstabulation table of the target variable by the feature married\_couple. The table contains the frequency counts.

**Solution:-**

married_couple		
A	0	1
0	117110	26581
1	333272	92795
2	75310	20181

- e) (10 points) Calculate the Cramer's V statistics for the above three crosstabulations tables. Based on these Cramer's V statistics, which feature has the largest association with the target A?

**Solution:-**

Cramer's V Value for group\_size is: **0.027102014055820786**

Cramer's V Value for homeowner is: **0.09708641964781962**

Cramer's V Value for married\_couple is: **0.03242164583520746**

**"homeowner"** has the largest association with the target A.

- f) (5 points) Based on the assumptions of the Naïve Bayes model, express the joint probability  $\text{Prob}(A = a, \text{group\_size} = g, \text{homeowner} = h, \text{married\_couple} = m)$  as a product of the appropriate probabilities.

**Solution:-**

$\text{Prob}(A=a, \text{group\_size}=g, \text{homeowner}=h, \text{married\_couple}=m) =$

$\text{Prob}(\text{group\_size}=g, \text{homeowner}=h, \text{married\_couple}=m \mid A=a) * \text{Prob}(A=a) =$

**$\text{Prob}(\text{group\_size}=g \mid A=a) * \text{Prob}(\text{homeowner}=h \mid A=a) * \text{Prob}(\text{married\_couple}=m \mid A=a) * \text{Prob}(A=a)$**

- g) (10 points) For each of the sixteen possible value combinations of the three features, calculate the predicted probabilities for  $A = 0, 1, 2$  based on the Naïve Bayes model. List your answers in a table with proper labelling.

**Solution:-**

Possibilities	Pr(A=0)	Pr(A=1)	Pr(A=2)
group_size_1_homeowner_0_married_couple_0	0.269721900 83648967	0.580133399 3691891	0.150144699 79432118
group_size_1_homeowner_0_married_couple_1	0.232789218 51630957	0.614218557 8024016	0.152992223 68128876
group_size_1_homeowner_1_married_couple_0	0.194037904 75559898	0.669659004 8821739	0.136303090 3622272
group_size_1_homeowner_1_married_couple_1	0.164935004 743777	0.698278045 9509148	0.136786949 30530805
group_size_2_homeowner_0_married_couple_0	0.231143327 3249531	0.616518459 7447714	0.152338212 93027552
group_size_2_homeowner_0_married_couple_1	0.198015591 405003	0.647906780 7659843	0.154077627 82901277
group_size_2_homeowner_1_married_couple_0	0.163627525 52123652	0.700287808 8359464	0.136084665 64281702
group_size_2_homeowner_1_married_couple_1	0.138274170 44457968	0.725954963 0220522	0.135770866 53336812
group_size_3_homeowner_0_married_couple_0	0.308219393 78427693	0.515924167 7311622	0.175856438 48456095
group_size_3_homeowner_0_married_couple_1	0.268311057 11605896	0.550950897 1155715	0.180738045 76836952
group_size_3_homeowner_1_married_couple_0	0.226971831 46374494	0.609611781 1433283	0.163416387 39292683
group_size_3_homeowner_1_married_couple_1	0.194369513 62831584	0.640409773 5081213	0.165220712 86356266
group_size_4_homeowner_0_married_couple_0	0.375490390 7259939	0.487810100 5336526	0.136699508 74035344
group_size_4_homeowner_0_married_couple_1	0.330743444 1365481	0.527098304 946624	0.142158250 91682782
group_size_4_homeowner_1_married_couple_0	0.282172679 6029393	0.588196454 8622688	0.129630865 5347919
group_size_4_homeowner_1_married_couple_1	0.243930339 20041854	0.623765964 2682374	0.132303696 53134402

- h) (5 points) Based on your model, what values of group\_size, homeowner, and married\_couple will maximize the odds value  $\text{Prob}(A=1) / \text{Prob}(A=0)$ ? What is that maximum odd value?

**Solution:-**

The maximum value occurs when group\_size = 2, homeowner = 1, married\_couple = 1.

The maximum value is: **5.250112589270714**