

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.

1. It contains 665,249 observations on 97,009 unique Customer ID.
2. The nominal target variable is **insurance** which has 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)?

Question 1 (35 points)

You will build a multinomial logistic model with the following model specifications.

1. Enter the six effects to the model in this sequence:
 - a. `group_size`
 - b. `homeowner`
 - c. `married_couple`
 - d. `group_size * homeowner`
 - e. `group_size * married_couple`
 - f. `homeowner * married_couple`
2. Include the Intercept term in the model
3. The optimization method is Newton
4. The maximum number of iterations is 100
5. The tolerance level is $1e-8$.
6. 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 columns that you found in your model matrix.

The aliased parameters found in model are

```
group_size_4  
homeowner_1  
married_couple_1  
group_size_1 * homeowner_1
```

group_size_2 * homeowner_1
 group_size_3 * homeowner_1
 group_size_4 * homeowner_0
 group_size_4 * homeowner_1
 group_size_1 * married_couple_1
 group_size_2 * married_couple_1
 group_size_3 * married_couple_1
 group_size_4 * married_couple_0
 group_size_4 * married_couple_1
 homeowner_0 * married_couple_1
 homeowner_1 * married_couple_0
 homeowner_1 * married_couple_1

b) (5 points) How many degrees of freedom does your model have?

Degrees of freedom of final model - 2

c) (20 points) After entering each 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.

d)

Step	Effect Entered	# Free Parameter	Log-Likelihood	Deviance	Degrees of Freedom	Significance
0	Intercept	2	-595406.76188 44224	Not Applicable		
1	group_size	8	-594912.97358 41593	987.57660052 62267	6	4.34787038902711 7e-210
2	homeowner	10	-591979.08283 39825	5867.7815003 53478	2	0.0
3	married_couple	12	-591936.79383 27907	84.578002383 69964	2	4.30645721853695 87e-19
4	group_size * homeowner	18	-591809.75477 01088	254.07812536 368147	6	5.51210596793472 1e-52
5	group_size * married_couple	24	-591105.49317 71928	1382.5423636 827618	6	1.45970012121037 11e-295

Step	Effect Entered	# Free Parameter	Log-Likelihood	Deviance	Degrees of Freedom	Significance
6	homeowner * married_couple	26	-591105.4931771928	25.980822149431333	2	2.28210778553294e-06

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

Effect Entered	Importance
Intercept	Not Applicable
group_size	209.36172341080683
homeowner	inf
married_couple	18.365879862820417
group_size * homeowner	51.25868244189017
group_size * married_couple	294.83573635591443
homeowner * married_couple	5.641663847403906

Question 2 (25 points)

Please answer the following questions based on your multinomial logistic model in Question 1.

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

group_size	homeowner	married_couple	Prob(insurance = 0)	Prob(insurance = 1)	Prob(insurance = 2)
1	0	0	0.257582	0.591653	0.150765
1	0	1	0.328060	0.510687	0.161253
1	1	0	0.180464	0.686085	0.133452
1	1	1	0.217257	0.628228	0.154515
2	0	0	0.279425	0.550953	0.169623
2	0	1	0.203284	0.647446	0.149269
2	1	0	0.249383	0.597778	0.152838
2	1	1	0.161437	0.701504	0.137059

group_size	homeowner	married_couple	Prob(insurance = 0)	Prob(insurance = 1)	Prob(insurance = 2)
3	0	0	0.237434	0.654601	0.107965
3	0	1	0.240406	0.597961	0.161632
3	1	0	0.282651	0.603586	0.113763
3	1	1	0.260167	0.562521	0.177312
4	0	0	0.304008	0.595211	0.100781
4	0	1	0.193714	0.673257	0.133029
4	1	0	0.505939	0.406206	0.087855
4	1	1	0.332066	0.531139	0.136796

- b) (5 points) Based on your answers in (a), what value combination of group_size, homeowner, and married_couple will maximize the odds value $\text{Prob}(\text{insurance} = 1) / \text{Prob}(\text{insurance} = 0)$? What is that maximum odd value?

Max odd value = 4.345371

- c) (5 points) Based on your model, what is the odds ratio for group_size = 3 versus group_size = 1, and insurance = 2 versus insurance = 0?
(Hint: The odds ratio is this odds ($\text{Prob}(\text{insurance} = 2) / \text{Prob}(\text{insurance} = 0) \mid \text{group_size} = 3$) divided by this odds ($(\text{Prob}(\text{insurance} = 2) / \text{Prob}(\text{insurance} = 0) \mid \text{group_size} = 1)$.)
 This value depends on homeowner

Homeowner = 0

$$\begin{aligned}
 &= (\text{Prob}(\text{insurance}=1)/\text{Prob}(\text{insurance}=0) \mid \text{group_size} = 3) / \\
 &((\text{Prob}(\text{insurance}=2)/\text{Prob}(\text{insurance}=0) \mid \text{group_size} = 1) \\
 &= (0.45471635 / 0.58530902) \\
 &= 0.776882532
 \end{aligned}$$

Homeowner = 1

$$\begin{aligned}
 &= (\text{Prob}(\text{insurance}=1)/\text{Prob}(\text{insurance}=0) \mid \text{group_size} = 3) / \\
 &((\text{Prob}(\text{insurance}=2)/\text{Prob}(\text{insurance}=0) \mid \text{group_size} = 1) \\
 &= 0.40248696272873963 / 0.5853090230224534 \\
 &= 0.687648655492023
 \end{aligned}$$

- d) (5 points) Based on your model, what is the odds ratio for homeowner = 1 versus homeowner = 0, and insurance = 0 versus insurance = 1?

This value depends on group size and married couple

Group size = 1 and married couple = 0

$$\begin{aligned}
 &((\text{Prob}(\text{insurance}=0)/\text{Prob}(\text{insurance}=1) \mid \text{homeowner} = 1) / \\
 &((\text{Prob}(\text{insurance}=0)/\text{Prob}(\text{insurance}=1) \mid \text{homeowner} = 0) \\
 &= 0.2630339724 / 0.43536032519
 \end{aligned}$$

=0.6041753399

Group size = 1 and married couple = 1

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

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

=0.3458253503422945 / 0.6423882982223854

=0.5383431661181581

Group size = 2 and married couple = 0

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

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

= 0.4171834144139587 / 0.5071660097898533

=0.822577630127107

Similarly, it is calculated for other values of group size and married couple.

Question 3 (40 points)

You will build a Naïve Bayes model without any smoothing. In other words, the Laplace/Lidstone alpha is zero. 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.

insurance	0	1	2
Frequency Count	143691	426067	95491
Class Probability	0.215996	0.640462	0.143542

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

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

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

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

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

Married_couple	insurance		
	0	1	2
0	117110	333272	75310
1	26581	92795	20181

- e) (5 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 insurance?

Group_size = 0.027102014055820786

Homeowner = 0.09708641964781962

Married_couple = 0.03242164583520746

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

group_size	homeowner	married_couple	Prob(insurance = 0)	Prob(insurance = 1)	Prob(insurance = 2)
1	0	0	0.269722	0.580133	0.150145
1	0	1	0.232789	0.614219	0.152992
1	1	0	0.194038	0.669659	0.136303
1	1	1	0.164935	0.698278	0.136303
2	0	0	0.231143	0.616518	0.152338
2	0	1	0.198016	0.647907	0.154078
2	1	0	0.163628	0.700288	0.136085
2	1	1	0.138274	0.725955	0.135771
3	0	0	0.308219	0.515924	0.175856
3	0	1	0.268311	0.550951	0.180738
3	1	0	0.226972	0.609612	0.163416
3	1	1	0.194370	0.640410	0.165221
4	0	0	0.375490	0.487810	0.136700
4	0	1	0.330743	0.527098	0.142158
4	1	0	0.282173	0.588196	0.129631
4	1	1	0.243930	0.623766	0.132304

- g) (5 points) Based on your model, what value combination of group_size, homeowner, and married_couple will maximize the odds value $\text{Prob}(\text{insurance} = 1) / \text{Prob}(\text{insurance} = 0)$? What is that maximum odd value?

Max odd value - 5.250112589270714

(2, 1, 1) - group_size = 2 homeowner = 1 married_couple = 1

