

SAA for many N level multilevel models

Welcome to an SAA for fitting many model types developed for Stat-JR v1.0.5

Input questions

Firstly on this page you will need to specify the dataset required from the list of available datasets.

Which dataset do you wish to use:

Submit

Next you need to choose many options including the response, estimation method, clustering variables and predictor variables (both continuous and categorical) from the chosen dataset. After choosing these variables the SAA will run and you will see a block of text describing how many observations are to be used at the bottom of this page. The rest of the analysis will appear in pages 2-12.

What estimation method do you want to use:

IGLS

What is the response variable:

use

What distribution are you going to assume:

Binomial

Which column contains the denominators:

cons

What link function do you wish to use:

logit

Please enter your possible (nested) classifications / levels (lowest first, not including level-1):	cons
Are there any continuous predictors that need including in all models:	No
Are there any categorical predictors that need including in all models:	No
Do you want to include any continuous predictors as candidates for inclusion in the models:	Yes
Which continuous predictors do you want to consider:	age
Do you want to include any categorical predictors as candidates for inclusion in the models:	Yes
Which categorical predictors do you want to consider:	lc
What selection type do you require:	Forward pass
Do you want to test for random slopes:	No
Do you want to test for interactions:	No

The Analysis Assistant you are currently using is designed to work on complete datasets only and so as a pre-processing step we have to remove any rows that contain missing data in columns used in the analysis that follows. For now the list

of columns to be considered is: use, cons, cons, age, lc. There are 0 (0.0%) rows that get deleted This results in a dataset of 1934 rows.

On the next page we will look at the shape of the response and, in the case of normal responses, decide whether to log transform.

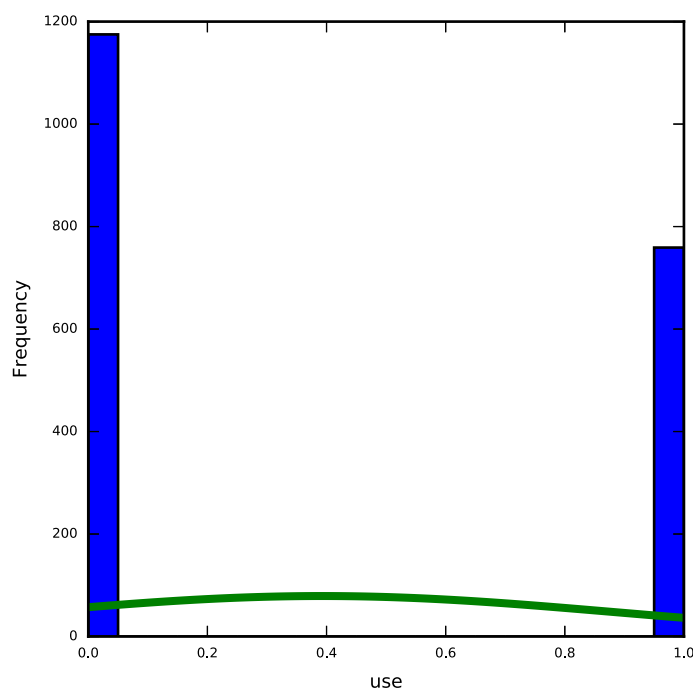
Exploring the response

We will begin our analysis of the dataset by doing some basic data exploration.

You have chosen use as your response variable and so a first step is to take a look at this variable and assess its suitability for modelling. The summary statistics for the variable are in the table below:

	Observations	1934
	Mean	0.392
	Standard Deviation	0.488
	Median	0.0

We also look at a histogram of use to see what it looks like - noting that for a Binomial model this is of less interest as it will simply look like a bar graph.



Here the median is smaller than the mean and there is significant skew to the right. The skewness value is 0.441. Here the statistical significance may be to some degree due to the large sample size as from a practical perspective values of skew less than 2 are not considered too big a skew.

There are no obvious outliers in use.

Exploring the predictors individually

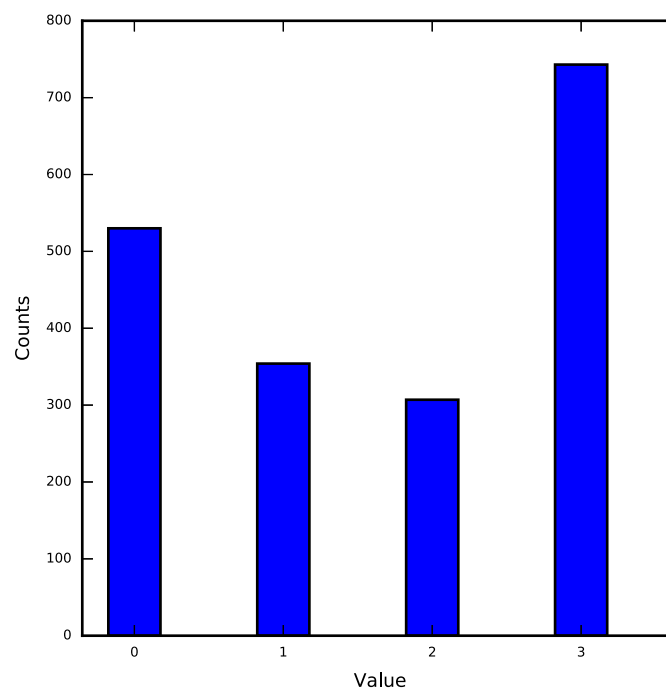
We can also look at each of the predictor variables in turn in isolation.

For categorical predictors we are looking at how common each category is in the dataset. In particular we are checking for rare categories which might cause difficulties in modelling and might therefore be usefully merged with other categories (though this would need to be done outside this SAA).

For predictor `lc` we see the following:

lc	N	Percentage
0	530	27.404
1	354	18.304
2	307	15.874
3	743	38.418
Total	1934	100

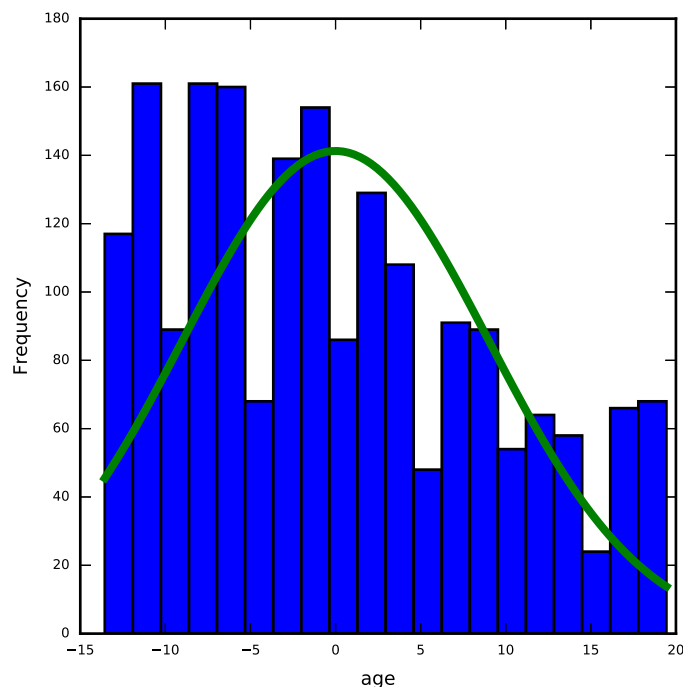
None of the categories of `lc` have fewer than 5 observations.



For continuous predictors we are interested in looking at summary statistics, the shape of the distribution and any unusual values. If the distribution is skewed then we might want to transform the variable before fitting it in the model although it is more important to consider transformations of the response variable and remember what is important is whether the relationship between the response and predictor is linear. If there are unusual values we will want to check that the unusual values are correct and not errors and also whether we may want to treat the variable differently. Another possibility for unusual shaped distributions is to instead categorise the variable into ranges of values.

For predictor age we see the following:

Name	age
Observations	1934
Mean	0.002
Standard Deviation	9.011
Median	-1.56



Here the median is smaller than the mean and there is significant skew to the right. The skewness value is 0.441. Here the statistical significance may be to some degree due to the large sample size as from a practical perspective values of skew less than 2 in absolute magnitude are not considered too big a skew.

There are no obvious outliers in age.

Assessing the relationship between the response and individual predictors

Once we are happy with our response variable and our set of predictors we now want to have a preliminary look at them together before progressing to the univariable modelling.

For the categorical predictors it is worth tabulating the response for each category to look at whether patterns differ. We can formally test this with a chi-squared test.

We will investigate categorical variable *lc*. To do a chi-squared test we start by tabulated observed counts and totals:

Observed	use=0	use=1	Total
lc=0	397	133	530
lc=1	190	164	354
lc=2	160	147	307
lc=3	428	315	743
Total	1175	759	1934

We can therefore work out the expected counts from the margins of the observed data.

And so we expect

$$E(\text{use} = 0, \text{lc} = 0) = \text{Total use} = 0 * \text{Total lc} = 0 / \text{grand total} = 1175 * 530 / 1934 = 322.0.$$

$$E(\text{use} = 1, \text{lc} = 0) = \text{Total use} = 1 * \text{Total lc} = 0 / \text{grand total} = 759 * 530 / 1934 = 208.0.$$

$$E(\text{use} = 0, \text{lc} = 1) = \text{Total use} = 0 * \text{Total lc} = 1 / \text{grand total} = 1175 * 354 / 1934 = 215.07.$$

$$E(\text{use} = 1, \text{lc} = 1) = \text{Total use} = 1 * \text{Total lc} = 1 / \text{grand total} = 759 * 354 / 1934 = 138.93.$$

$$E(\text{use} = 0, \text{lc} = 2) = \text{Total use} = 0 * \text{Total lc} = 2 / \text{grand total} = 1175 * 307 / 1934 = 186.52.$$

$$E(\text{use} = 1, \text{lc} = 2) = \text{Total use} = 1 * \text{Total lc} = 2 / \text{grand total} = 759 * 307 / 1934 = 120.48.$$

$$E(\text{use} = 0, \text{lc} = 3) = \text{Total use} = 0 * \text{Total lc} = 3 / \text{grand total} = 1175 * 743 / 1934 =$$

451.41.

$E(\text{use} = 1, \text{lc} = 3) = \text{Total use} = 1 * \text{Total lc} = 3 / \text{grand total} = 759 * 743 / 1934 = 291.59.$

So the table of expected counts is:

Expected	use=0	use=1	Total
lc=0	322.0	208.0	530.0
lc=1	215.07	138.93	354.0
lc=2	186.52	120.48	307.0
lc=3	451.41	291.59	743.0
Total	1175.0	759.0	1934.0

We next look at differences between what we observe and expect in each cell. We square these values so that every difference is positive and scale by the expected counts so that more frequently expected cells aren't overly influential. So for example for use=0, lc=0 $(O-E)^2/E = (397-322.0)^2/322.0=17.47$. This statistic is shown in tabular form below:

(O-E)^2/E	use=0	use=1
lc=0	17.47	27.04
lc=1	2.92	4.52
lc=2	3.77	5.84
lc=3	1.21	1.88

The test statistic for a chi-squared test is found by summing the values of this table so:

$\text{Chisq}=17.47+27.04+2.92+4.52+3.77+5.84+1.21+1.88=64.66.$

This is compared with a chi-squared table with degrees of freedom = (number of columns -1)x(number of rows - 1) =

$(4-1)x(2-1)=3.$

Looking up the chi-squared table the value for p=0.05 is 7.81 and for p=0.01 = 11.34

As $64.66 > 11.34$ our p value is less than 0.01 and we have strong evidence to reject the null hypothesis (at the $p=0.01$) level.

The p-value is in fact less than 0.0001.

For the continuous predictors it is worth looking at the mean value of each predictor for the 0 and 1 responses to assess if there is any difference. We can formally test this with a t-test.

Here is a tabulation of the predictor, age for response use with category 1 having the largest mean and category 0 the smallest.

Category	N	Mean	Standard Deviation	Median
0	1175	-0.208	9.707	-1.56
1	759	0.327	7.802	-0.56

The formal test is as follows:

There are two groups in the data:

The first group has 1175 observations with mean -0.208 standard deviation 9.711.

The second group has 759 observations with mean 0.327 standard deviation 7.807.

We are trying to test a hypothesis as to whether the two groups differ in their (population) means by a statistically significant amount. Statistical significance is related to how likely a result is to be a chance occurrence. Here we are trying to differentiate between a real difference (no matter how small) and a difference that may have occurred due to the samples we have chosen.

The mean difference is 0.534 with the second group having the larger sample mean.

We need to quantify if this difference is large relative to the variability in the data. To do this we calculate the standard error of the difference. This is a function of the variabilities in the samples from group A and group B combined with their sample sizes. The bigger the 2 variabilities the larger the standard error, whilst the smaller the variability the smaller the standard error.

For our data the standard error of the mean difference is 0.401 and we divide our observed difference by this standard error to give a test statistic with value 1.334.

This test statistic is then compared to a t distribution with degrees of freedom equal to the sum of the sample sizes in each group $(1934) - 2$. In this case a t distribution with 1932. This t table has values of 1.961 for $p=0.05$ and 2.578 for $p=0.01$.

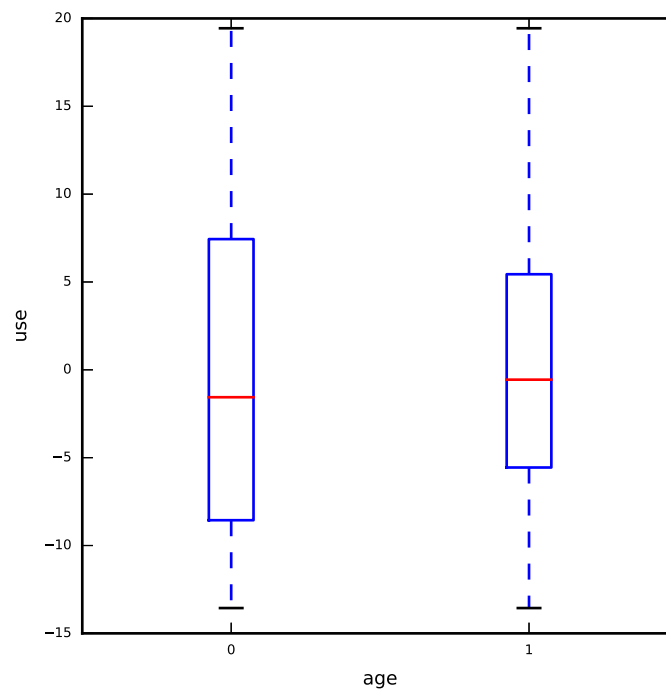
As our test statistic is $1.334 < 1.961$ this means that the p value is > 0.05 and so we cannot reject the null hypothesis.

The p-value is in fact 0.1825. .

The t test assumes that the distribution of the response in each group follows a Normal distribution. We could check this by looking at histograms of the variable in each group. If we were concerned about the normality assumption then we could instead use a Mann Whitney (MW) test.

A Mann Whitney test works simply on the order (or ranks) of the responses across the two groups. So the response variable is firstly sorted and then each value is ranked. The ranks for each group are then summed and the value that is larger is compared with what would be expected if there was no difference between the groups.

In this case the MW U statistic is 413204 which for samples of size 1175 and 759 corresponds to a p value of 0.0127.



Choosing appropriate random classifications

We begin this section by deciding which of the possible random classifications to include in the modelling.

This is done by fitting combinations in turn and picking more complicated models if they make a significant improvement via a Wald test. All models are displayed along with their chi-squared test statistic in the table below:

Higher-level classifications	Significance
cons	nan

The best model based on the Likelihood has levels:

As this is a multilevel modelling SAA we will also want to look at how the response is distributed across the levels of the model.

For this we will use the best model chosen above and look at how the variance is distributed across levels.

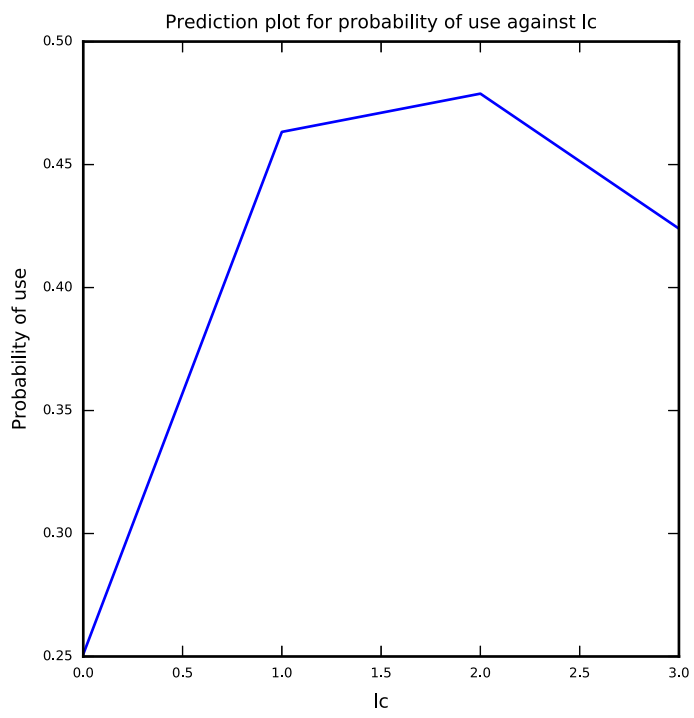
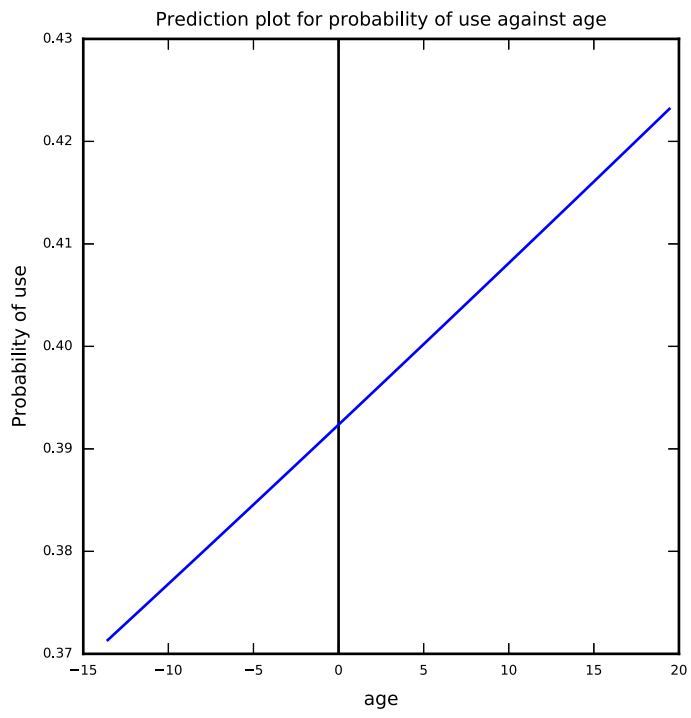
Variable	Coefficient	SE
Intercept	-0.437	0.0466

Performing univariable modelling

Our next step in modelling now that we have a set of potential predictors is to consider models for each predictor in turn along with a random intercept at each chosen classification from the best model in the last section. In the fixed part these models simply contain an intercept and the particular predictor and so for continuous predictors will be multilevel linear regressions and for categorical predictors will be multilevel generalisations of ANOVAs. In the table below we summarise the modelling by showing the coefficients for each predictor along with the p value comparing the model with that predictor with a Null model. This Univariable modelling step will identify a set of candidate predictors to be taken forward into the next stage of modelling.

Variable	Coefficient	SE	p value	Significance
age	0.00657	0.00516	0.203	N/S
lc_1	0.946	0.146	< 0.001	***
lc_2	1.009	0.152		
lc_3	0.787	0.125		

Which predictors we consider for the next stage of analysis will depend on their significance in the above table (but may in practice also depend on the size the effect and substantive interest of the variable though this is hard to automate). We will use a threshold on the p values associated with the predictors to decide whether to include the predictors in the next stage. Here we are currently using a threshold of 0.05. so the predictors to carry forward are: lc.



Looking at correlations between predictors

Our next step is to check that none of the correlations between the predictor variables are too great as this could cause estimation problems when we add the predictors to the model together. To do this we look at all correlations between the predictor variables that have been identified as significant univariably and are thus candidates to be added to the model.

The correlations are as follows:

Variables	Correlation
(lc_1, age)	-0.206
(lc_2, age)	0.013
(lc_2, lc_1)	-0.206
(lc_3, age)	0.632
(lc_3, lc_1)	-0.374
(lc_3, lc_2)	-0.343

Correlations greater than 0.8 (in magnitude) are worth looking at as they may result in model fitting problems when both predictors are included.

Performing multivariable model selection - random intercept models

In this next stage we will look at the best random intercepts model using only main effects for the variables to be considered. You have chosen to perform forward pass which is a quicker method than full forward selection. It may therefore not explore as many possible models. The predictor variables are considered in turn based on their significance in the univariable analysis and each is added to the current model. If the resulting model is a significant improvement then the predictor is kept in the model otherwise it is removed. Attention then moves on to the next predictor until all predictors are considered.

You have chosen to use Wald tests to compare models. These work by looking at estimates and standard error matrices for each predictor to assess significance and run quicker than the alternative methods as they do not need to run submodels.

The most significant predictor in the univariable analysis was lc so our starting point in multivariable modelling is the model:

$$use_i \sim \text{Binomial}(cons_i, p_i), \text{logit}(p_i) = \beta_0 lc_1_i + \beta_1 lc_2_i + \beta_2 lc_3_i + \beta_3 \text{intercept}_i$$

Variable	Coefficient	SE	p value	Significance
lc_1	0.946	0.146	< 0.001	***
lc_2	1.009	0.152		
lc_3	0.787	0.125		
Intercept	-1.094	0.1		

Adding variable lc was a significant improvement and so we retain it in the model.

Our next step is to consider adding variable age to the current model.

$$use_i \sim \text{Binomial}(cons_i, p_i), \text{logit}(p_i) = \beta_0 lc_1_i + \beta_1 lc_2_i + \beta_2 lc_3_i + \beta_3 \text{age}_i + \beta_4 \text{intercept}_i$$

Variable	Coefficient	SE	p value	Significance
lc_1	1.031	0.15	< 0.001	***
lc_2	1.184	0.164		
lc_3	1.112	0.168		
age	-0.0217	0.00741	0.003	**
Intercept	-1.264	0.117		

Adding variable age was a significant improvement and so we retain it in the model.

This is our final model.

Choosing interactions

You have chosen not to investigate interactions and so this page is empty.

Adding random slopes

You have chosen not to look at random slopes and so this page is blank.

Analysing the residuals

Here we look at the residuals from the model and plot them in various ways.

Looking at predictions

Having fitted a model with several predictors we might like to represent this model graphically. This is more difficult than when we have only one predictor and so for now we consider each predictor in turn and set all other predictors to their mean values.

