

# Homework Assignment - 05

*Critical Thinking Group 5*

*Arindam Barman*

*Mohamed Elmoudni*

*Shazia Khan*

*Kishore Prasad*

## Contents

<b>Overview</b>	<b>3</b>
<b>Objective</b>	<b>3</b>
<b>1 Data Exploration Analysis</b>	<b>3</b>
1.1 Variable identification . . . . .	3
1.2 Variable Relationships . . . . .	4
1.3 Data summary analysis . . . . .	5
1.4 Outliers Identification . . . . .	8
<b>2. Data Preparation</b>	<b>8</b>
2.1 Missing Flags . . . . .	8
2.2 Missing values treatment . . . . .	9
2.3 Outliers treatment . . . . .	9
2.4 Dummy Variables . . . . .	10
2.5 Correlation for new variables . . . . .	10
<b>3. Build Models</b>	<b>11</b>
3.1 Poisson models . . . . .	13
3.1.1 Poisson Model 1 . . . . .	13
3.1.1.2 Interpretation Poisson Model 1 . . . . .	13
3.1.1.3 Coefficient Analysis: . . . . .	14
3.1.1.3 Overdispersion Analysis: . . . . .	15
3.1.2 Quasi-Poisson model . . . . .	15

3.1.2.1 Interpretation Quasi-Poisson model	16
3.1.3 zero-inflation model	16
3.1.3.1 Coefficient Analysis:	17
3.1.3.2 Overdispersion Analysis	17
3.2 Poisson Model 2	18
3.2.1 Interpretation Poisson Model 2	19
3.2.1.1 Overdispersion Analysis	19
3.2.2 Quasi-Poisson model 2	20
3.2.2.1 Interpretation Quasi-Poisson model 2	21
3.2.3 zero-inflation model	21
3.2 Negative Binomial models	23
3.2.1 Negative Binomial model 3	23
3.3 Linear Regression models	32
3.3.1 Linear Regression Model 5	32
3.3.1 Linear Regression Model 6	34
<b>4 Model Selection</b>	<b>36</b>
<b>5 Prediction Using Evaluation Data</b>	<b>37</b>
5.1 Transformation of Evaluation Data	38
5.2 Model Output	38
5.3 Conclusion	39
<b>Appendix A: DATA621 Homework 05 R Code</b>	<b>40</b>

# Overview

In this homework assignment, we will explore, analyze and model a data set containing information on approximately 12795 commercially available wines using 16 variables. The variables are mostly related to the chemical properties of the wine being sold. The response variable is the number of sample cases of wine that were purchased by wine distribution companies after sampling a wine. These cases would be used to provide tasting samples to restaurants and wine stores around the United States. The more sample cases purchased, the more likely is a wine to be sold at a high end restaurant. A large wine manufacturer is studying the data in order to predict the number of wine cases ordered based upon the wine characteristics. If the wine manufacturer can predict the number of cases, then that manufacturer will be able to adjust their wine offering to maximize sales.

## Objective

Our objective is to build a count regression model to predict the number of cases of wine that will be sold given certain properties of the wine. Using the training data set, we will build at least two different Poisson regression models, at least two different negative binomial regression models, and at least two multiple linear regression models, using different variables (or the same variables with different transformations).

To attain our objective, we will be following the below best practice steps and guidelines:

- 1 -Data Exploration
- 2 -Data Preparation
- 3 -Build Models
- 4 -Select Models

## 1 Data Exploration Analysis

In section we will explore and gain some insights into the dataset by pursuing the below high level steps and inquiries:

- Variable identification
- Variable Relationships
- Data summary analysis
- Outliers and Missing Values Identification

### 1.1 Variable identification

First we look the variables' datatypes and their roles.

Variable	Datatype	Role
INDEX	int	none
TARGET	int	response
FixedAcidity	num	predictor
VolatileAcidity	num	predictor
CitricAcid	num	predictor
ResidualSugar	num	predictor

Variable	Datatype	Role
Chlorides	num	predictor
FreeSulfurDioxide	num	predictor
TotalSulfurDioxide	num	predictor
Density	num	predictor
pH	num	predictor
Sulphates	num	predictor
Alcohol	num	predictor
LabelAppeal	int	predictor
AcidIndex	int	predictor
STARS	int	predictor

From the Table 1 above, we see that that all variables are quantitative mainly of numeric and integer datatype. Also, we will ignore the INDEX variable as it is just a unique identifier for each row. However, we will use the TARTGET variable as response variable and the remaining variables as predictors.

## 1.2 Variable Relationships

Next let's display and examine the variable relationships as shown in table 2.

Table 2: Variable Description

VARIABLE	DEFINITION	THEORETICAL.EFFECT
INDEX	Identification Variable (do not use)	None
TARGET	Number of Cases Purchased	None
AcidIndex	Proprietary method of testing total acidity of wine by using a weighted average	
Alcohol	Alcohol Content	
Chlorides	Chloride content of wine	
CitricAcid	Citric Acid Content	
Density	Density of Wine	
FixedAcidity	Fixed Acidity of Wine	
FreeSulfurDioxide	Sulfur Dioxide content of wine	
LabelAppeal	Marketing Score indicating the appeal of label design for consumers. High numbers suggest customers like the label design. Negative numbers suggest customers don't like the design.	Many consumers purchase based on the visual appeal of the wine label design. Higher numbers suggest better sales.
ResidualSugar	Residual Sugar of wine	
STARS	Wine rating by a team of experts. 4 Stars = Excellent, 1 Star = Poor	A high number of stars suggests high sales
Sulphates	Sulfate conten of wine	
TotalSulfurDioxide	Total Sulfur Dioxide of Wine	
VolatileAcidity	Volatile Acid content of wine	
pH	pH of wine	

At first glance, we can easily deduce that that the FreeSulfurDioxide (Sulfur Dioxide content of wine) can be

derived from the TotalSulfurDioxide (Total Sulfur Dioxide of Wine). However, looking closer at the role of the sulfur dioxide  $SO_2$ , as it is used as a preservative because of its anti-oxidative and anti-microbial properties in wine and also as a cleaning agent for barrels and winery facilities, we realize that when a winemaker says his/her wine has 100 ppm (part per million) of  $SO_2$ , he/she is most probably referring to the total amount of  $SO_2$  in his wine, and that means:

total  $SO_2$  = free  $SO_2$  + bound  $SO_2$ .

free  $SO_2$ : molecular  $SO_2$  + bisulfites + sulfites

bound  $SO_2$ : sulfites attached to either sugars, acetaldehyde or phenolic compounds

In this case the free  $SO_2$  portion (not associated with wine molecules) is effectively the buffer against microbes and oxidation... Hence without knowing the bound  $SO_2$ , we won't be able to derive FreeSulfurDioxide from TotalSulfurDioxide.

Also, looking briefly at the VolatileAcidity (Volatile Acid content of wine) and FixedAcidity (Fixed Acidity of Wine), we can easily deduce AcidIndex as the Acid index = Total acid (g/L) - pH. where Total acidity = Volatile Acid + Fixed Acidity. However, in our case the index is weighted average and we don't know the weighted average of either Volatile Acid or Fixed Acidity. Hence we will assume these variable do not have strict arithmetic relationships.

### 1.3 Data summary analysis

In this section, we will create summary data to better understand the initial relationship variables have with our dependent variable using correlation, central tendency, and dispersion As shown in table 3.

```
## 'data.frame': 12795 obs. of 15 variables:
## $ TARGET : int 3 3 5 3 4 0 0 4 3 6 ...
## $ FixedAcidity : num 3.2 4.5 7.1 5.7 8 11.3 7.7 6.5 14.8 5.5 ...
## $ VolatileAcidity : num 1.16 0.16 2.64 0.385 0.33 0.32 0.29 -1.22 0.27 -0.22 ...
## $ CitricAcid : num -0.98 -0.81 -0.88 0.04 -1.26 0.59 -0.4 0.34 1.05 0.39 ...
## $ ResidualSugar : num 54.2 26.1 14.8 18.8 9.4 ...
## $ Chlorides : num -0.567 -0.425 0.037 -0.425 NA 0.556 0.06 0.04 -0.007 -0.277 ...
## $ FreeSulfurDioxide : num NA 15 214 22 -167 -37 287 523 -213 62 ...
## $ TotalSulfurDioxide: num 268 -327 142 115 108 15 156 551 NA 180 ...
## $ Density : num 0.993 1.028 0.995 0.996 0.995 ...
## $ pH : num 3.33 3.38 3.12 2.24 3.12 3.2 3.49 3.2 4.93 3.09 ...
## $ Sulphates : num -0.59 0.7 0.48 1.83 1.77 1.29 1.21 NA 0.26 0.75 ...
## $ Alcohol : num 9.9 NA 22 6.2 13.7 15.4 10.3 11.6 15 12.6 ...
## $ LabelAppeal : int 0 -1 -1 -1 0 0 0 1 0 0 ...
## $ AcidIndex : int 8 7 8 6 9 11 8 7 6 8 ...
## $ STARS : int 2 3 3 1 2 NA NA 3 NA 4 ...
```

Table 3: Data Summary

	mean	sd	median	trimmed
TARGET	3.0290739	1.9263682	3.00000	3.0538244
FixedAcidity	7.0757171	6.3176435	6.90000	7.0736739
VolatileAcidity	0.3241039	0.7840142	0.28000	0.3243890
CitricAcid	0.3084127	0.8620798	0.31000	0.3102520
ResidualSugar	5.4187331	33.7493790	3.90000	5.5800410
Chlorides	0.0548225	0.3184673	0.04600	0.0540159
FreeSulfurDioxide	30.8455713	148.7145577	30.00000	30.9334877
TotalSulfurDioxide	120.7142326	231.9132105	123.00000	120.8895367

	mean	sd	median	trimmed
Density	0.9942027	0.0265376	0.99449	0.9942130
pH	3.2076282	0.6796871	3.20000	3.2055706
Sulphates	0.5271118	0.9321293	0.50000	0.5271453
Alcohol	10.4892363	3.7278190	10.40000	10.5018255
LabelAppeal	-0.0090660	0.8910892	0.00000	-0.0099639
AcidIndex	7.7727237	1.3239264	8.00000	7.6431572
STARS	2.0417550	0.9025400	2.00000	1.9711258

Below is the missing values and correlation table of the predictor variables to the response variables.

Table 4: Missing Data and Data Correlation

	Missing	Correlation
TARGET	0	1.0000000
FixedAcidity	0	-0.0490109
VolatileAcidity	0	-0.0887932
CitricAcid	0	0.0086846
ResidualSugar	616	0.0164913
Chlorides	638	-0.0382631
FreeSulfurDioxide	647	0.0438241
TotalSulfurDioxide	682	0.0514784
Density	0	-0.0355175
pH	395	-0.0094448
Sulphates	1210	-0.0388496
Alcohol	653	0.0620616
LabelAppeal	0	0.3565005
AcidIndex	0	-0.2460494
STARS	3359	0.5587938

### *Missing Values and Correlation Interpretation*

From tables 3 and 4 above, we observe the followings:

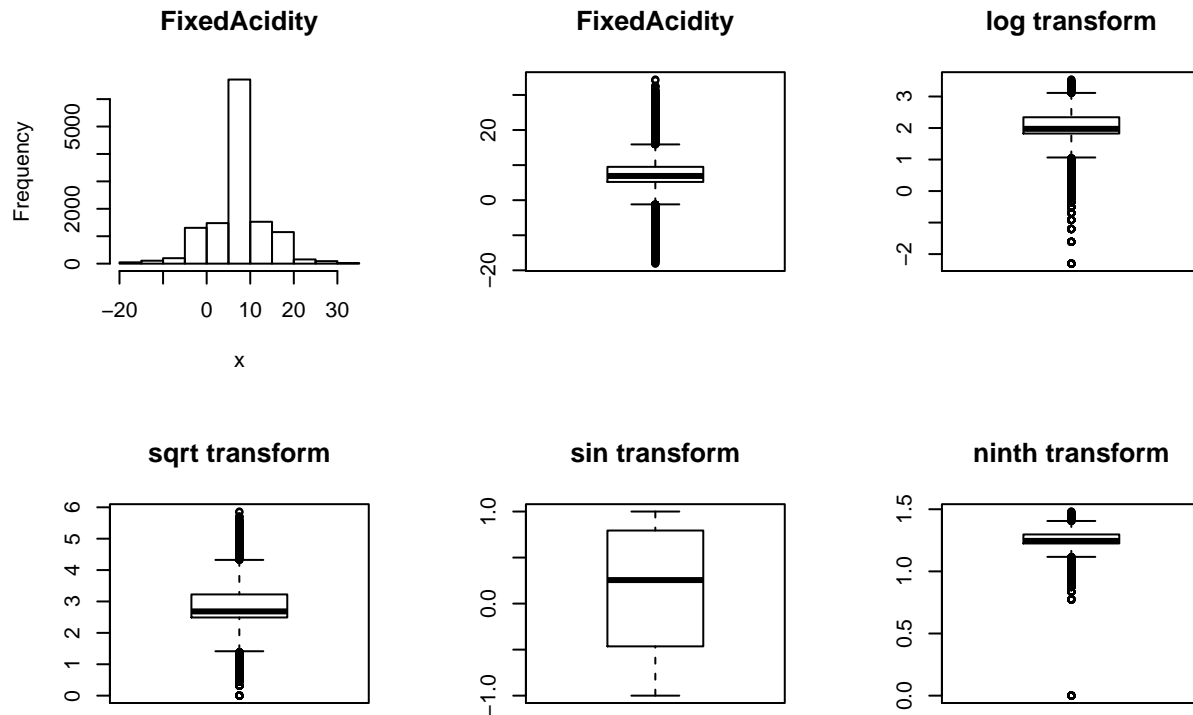
- Variable ResidualSugar has 616 and 0.0164913 correlation. Given the low correlation we will try some imputation techniques to handle the missing the values and replace missing values with their respective value.
- variable Chlorides 638 -0.0382631 correlation. . Given the low negative correlation we will try we would replace missing values with their respective value
- Variable FreeSulfurDioxide 647 0.0438241. Given the low correlation we will impute the missing values with their respective value
- Variable TotalSulfurDioxide has 682 missing values with 0.0514784 correlation. Given the low correlation we will impute the missing values with their respective value.
- Variable Alcohol has 653 missing values with 0.0620616 correlation. Given the low correlation we will impute the missing values with their respective value.

Please note that ResidualSugar, Chlorides, FreeSulfurDioxide, Alcohol, and TotalSulfurDioxide variables have similar number of missing values. They are chemically related. However, we don't think they are arithmetically related.

- In addition, variable pH has 395 missing values with negative correlation of -0.0094448. Again we may just ignore these missing values especially that it has very low negative correlation to the target variable.
- Variable Sulphates has much higher missing values of 1210 with low negative correlation of -0.0388496. We will be imputing this values with their respective value
- Now, variable STARS has the highest missing values of 3359 and highest correlation of 0.5587938. This is very important variable and it drives sales and consequently heavily impacts our response variable. We have to be careful in fixing the missing values as this variable STARS is rating score variable with 1 being the lowest and 4 the highest

## 1.4 Outliers Identification

In this section we look at boxplots to determine the outliers in variables and decide on whether to act on the outliers. Lets do some univariate analysis. We will look at the Histogram and Boxplot for each variable to detect outliers if any and treat it accordingly.



\*\*\*Please note that we generated the above plots for all other variables. However we hid the results for ease of streamlining our report.

## 2. Data Preparation

Now that we have completed the preliminary analysis, we will be cleaning and consolidating data into one dataset for use in analysis and modeling. We will be purging the below steps as guidelines:

- Missing Flags
- Missing values treatment
- Outliers treatment
- Dummy Variables

### 2.1 Missing Flags

We create flag variables to indicate whether some of the fields are missing any values. If the value is missing, we code it with 1 and if the value is present we code it with 0. The following are the variables that are created:



- ResidualSugar\_MISS
- Chlorides\_MISS
- FreeSulfurDioxide\_MISS
- TotalSulfurDioxide\_MISS
- pH\_MISS
- Sulphates\_MISS
- Alcohol\_MISS
- STARS\_MISS

## 2.2 Missing values treatment

Next we impute missing values. We can go ahead and use the mean as impute values. We will replace the missing values in the original variables. However, for STARS, we will code the missing value as a '0' instead of a mean. The following are the variables that are impacted:

- ResidualSugar
- Chlorides
- FreeSulfurDioxide
- TotalSulfurDioxide
- pH
- Sulphates
- Alcohol
- STARS

## 2.3 Outliers treatment

For outliers, we will use the capping method. In this method, we will replace all outliers that lie outside the 1.5 times of IQR limits. We will cap it by replacing those observations less than the lower limit with the value of 5th %ile and those that lie above the upper limit with the value of 95th %ile.

Accordingly we create the following new variables while retaining the original variables.

- FixedAcidity\_CAP
- VolatileAcidity\_CAP
- CitricAcid\_CAP
- ResidualSugar\_CAP
- Chlorides\_CAP
- FreeSulfurDioxide\_CAP
- TotalSulfurDioxide\_CAP
- Density\_CAP
- pH\_CAP
- Sulphates\_CAP
- Alcohol\_CAP
- AcidIndex\_CAP

## 2.4 Dummy Variables

Finally, we will also create dummy variables for the following variables:

- LabelAppeal : For this variable, we create a dummy variable to indicate if the value is Zero / Positive or Negative.
- STARS - We create a Dummy Variable for each of the star ratings - 1,2,3,4. The value is 1 in the respective variable based on the STARS value. A Zero value in all of the STARS dummy vars indicate that the value was missing in the original variable.

## 2.5 Correlation for new variables

Lets see how the new variables stack up against the TARGET.

Table 5: Correlation between TARGET and predictor variables

	Correlation
STARS_3	0.3597277
STARS_4	0.2783731
STARS_2	0.2484240
Alcohol_CAP	0.0634633
TotalSulfurDioxide_CAP	0.0503492
FreeSulfurDioxide_CAP	0.0417585
LabelAppeal_Positive	0.0206261
ResidualSugar_CAP	0.0204409
CitricAcid_CAP	0.0120351
ResidualSugar_MISS	0.0111995
TotalSulfurDioxide_MISS	0.0061720
Chlorides_MISS	0.0026937
Alcohol_MISS	0.0014776
FreeSulfurDioxide_MISS	-0.0001501
pH_MISS	-0.0099654
pH_CAP	-0.0102565
Sulphates_MISS	-0.0125039
Chlorides_CAP	-0.0304686
Density_CAP	-0.0315375
Sulphates_CAP	-0.0359312
FixedAcidity_CAP	-0.0510757
VolatileAcidity_CAP	-0.0891214
STARS_1	-0.1300422
AcidIndex_CAP	-0.2353997
STARS_MISS	-0.5715792

From the above Correlations, we can make the following observations:

- The following variables have a positive correlation with TARGET: STARS\_3, STARS\_4, STARS\_2, Alcohol\_CAP, TotalSulfurDioxide\_CAP, FreeSulfurDioxide\_CAP, LabelAppeal\_Positive, ResidualSugar\_CAP, CitricAcid\_CAP, ResidualSugar\_MISS, TotalSulfurDioxide\_MISS, Chlorides\_MISS, Alcohol\_MISS.
- The following variables have a negative correlation with TARGET: FreeSulfurDioxide\_MISS, pH\_MISS,

pH\_CAP, Sulphates\_MISS, Chlorides\_CAP, Density\_CAP, Sulphates\_CAP, FixedAcidity\_CAP, VolatileAcidity\_CAP, STARS\_1, AcidIndex\_CAP, STARS\_MISS.

- Not all variable have a strong correlation in either direction. However, the following stand out for having a stronger correlation: STARS\_MISS, STARS\_3, STARS\_4, STARS\_2, AcidIndex\_CAP, STARS\_1, VolatileAcidity\_CAP, Alcohol\_CAP, FixedAcidity\_CAP, TotalSulfurDioxide\_CAP.

### 3. Build Models

Since we are dealing with count variables, our modeling technique will mainly focus on using variation of the Generalized Linear Model (GLM) family functions. We will start with the classical Poisson regression; then we will enhance it using model Negative binominal model.

In addition, we will also create models using linear regression.

Using original and transformed datasets, we will build at least ten models as follow:

- Two Poisson models
- Two Quasi-Poisson models
- Two Negative binomial models
- Two Zero-inflated models
- Two Linear regression models

Below is a summary table showing models and their respective variables.

Table 6: Models and their Respective Variables

Variable	Model.1	Model.2	Comments
TARGET	Y	Y	The TARGET variable
FixedAcidity	Y		Imputed with Mean
VolatileAcidity	Y		Imputed with Mean
CitricAcid	Y		Imputed with Mean
ResidualSugar	Y		Imputed with Mean
Chlorides	Y		Imputed with Mean
FreeSulfurDioxide	Y		Imputed with Mean
TotalSulfurDioxide	Y		Imputed with Mean
Density	Y		Imputed with Mean
pH	Y		Imputed with Mean
Sulphates	Y		Imputed with Mean
Alcohol	Y		Imputed with Mean
LabelAppeal	Y		Original Variable
AcidIndex	Y		Imputed with Mean
STARS	Y		Original Variable
ResidualSugar_MISS		Y	Missing Flag
Chlorides_MISS		Y	Missing Flag
FreeSulfurDioxide_MISS		Y	Missing Flag
TotalSulfurDioxide_MISS		Y	Missing Flag
pH_MISS		Y	Missing Flag
Sulphates_MISS		Y	Missing Flag
Alcohol_MISS		Y	Missing Flag
STARS_MISS		Y	Missing Flag
FixedAcidity_CAP		Y	Imputed with Mean and Outliers capped
VolatileAcidity_CAP		Y	Imputed with Mean and Outliers capped
CitricAcid_CAP		Y	Imputed with Mean and Outliers capped
ResidualSugar_CAP		Y	Imputed with Mean and Outliers capped

Variable	Model.1	Model.2	Comments
Chlorides_CAP		Y	Imputed with Mean and Outliers capped
FreeSulfurDioxide_CAP		Y	Imputed with Mean and Outliers capped
TotalSulfurDioxide_CAP		Y	Imputed with Mean and Outliers capped
Density_CAP		Y	Imputed with Mean and Outliers capped
pH_CAP		Y	Imputed with Mean and Outliers capped
Sulphates_CAP		Y	Imputed with Mean and Outliers capped
Alcohol_CAP		Y	Imputed with Mean and Outliers capped
AcidIndex_CAP		Y	Imputed with Mean and Outliers capped
LabelAppeal_Positive		Y	Positive or Negative Dummy Variable
STARS_1		Y	Dummy Variable
STARS_2		Y	Dummy Variable
STARS_3		Y	Dummy Variable
STARS_4		Y	Dummy Variable

## 3.1 Poisson models

Our first attempt to capture the relationship between the wine chemical properties and number of cases of the wine being sold in a parametric regression model, we fit the basic Poisson regression model

### 3.1.1 Poisson Model 1

We will explore the Poisson regression model Using original data with replacing all missing data with the means.

```
##
## Call:
## glm(formula = TARGET ~ ., family = poisson, data = winedata_orig)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.9733  -0.7200   0.0694   0.5785   3.2315
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.526e+00  1.955e-01   7.807 5.87e-15 ***
## FixedAcidity   -3.045e-04  8.205e-04  -0.371 0.710502
## VolatileAcidity -3.343e-02  6.516e-03  -5.131 2.88e-07 ***
## CitricAcid      7.773e-03  5.892e-03   1.319 0.187124
## ResidualSugar    5.676e-05  1.546e-04   0.367 0.713588
## Chlorides       -4.141e-02  1.645e-02  -2.518 0.011816 *
## FreeSulfurDioxide 1.254e-04  3.512e-05   3.571 0.000356 ***
## TotalSulfurDioxide 8.296e-05  2.275e-05   3.647 0.000266 ***
## Density         -2.823e-01  1.920e-01  -1.471 0.141348
## pH              -1.572e-02  7.638e-03  -2.058 0.039554 *
## Sulphates       -1.267e-02  5.749e-03  -2.205 0.027480 *
## Alcohol         2.201e-03  1.410e-03   1.561 0.118446
## LabelAppeal     1.332e-01  6.063e-03  21.968 < 2e-16 ***
## AcidIndex       -8.705e-02  4.548e-03 -19.139 < 2e-16 ***
## STARS           3.113e-01  4.531e-03  68.700 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 22861  on 12794  degrees of freedom
## Residual deviance: 14728  on 12780  degrees of freedom
## AIC: 46700
##
## Number of Fisher Scoring iterations: 5
```

#### 3.1.1.2 Interpretation Poisson Model 1

From this output, we have the following estimated model:

$$\hat{y} = e^{B_0x_0+B_1x_1+B_2x_2+B_3x_3+B_4x_4+B_5x_5+B_6x_6+B_7x_7+B_8x_8+B_9x_9+B_{10}x_{10}+B_{11}x_{11}+B_{12}x_{12}+B_{13}x_{13}+B_{14}x_{14}}$$

where

$$\begin{aligned}B_0 &= 1.526 \\B_1 &= -3.045e - 04 \\B_2 &= -3.343e - 02 \\B_3 &= 7.773e - 03 \\B_4 &= 5.676e - 05 \\B_5 &= -4.141e - 02 \\B_6 &= 1.254e - 04 \\B_7 &= 8.296e - 05 \\B_8 &= -2.823e - 01 \\B_9 &= -1.572e - 02 \\B_{10} &= -1.267e - 02 \\B_{11} &= 2.201e - 03 \\B_{12} &= 1.332e - 01 \\B_{13} &= -8.705e - 02 \\B_{14} &= 3.113e - 0\end{aligned}$$

and

$$\begin{aligned}x_0 &= 1 \\x_1 &= \textit{FixedAcidity} \\x_2 &= \textit{VolatileAcidity} \\x_3 &= \textit{CitricAcid} \\x_4 &= \textit{ResidualSugar} \\x_5 &= \textit{Chlorides} \\x_6 &= \textit{FreeSulfurDioxide} \\x_7 &= \textit{TotalSulfurDioxide} \\x_8 &= \textit{Density} \\x_9 &= \textit{pH} \\x_{10} &= \textit{Sulphates} \\x_{11} &= \textit{Alcohol} \\x_{12} &= \textit{LabelAppeal} \\x_{13} &= \textit{AcidIndex} \\x_{14} &= \textit{STARS}\end{aligned}$$

### 3.1.1.3 Coefficient Analysis:

In addition, the coefficient for VolatileAcidity, FreeSulfurDioxide, TotalSulfurDioxide, LabelAppeal, AcidIndex, and STARS are highly significant.

Unlike the linear model, in order to interpret the slope coefficient in a Poisson regression, it makes better sense to look at the ratio of predicted responses (instead of the difference) for a unit increase in x. for instance:

$$\frac{e^{b_0+B_1(x+1)}}{e^{b_0+B_1x}} = e^{B_1}$$

For instance, for with  $B_1 = -(0.0003045)$ , we have  $e^{B_1} = e^{-(0.0003045)} = 0.999695$

Thus, for a unit increase in the FixedAcidity, we would expect to see the number of cases of wine that will be sold given certain properties of the wine to decrease by a factor of  $= 0.999695$ .

Hence, for a unit increase in our highly significant variables:

- VolatileAcidity, we expect a decrease of  $e^{-(0.0343)} = 0.9662816$  the number of cases of wine that will be sold

- FreeSulfurDioxide, we expect an increase of  $e^{0.0000829} = 1.000083$  the number of cases of wine that will be sold
- TotalSulfurDioxide, we expect a decrease of  $e^{-(0.2823)} = 0.7540474$  the number of cases of wine that will be sold
- LabelAppeal, we expect a increase of  $e^{(.1332)} = 1.142478$  the number of cases of wine that will be sold
- AcidIndex, we expect a decrease of  $e^{-(.08705)} = 0.9166313$  the number of cases of wine that will be sold
- STARS, we expect a increase of  $e^{(3.113)} = 22.48841$  the number of cases of wine that will be sold

### 3.1.1.3 Overdispersion Analysis:

Another common problem with Poisson regression is that the response is more variable than what is expected by the model; this is called overdispersion. Thus checking for overdispersion, we will examine if the residual deviance greatly exceeds the residual degrees of freedom, then that is an indication of an overdispersion problem.

For our model(1), we see that our Residual deviance is 14728 and degrees of freedom is 12780; our Residual deviance 1.15 greater than our Residual degrees of freedom. Hence, the response is little more variable than what is expected by model (1). However, we won't address this issue as the Residual deviance does not greatly exceed residual degrees of freedom.

Sine we see that we have over dispersion, let's find out the dispersion parameter  $\phi$ . Since the variance in the Poisson model is identical to the mean, the expectations are to have  $\phi = 1$ .

```
## [1] 0.851513
```

Our dispersion parameter is 0.851513; obviously it is not 1.

### 3.1.2 Quasi-Poisson model

Another way of dealing with over-dispersion is to use Quasi-Poisson model which uses the mean regression function and the variance function from the Poisson GLM but to leave the dispersion parameter  $\phi$  unrestricted. Thus,  $\phi$  is not assumed to be fixed at 1 but is estimated from the data. This strategy leads to the same coefficient estimates as the standard Poisson model but inference is adjusted for over-dispersion.

```
##
## Call:
## glm(formula = TARGET ~ ., family = quasipoisson, data = winedata_orig)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.9733  -0.7200   0.0694   0.5785   3.2315
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    1.526e+00  1.804e-01   8.460 < 2e-16 ***
## FixedAcidity   -3.045e-04  7.571e-04  -0.402  0.68751
## VolatileAcidity -3.343e-02  6.013e-03  -5.560 2.75e-08 ***
## CitricAcid      7.773e-03  5.437e-03   1.430  0.15288
## ResidualSugar   5.676e-05  1.427e-04   0.398  0.69082
## Chlorides      -4.141e-02  1.518e-02  -2.728  0.00638 **
## FreeSulfurDioxide 1.254e-04  3.241e-05   3.869  0.00011 ***
## TotalSulfurDioxide 8.296e-05  2.099e-05   3.952  7.80e-05 ***
```

```
## Density          -2.823e-01  1.771e-01  -1.594  0.11099
## pH               -1.572e-02  7.048e-03  -2.231  0.02572 *
## Sulphates        -1.267e-02  5.305e-03  -2.389  0.01690 *
## Alcohol           2.201e-03  1.301e-03   1.692  0.09067 .
## LabelAppeal       1.332e-01  5.595e-03  23.806 < 2e-16 ***
## AcidIndex         -8.705e-02  4.197e-03 -20.741 < 2e-16 ***
## STARS             3.113e-01  4.181e-03  74.449 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for quasipoisson family taken to be 0.85152)
##
## Null deviance: 22861 on 12794 degrees of freedom
## Residual deviance: 14728 on 12780 degrees of freedom
## AIC: NA
##
## Number of Fisher Scoring iterations: 5
```

### 3.1.2.1 Interpretation Quasi-Poisson model

```
qpr <- residuals(fm_qpois,"pearson")
qphi <- sum(qpr^2)/df.residual(fm_qpois)
qphi
```

```
## [1] 0.851513
```

Please note that the Quasi-Poisson model leads to the same coefficient estimates as the standard Poisson model but inference is adjusted for over-dispersion. Hence please refer to Poisson model Coefficient Analysis for details.

Please note that dispersion parameter in the Quasi-Poisson model is 0.851513; which is similar to that of the classical Poisson Model (1)

### 3.1.3 zero-inflation model

Next we will proceed with zero-inflation model as another very common occurrence when working with count data is that there will be an overabundance of zero counts which is not consistent with the Poisson model.

```
##
## Call:
## zeroinfl(formula = TARGET ~ ., data = winedata_orig, dist = "poisson")
##
## Pearson residuals:
##      Min      1Q   Median      3Q      Max
## -2.122598 -0.404868 -0.007538  0.371282  5.769256
##
## Count model coefficients (poisson with log link):
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.443e+00  2.020e-01   7.146 8.91e-13 ***
## FixedAcidity    3.383e-04  8.420e-04   0.402 0.687855
```



```

## VolatileAcidity    -1.211e-02  6.721e-03  -1.801  0.071625 .
## CitricAcid         4.926e-04  6.024e-03   0.082  0.934822
## ResidualSugar      -7.702e-05  1.586e-04  -0.485  0.627336
## Chlorides          -2.241e-02  1.691e-02  -1.325  0.185076
## FreeSulfurDioxide  2.546e-05  3.547e-05   0.718  0.472877
## TotalSulfurDioxide -1.783e-05  2.265e-05  -0.787  0.431015
## Density            -2.845e-01  1.983e-01  -1.435  0.151310
## pH                 5.931e-03  7.859e-03   0.755  0.450387
## Sulphates          1.726e-04  5.919e-03   0.029  0.976735
## Alcohol             6.886e-03  1.440e-03   4.783  1.72e-06 ***
## LabelAppeal        2.330e-01  6.303e-03  36.962  < 2e-16 ***
## AcidIndex          -1.858e-02  4.898e-03  -3.794  0.000148 ***
## STARS              1.009e-01  5.201e-03  19.403  < 2e-16 ***
##
## Zero-inflation model coefficients (binomial with logit link):
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -4.4483881  1.3374162  -3.326  0.000881 ***
## FixedAcidity     0.0007591  0.0055469   0.137  0.891146
## VolatileAcidity  0.1937198  0.0438512   4.418  9.98e-06 ***
## CitricAcid      -0.0296037  0.0399713  -0.741  0.458922
## ResidualSugar   -0.0011765  0.0010429  -1.128  0.259307
## Chlorides        0.0921158  0.1093491   0.842  0.399564
## FreeSulfurDioxide -0.0007419  0.0002422  -3.063  0.002190 **
## TotalSulfurDioxide -0.0009866  0.0001523  -6.476  9.41e-11 ***
## Density          0.4900517  1.3159510   0.372  0.709600
## pH               0.2160935  0.0512207   4.219  2.46e-05 ***
## Sulphates        0.1323441  0.0387670   3.414  0.000641 ***
## Alcohol          0.0279120  0.0095782   2.914  0.003567 **
## LabelAppeal      0.7229711  0.0429468  16.834  < 2e-16 ***
## AcidIndex        0.4347418  0.0258387  16.825  < 2e-16 ***
## STARS            -2.3768721  0.0603161 -39.407  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Number of iterations in BFGS optimization: 39
## Log-likelihood: -2.041e+04 on 30 Df

```

### 3.1.3.1 Coefficient Analysis:

We noticed that some variables have their coefficient sign changed from negative to positive and vice versa. For instance;

FixedAcidity changed from -3.045e-04 in model 1 to 3.383e-04 in the zip model ResidualSugar changed from 5.676e-05 in model 1 to -7.702e-05 in the zip model TotalSulfurDioxide changed from 8.296e-05 in model 1 to -1.783e-05 in the zip model. pH changed from -1.572e-02 in model 1 to pH 5.931e-03 in the zip model. Sulphates changed from -1.267e-02 in model 1 to 1.726e-04 in the zip model.

### 3.1.3.2 Overdispersion Analysis

Please note that dispersion parameter in the zero-inflation model is 0.4636815; which is lower than of the classical Poisson Model (1)

```
zippr <- residuals(mod1zip,"pearson")
zipphi <- sum(zippr^2)/df.residual(mod1zip)
zipphi
```

```
## [1] 0.4636815
```

Note that the zip model output above does not indicate in any way if our zero-inflated model is an improvement over a standard Poisson regression. We can determine this by running the corresponding standard negative Poisson model and then performing a Vuong test of the two models.

```
## Vuong Non-Nested Hypothesis Test-Statistic:
## (test-statistic is asymptotically distributed N(0,1) under the
## null that the models are indistinguishable)
## -----
##              Vuong z-statistic              H_A      p-value
## Raw              47.98330 model1 > model2 < 2.22e-16
## AIC-corrected    47.73759 model1 > model2 < 2.22e-16
## BIC-corrected    46.82150 model1 > model2 < 2.22e-16
```

The Vuong test suggests that the zero-inflated Poisson model is slight improvement over a standard Poisson model.

```
#vuong(fm_zinb0,poismod1)
```

## 3.2 Poisson Model 2

In this model we will be using the basic Poisson regression model; however using transformed data.

```
##
## Call:
## glm(formula = TARGET ~ ., family = poisson, data = winedata_trans)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.9184  -0.8511  -0.0111   0.5226   4.0826
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    2.570e+00  2.001e-01  12.841 < 2e-16 ***
## ResidualSugar_MISS  2.283e-02  2.340e-02   0.976  0.32923
## Chlorides_MISS    3.017e-03  2.330e-02   0.130  0.89694
## FreeSulfurDioxide_MISS  2.300e-02  2.366e-02   0.972  0.33101
## TotalSulfurDioxide_MISS  1.883e-02  2.246e-02   0.838  0.40176
## pH_MISS        -3.495e-02  2.991e-02  -1.169  0.24258
## Sulphates_MISS  -6.758e-03  1.757e-02  -0.385  0.70054
## Alcohol_MISS    2.136e-02  2.306e-02   0.926  0.35434
## STARS_MISS      -1.471e+00  2.371e-02 -62.039 < 2e-16 ***
## FixedAcidity_CAP -5.712e-04  9.179e-04  -0.622  0.53372
## VolatileAcidity_CAP -3.550e-02  7.248e-03  -4.898  9.66e-07 ***
```

```
## CitricAcid_CAP      7.430e-03  6.527e-03   1.138  0.25492
## ResidualSugar_CAP   1.348e-04  1.538e-04   0.876  0.38090
## Chlorides_CAP       -2.664e-02  1.618e-02  -1.646  0.09977 .
## FreeSulfurDioxide_CAP 1.600e-04  5.265e-05   3.039  0.00237 **
## TotalSulfurDioxide_CAP 8.381e-05  2.599e-05   3.224  0.00126 **
## Density_CAP         -2.848e-01  1.946e-01  -1.464  0.14332
## pH_CAP              -1.361e-02  8.672e-03  -1.569  0.11667
## Sulphates_CAP       -1.194e-02  5.908e-03  -2.020  0.04334 *
## Alcohol_CAP         3.956e-03  1.646e-03   2.404  0.01622 *
## AcidIndex_CAP       -7.801e-02  5.258e-03 -14.835 < 2e-16 ***
## LabelAppeal_Positive -2.560e-02  1.854e-02  -1.380  0.16746
## STARS_1             -7.179e-01  2.081e-02 -34.504 < 2e-16 ***
## STARS_2             -3.427e-01  1.944e-02 -17.628 < 2e-16 ***
## STARS_3             -1.734e-01  2.006e-02  -8.646 < 2e-16 ***
## STARS_4              NA          NA      NA      NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 22861  on 12794  degrees of freedom
## Residual deviance: 14376  on 12770  degrees of freedom
## AIC: 46368
##
## Number of Fisher Scoring iterations: 6
```

### 3.2.1 Interpretation Poisson Model 2

Most of the coefficients stayed still significant in the model. However, some variables experienced a decrease in p values especially the ones that have capped; which was expected as in the original they had untreated outliers. For instance FixedAcidity p-value went from 0.710502 to 0.53372. The same for ResidualSugar variable went from 0.713588 to 0.38090. Again this is due to outliers' treatment.

In addition, the Poisson model with transformed data has a slight improved as its AIC, 46368, is slightly lower than the model 1 AIC (46700.); which was run against the original data.

#### 3.2.1.1 Overdispersion Analysis

For our model(2), we see that our Residual deviance is 14376 and degrees of freedom is 12770; our Residual deviance 1.12 greater than our Residual degrees of freedom. Hence, the response is little more variable than what is expected by model (2). Please note that this is a slight improvement from model 1 with original data which was 1.15.

Sine we see that we have over dispersion, let's find out the dispersion parameter  $\phi$ . Since the variance in the Poisson model is identical to the mean, the expectations are to have  $\phi = 1$ .

```
pr2 <- residuals(poismod2,"pearson")
phi2 <- sum(pr2^2)/df.residual(poismod2)
phi2
```

```
## [1] 0.9667917
```

Our dispersion parameter for Modle (2) is 0.9667917 which is much closer to 1 than the dispersion parameter of our Modle (1).

### 3.2.2 Quasi-Poisson model 2

```
mod2qpois <- glm(TARGET ~ ., data = winedata_trans, family = quasipoisson)
summary(mod2qpois)
```

```
##
## Call:
## glm(formula = TARGET ~ ., family = quasipoisson, data = winedata_trans)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.9184  -0.8511  -0.0111   0.5226   4.0826
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    2.570e+00  1.968e-01  13.060 < 2e-16 ***
## ResidualSugar_MISS  2.283e-02  2.301e-02   0.992  0.32108
## Chlorides_MISS     3.017e-03  2.291e-02   0.132  0.89520
## FreeSulfurDioxide_MISS  2.300e-02  2.326e-02   0.989  0.32286
## TotalSulfurDioxide_MISS  1.883e-02  2.208e-02   0.853  0.39380
## pH_MISS          -3.495e-02  2.941e-02  -1.188  0.23468
## Sulphates_MISS    -6.758e-03  1.728e-02  -0.391  0.69570
## Alcohol_MISS      2.136e-02  2.267e-02   0.942  0.34622
## STARS_MISS        -1.471e+00  2.332e-02 -63.095 < 2e-16 ***
## FixedAcidity_CAP   -5.712e-04  9.025e-04  -0.633  0.52679
## VolatileAcidity_CAP -3.550e-02  7.126e-03  -4.982  6.38e-07 ***
## CitricAcid_CAP     7.430e-03  6.417e-03   1.158  0.24694
## ResidualSugar_CAP   1.348e-04  1.512e-04   0.891  0.37286
## Chlorides_CAP      -2.664e-02  1.591e-02  -1.674  0.09415 .
## FreeSulfurDioxide_CAP  1.600e-04  5.177e-05   3.091  0.00200 **
## TotalSulfurDioxide_CAP  8.381e-05  2.556e-05   3.279  0.00104 **
## Density_CAP        -2.848e-01  1.913e-01  -1.488  0.13665
## pH_CAP            -1.361e-02  8.527e-03  -1.596  0.11059
## Sulphates_CAP      -1.194e-02  5.809e-03  -2.055  0.03991 *
## Alcohol_CAP        3.956e-03  1.618e-03   2.445  0.01451 *
## AcidIndex_CAP      -7.801e-02  5.170e-03 -15.087 < 2e-16 ***
## LabelAppeal_Positive -2.560e-02  1.823e-02  -1.404  0.16036
## STARS_1            -7.179e-01  2.046e-02 -35.091 < 2e-16 ***
## STARS_2            -3.427e-01  1.911e-02 -17.928 < 2e-16 ***
## STARS_3            -1.734e-01  1.972e-02  -8.793 < 2e-16 ***
## STARS_4              NA         NA      NA      NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for quasipoisson family taken to be 0.9667917)
##
```

```
## Null deviance: 22861 on 12794 degrees of freedom
## Residual deviance: 14376 on 12770 degrees of freedom
## AIC: NA
##
## Number of Fisher Scoring iterations: 6
```

### 3.2.2.1 Interpretation Quasi-Poisson model 2

Please note that the Quasi-Poisson model leads to the same coefficient estimates as the standard Poisson model but inference is adjusted for over-dispersion. Hence please refer to Poisson model Coefficient Analysis for details.

Also, please note that dispersion parameter in the Quasi-Poisson model is 0.9667917; which is similar to that of the classical Poisson Model (2)

### 3.2.3 zero-inflation model

Next we will proceed with zero-inflation model as another very common occurrence when working with count data is that there will be an overabundance of zero counts which is not consistent with the Poisson model.

```
library(sandwich)
library(msm)
library(pscl)
#mod2zip <- zeroinfl(TARGET~ ., data = winedata_trans, dist = "poisson")
#summary(mod2zip)

#####

quine3 <- as.data.frame(model.matrix(poismod2)) ## all regressors
quine3 <- quine3[, !is.na(coef(poismod2))]      ## only identified
quine3 <- quine3[, -1]                        ## omit intercept
quine3$TARGET <- winedata_trans$TARGET        ## add response

## re-fit glm.nb()
fm1a <- glm(TARGET ~ ., data = quine3, family="poisson")
## equivalent to previous fit
logLik(fm1a) - logLik(poismod2)

## 'log Lik.' 0 (df=25)

coef(fm1a) - na.omit(coef(poismod2))
```

```
##          (Intercept)      ResidualSugar_MISS      Chlorides_MISS
##              0              0              0
## FreeSulfurDioxide_MISS TotalSulfurDioxide_MISS      pH_MISS
##              0              0              0
##      Sulphates_MISS      Alcohol_MISS      STARS_MISS
##              0              0              0
##      FixedAcidity_CAP      VolatileAcidity_CAP      CitricAcid_CAP
##              0              0              0
```

```
##      ResidualSugar_CAP      Chlorides_CAP      FreeSulfurDioxide_CAP
##              0              0              0
## TotalSulfurDioxide_CAP      Density_CAP      pH_CAP
##              0              0              0
##      Sulphates_CAP      Alcohol_CAP      AcidIndex_CAP
##              0              0              0
## LabelAppeal_Positive      STARS_1      STARS_2
##              0              0              0
##      STARS_3
##              0
## attr("na.action")
## STARS_4
##      26
## attr("class")
## [1] "omit"
```

```
## fit zeroinfl(), now works
mod2zip<- zeroinfl(TARGET ~ . | 1, data = quine3, dist = "poisson")

summary(mod2zip)
```

```
##
## Call:
## zeroinfl(formula = TARGET ~ . | 1, data = quine3, dist = "poisson")
##
## Pearson residuals:
##      Min      1Q   Median      3Q      Max
## -1.80485 -0.61159  0.06038  0.53998  5.74825
##
## Count model coefficients (poisson with log link):
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      2.474e+00  2.059e-01  12.014 < 2e-16 ***
## ResidualSugar_MISS 2.186e-02  2.400e-02   0.911  0.36234
## Chlorides_MISS     7.389e-03  2.394e-02   0.309  0.75758
## FreeSulfurDioxide_MISS 2.014e-02  2.421e-02   0.832  0.40553
## TotalSulfurDioxide_MISS 2.353e-02  2.306e-02   1.020  0.30756
## pH_MISS           -2.770e-02  3.076e-02  -0.901  0.36781
## Sulphates_MISS    -6.184e-03  1.805e-02  -0.343  0.73193
## Alcohol_MISS       1.696e-02  2.363e-02   0.718  0.47287
## STARS_MISS         -1.360e+00  2.627e-02 -51.786 < 2e-16 ***
## FixedAcidity_CAP   -4.496e-04  9.424e-04  -0.477  0.63328
## VolatileAcidity_CAP -3.032e-02  7.460e-03  -4.064  4.82e-05 ***
## CitricAcid_CAP      5.588e-03  6.699e-03   0.834  0.40420
## ResidualSugar_CAP   7.921e-05  1.577e-04   0.502  0.61543
## Chlorides_CAP      -2.118e-02  1.659e-02  -1.277  0.20171
## FreeSulfurDioxide_CAP 1.529e-04  5.382e-05   2.841  0.00449 **
## TotalSulfurDioxide_CAP 5.926e-05  2.641e-05   2.244  0.02482 *
## Density_CAP        -2.951e-01  2.000e-01  -1.475  0.14018
## pH_CAP             -8.055e-03  8.914e-03  -0.904  0.36617
## Sulphates_CAP      -9.397e-03  6.070e-03  -1.548  0.12164
## Alcohol_CAP         4.775e-03  1.687e-03   2.831  0.00464 **
## AcidIndex_CAP      -6.702e-02  5.560e-03 -12.055 < 2e-16 ***
## LabelAppeal_Positive -2.722e-02  1.903e-02  -1.430  0.15271
## STARS_1            -6.212e-01  2.191e-02 -28.348 < 2e-16 ***
```

```
## STARS_2          -3.267e-01  1.948e-02 -16.773 < 2e-16 ***
## STARS_3          -1.730e-01  2.006e-02  -8.626 < 2e-16 ***
##
## Zero-inflation model coefficients (binomial with logit link):
##           Estimate Std. Error z value Pr(>|z|)
## (Intercept) -2.80683    0.08466  -33.16  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Number of iterations in BFGS optimization: 34
## Log-likelihood: -2.306e+04 on 26 Df
```

```
zippr2 <- residuals(mod2zip,"pearson")
zipphi2 <- sum(zippr2^2)/df.residual(mod2zip)
zipphi2
```

```
## [1] 0.8386535
```

```
vuong(mod2zip,poismod2)
```

```
## Vuong Non-Nested Hypothesis Test-Statistic:
## (test-statistic is asymptotically distributed N(0,1) under the
## null that the models are indistinguishable)
## -----
##           Vuong z-statistic           H_A    p-value
## Raw                6.151478 model1 > model2 3.8382e-10
## AIC-corrected       6.151478 model1 > model2 3.8382e-10
## BIC-corrected       6.151478 model1 > model2 3.8382e-10
```

The Vuong test suggests that the zero-inflated Poisson model is a slight improvement over a standard Poisson model using transformed data.

## 3.2 Negative Binomial models

A more formal way to accommodate over-dispersion in a count data regression model is to use a negative binomial model. Hence we will explore the negative binomial model both in original data as well as transformed data.

### 3.2.1 Negative Binomial model 3

We will explore the Negative Binomial model Using original data with replacing all missing data with the means.

```
##
## Call:
## glm.nb(formula = TARGET ~ ., data = winedata_orig, init.theta = 48974.65509,
##       link = log)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -2.9732 -0.7200 0.0694 0.5785 3.2314
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.526e+00  1.955e-01  7.806 5.88e-15 ***
## FixedAcidity   -3.046e-04  8.205e-04 -0.371 0.710504
## VolatileAcidity -3.343e-02  6.516e-03 -5.131 2.89e-07 ***
## CitricAcid      7.773e-03  5.892e-03  1.319 0.187136
## ResidualSugar   5.676e-05  1.546e-04  0.367 0.713573
## Chlorides       -4.142e-02  1.645e-02 -2.518 0.011817 *
## FreeSulfurDioxide 1.254e-04  3.512e-05  3.571 0.000356 ***
## TotalSulfurDioxide 8.296e-05  2.275e-05  3.647 0.000266 ***
## Density        -2.824e-01  1.920e-01 -1.471 0.141356
## pH             -1.572e-02  7.638e-03 -2.058 0.039552 *
## Sulphates       -1.267e-02  5.749e-03 -2.205 0.027480 *
## Alcohol         2.201e-03  1.410e-03  1.561 0.118467
## LabelAppeal     1.332e-01  6.064e-03 21.967 < 2e-16 ***
## AcidIndex       -8.705e-02  4.548e-03 -19.139 < 2e-16 ***
## STARS           3.113e-01  4.531e-03 68.698 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(48974.66) family taken to be 1)
##
##      Null deviance: 22860  on 12794  degrees of freedom
## Residual deviance: 14728  on 12780  degrees of freedom
## AIC: 46703
##
## Number of Fisher Scoring iterations: 1
##
##              Theta: 48975
##              Std. Err.: 50715
## Warning while fitting theta: iteration limit reached
##
## 2 x log-likelihood: -46670.5
```

\*\*\* Interpretation Negative Binomial Model 3\*\*\*

As per the below table, it is worth noting that the classical Poisson Coefficients are similar to that of the Negative Binomial's.

One possible explanation is that if all we care about is fitting separate means to disjoint subsets of our sample, then GLMs will always yield  $\hat{\mu}_j = \hat{y}_j$  for each subset  $j$ , so the actual error structure and parametrization of the density both become irrelevant to the estimation. In other words, Fitting orthogonal categorical factors by maximum likelihood is equivalent to fitting separate means to disjoint subsets of our sample, so this explains why Poisson and negative binomial GLMs yield the same parameter estimates

In addition, Negative Binomial Model with original data has an AIC value, 46703, is slightly higher than of model 1 AIC (46700.); which was run against the original data.

```
kable(rbind(data.frame("Poisson Coeff"= poismod1$coefficients,"Negative Binom Coeffi" = nbmod3$coefficients
```

	Poisson.Coeff	Negative.Binom.Coeffi
(Intercept)	1.5259824	1.5259982



	Poisson.Coeff	Negative.Binom.Coeffi
FixedAcidity	-0.0003045	-0.0003045
VolatileAcidity	-0.0334329	-0.0334338
CitricAcid	0.0077726	0.0077727
ResidualSugar	0.0000568	0.0000568
Chlorides	-0.0414139	-0.0414151
FreeSulfurDioxide	0.0001254	0.0001254
TotalSulfurDioxide	0.0000830	0.0000830
Density	-0.2823481	-0.2823537
pH	-0.0157219	-0.0157226
Sulphates	-0.0126738	-0.0126742
Alcohol	0.0022014	0.0022014
LabelAppeal	0.1331963	0.1331958
AcidIndex	-0.0870512	-0.0870531
STARS	0.3112869	0.3112910

### Overdispersion Analysis Negative Binomial Model 3

For our model(3), we see that our Residual deviance is 14728 and degrees of freedom is 12780; our Residual deviance 1.15 greater than our Residual degrees of freedom, which similar to that of classical Poisson model (1) with original data which was also 1.15.

Sine we see that we have over dispersion, let's find out the dispersion parameter  $\phi$ .

```
nbpr3 <- residuals(nbmod3,"pearson")
nbphi3 <- sum(nbpr3^2)/df.residual(nbmod3)
nbphi3
```

```
## [1] 0.851477
```

The Negative Binomial dispersion parameter for Modle (3) is 0.851477 which is similar to that of the classical Poisson Model (1). Hence theta value of the of the Negative binomial has not had much impact in improving in having the variance approximates to the mean.

### zero-inflation model Negative Binomial Model 3

Next we will proceed with the Negative Binomial zero-inflation model as it is another very common occurrence when working with count data using original data.

```
library(sandwich)
library(msm)
library(pscl)
nbmod3zip <- zeroinfl(TARGET~ ., data = winedata_orig, dist = "negbin")
summary(nbmod3zip)
```

```
##
## Call:
## zeroinfl(formula = TARGET ~ ., data = winedata_orig, dist = "negbin")
##
## Pearson residuals:
##      Min      1Q   Median      3Q      Max
```

```
## -2.122603 -0.404876 -0.007536 0.371265 5.768511
##
## Count model coefficients (negbin with log link):
##           Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.444e+00  2.020e-01   7.147 8.87e-13 ***
## FixedAcidity    3.379e-04  8.420e-04   0.401 0.688183
## VolatileAcidity -1.211e-02  6.721e-03  -1.801 0.071627 .
## CitricAcid      4.923e-04  6.024e-03   0.082 0.934864
## ResidualSugar   -7.703e-05  1.586e-04  -0.486 0.627289
## Chlorides       -2.241e-02  1.691e-02  -1.325 0.185111
## FreeSulfurDioxide 2.546e-05  3.547e-05   0.718 0.472894
## TotalSulfurDioxide -1.783e-05  2.265e-05  -0.787 0.431029
## Density         -2.847e-01  1.983e-01  -1.436 0.151124
## pH              5.929e-03  7.859e-03   0.754 0.450588
## Sulphates       1.728e-04  5.919e-03   0.029 0.976714
## Alcohol         6.886e-03  1.440e-03   4.784 1.72e-06 ***
## LabelAppeal     2.330e-01  6.303e-03  36.962 < 2e-16 ***
## AcidIndex       -1.858e-02  4.898e-03  -3.794 0.000148 ***
## STARS           1.009e-01  5.201e-03  19.403 < 2e-16 ***
## Log(theta)      1.696e+01  2.724e+00   6.227 4.75e-10 ***
##
## Zero-inflation model coefficients (binomial with logit link):
##           Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -4.4381583  1.3373702  -3.319 0.000905 ***
## FixedAcidity    0.0007553  0.0055468   0.136 0.891681
## VolatileAcidity  0.1937171  0.0438506   4.418 9.98e-06 ***
## CitricAcid     -0.0296094  0.0399708  -0.741 0.458829
## ResidualSugar   -0.0011762  0.0010429  -1.128 0.259390
## Chlorides       0.0921622  0.1093477   0.843 0.399320
## FreeSulfurDioxide -0.0007420  0.0002422  -3.063 0.002188 **
## TotalSulfurDioxide -0.0009866  0.0001523  -6.476 9.39e-11 ***
## Density         0.4801296  1.3159245   0.365 0.715215
## pH              0.2160267  0.0512199   4.218 2.47e-05 ***
## Sulphates       0.1323368  0.0387665   3.414 0.000641 ***
## Alcohol         0.0279102  0.0095780   2.914 0.003568 **
## LabelAppeal     0.7229464  0.0429458  16.834 < 2e-16 ***
## AcidIndex       0.4347283  0.0258382  16.825 < 2e-16 ***
## STARS          -2.3767989  0.0603130 -39.408 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Theta = 23249933.4907
## Number of iterations in BFGS optimization: 59
## Log-likelihood: -2.041e+04 on 31 Df

nbzpr3 <- residuals(nbmod3zip,"pearson")
nbzphi3 <- sum(nbzpr3^2)/df.residual(nbmod3zip)
nbzphi3

## [1] 0.4637071
```

Note that the zip model output above does not indicate in any way if our zero-inflated model is an improvement over a standard Negative Binomial regression. We can determine this by running the corresponding standard Negative Binomial model and then performing a Vuong test of the two models.

```
vuong(nbmod3zip,nbmod3)
```

```
## Vuong Non-Nested Hypothesis Test-Statistic:
## (test-statistic is asymptotically distributed N(0,1) under the
## null that the models are indistinguishable)
## -----
##              Vuong z-statistic              H_A      p-value
## Raw              47.98803 model1 > model2 < 2.22e-16
## AIC-corrected    47.74231 model1 > model2 < 2.22e-16
## BIC-corrected    46.82618 model1 > model2 < 2.22e-16
```

The Vuong test suggests that the zero-inflated Negative Binomial model is slight improvement over a standard Negative Binomial model. Please note that The model1 from the `vuong()` function output in this case refers to the first argument in our `vuong(mod3zip,nbmod3)` function which is the zero-inflation model Negative Binomial Model (3)

### 3.2.1 Negative Binomial model 4

In this model we will be using the basic Negative Binomial model; however using transformed data.

```
#transformed data. Negative Binomial model 4
```

```
nbmod4 = glm.nb(TARGET ~ ., data = winedata_trans)
```

```
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
```

```
summary(nbmod4)
```

```
##
## Call:
## glm.nb(formula = TARGET ~ ., data = winedata_trans, init.theta = 36223.3581,
## link = log)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.9183  -0.8511  -0.0110   0.5226   4.0824
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    2.570e+00  2.002e-01  12.841  < 2e-16 ***
## ResidualSugar_MISS  2.283e-02  2.341e-02   0.976  0.32925
## Chlorides_MISS     3.017e-03  2.330e-02   0.129  0.89697
## FreeSulfurDioxide_MISS  2.300e-02  2.366e-02   0.972  0.33102
## TotalSulfurDioxide_MISS  1.883e-02  2.246e-02   0.838  0.40176
## pH_MISS          -3.496e-02  2.991e-02  -1.169  0.24257
## Sulphates_MISS    -6.759e-03  1.757e-02  -0.385  0.70051
## Alcohol_MISS      2.136e-02  2.306e-02   0.926  0.35436
```

```
## STARS_MISS -1.471e+00 2.371e-02 -62.036 < 2e-16 ***
## FixedAcidity_CAP -5.713e-04 9.179e-04 -0.622 0.53371
## VolatileAcidity_CAP -3.550e-02 7.248e-03 -4.898 9.67e-07 ***
## CitricAcid_CAP 7.431e-03 6.527e-03 1.138 0.25493
## ResidualSugar_CAP 1.348e-04 1.538e-04 0.876 0.38090
## Chlorides_CAP -2.664e-02 1.618e-02 -1.646 0.09978 .
## FreeSulfurDioxide_CAP 1.600e-04 5.266e-05 3.039 0.00237 **
## TotalSulfurDioxide_CAP 8.381e-05 2.599e-05 3.224 0.00126 **
## Density_CAP -2.848e-01 1.946e-01 -1.463 0.14334
## pH_CAP -1.361e-02 8.673e-03 -1.569 0.11665
## Sulphates_CAP -1.194e-02 5.908e-03 -2.020 0.04333 *
## Alcohol_CAP 3.956e-03 1.646e-03 2.404 0.01623 *
## AcidIndex_CAP -7.801e-02 5.259e-03 -14.834 < 2e-16 ***
## LabelAppeal_Positive -2.560e-02 1.855e-02 -1.380 0.16746
## STARS_1 -7.179e-01 2.081e-02 -34.501 < 2e-16 ***
## STARS_2 -3.427e-01 1.944e-02 -17.627 < 2e-16 ***
## STARS_3 -1.734e-01 2.006e-02 -8.645 < 2e-16 ***
## STARS_4 NA NA NA NA
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(36223.36) family taken to be 1)
##
## Null deviance: 22860 on 12794 degrees of freedom
## Residual deviance: 14375 on 12770 degrees of freedom
## AIC: 46370
##
## Number of Fisher Scoring iterations: 1
##
##
## Theta: 36223
## Std. Err.: 31421
## Warning while fitting theta: iteration limit reached
##
## 2 x log-likelihood: -46318.31
```

\*\*\* Interpretation Negative Binomial Model 4\*\*\*

As per the below table, even for transformed daata, it is worth noting that the classical Poisson Coefficients are similar to that of the Negative Binomial's for teh same reason as was teh case for original data. Pease refer to “Interpretation Negative Binomial Model 3” for more details.

In addition, the Negative Binomial model with transformed data has an improved AIC of 46370, as it is lower than the Negative Binomial model 3 AIC (46703); which was run against the original data.

```
kable(rbind(data.frame("Poisson Coeff"= poismod2$coefficients,"Negative Binom Coeffi" = nbmod4$coefficients
```

	Poisson.Coeff	Negative.Binom.Coeffi
(Intercept)	2.5701252	2.5701601
ResidualSugar_MISS	0.0228341	0.0228344
Chlorides_MISS	0.0030173	0.0030168

	Poisson.Coeff	Negative.Binom.Coeffi
FreeSulfurDioxide_MISS	0.0230001	0.0230007
TotalSulfurDioxide_MISS	0.0188307	0.0188313
pH_MISS	-0.0349529	-0.0349554
Sulphates_MISS	-0.0067580	-0.0067590
Alcohol_MISS	0.0213581	0.0213583
STARS_MISS	-1.4710696	-1.4710700
FixedAcidity_CAP	-0.0005712	-0.0005713
VolatileAcidity_CAP	-0.0355011	-0.0355022
CitricAcid_CAP	0.0074304	0.0074305
ResidualSugar_CAP	0.0001348	0.0001348
Chlorides_CAP	-0.0266371	-0.0266378
FreeSulfurDioxide_CAP	0.0001600	0.0001600
TotalSulfurDioxide_CAP	0.0000838	0.0000838
Density_CAP	-0.2847644	-0.2847684
pH_CAP	-0.0136064	-0.0136077
Sulphates_CAP	-0.0119359	-0.0119366
Alcohol_CAP	0.0039558	0.0039557
AcidIndex_CAP	-0.0780062	-0.0780093
LabelAppeal_Positive	-0.0255998	-0.0256008
STARS_1	-0.7179018	-0.7179026
STARS_2	-0.3426734	-0.3426738
STARS_3	-0.1733976	-0.1733981
STARS_4	NA	NA

#### *Overdispersion Analysis Negative Binomial Model 4*

For our model(4), we see that our Residual deviance is 14375 and degrees of freedom is 12770; our Residual deviance 1.12 greater than our Residual degrees of freedom, which is similar to that of classical Poisson model (1) with transformed data which was also 1.12.

Sine we see that we have over dispersion, let's find out the dispersion parameter  $\phi$ .

```
nbpr4 <- residuals(nbmod4,"pearson")
nbphi4 <- sum(nbpr4^2)/df.residual(nbmod4)
nbphi4
```

```
## [1] 0.9667395
```

Our dispersion parameter for Modle (4) is 0.9667395 which is much closer to 1 than the dispersion parameter of our Modle (3). However, it is slightly lower than of the classical Poisson model using transformed data.

#### **zero-inflation model Negative Binomial Model 4**

Next we will proceed with the Negative Binomial zero-inflation model as it is another very common occurrence when working with count data using transformed data.

```
library(sandwich)
library(msm)
library(pscl)
```

```

quine4 <- as.data.frame(model.matrix(nbmod4)) ## all regressors
quine4 <- quine4[, !is.na(coef(nbmod4))]      ## only identified
quine4 <- quine4[, -1]                       ## omit intercept
quine4$TARGET <- winedata_trans$TARGET      ## add response

## re-fit glm.nb()
fm1a <- glm.nb(TARGET ~ ., data = quine4)

```

```

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached

```

```

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached

```

```

## equivalent to previous fit
logLik(fm1a) - logLik(nbmod4)

```

```

## 'log Lik.' 0 (df=26)

```

```

coef(fm1a) - na.omit(coef(nbmod4))

```

```

##           (Intercept)      ResidualSugar_MISS      Chlorides_MISS
##                0                0                0
## FreeSulfurDioxide_MISS TotalSulfurDioxide_MISS      pH_MISS
##                0                0                0
##      Sulphates_MISS      Alcohol_MISS      STARS_MISS
##                0                0                0
##      FixedAcidity_CAP      VolatileAcidity_CAP      CitricAcid_CAP
##                0                0                0
##      ResidualSugar_CAP      Chlorides_CAP      FreeSulfurDioxide_CAP
##                0                0                0
## TotalSulfurDioxide_CAP      Density_CAP      pH_CAP
##                0                0                0
##      Sulphates_CAP      Alcohol_CAP      AcidIndex_CAP
##                0                0                0
##      LabelAppeal_Positive      STARS_1      STARS_2
##                0                0                0
##      STARS_3
##                0
## attr(,"na.action")
## STARS_4
##      26
## attr(,"class")
## [1] "omit"

```

```

## fit zeroinfl(), now works
nbmod4zip<- zeroinfl(TARGET ~ . | 1, data = quine4, dist = "negbin")

```

```

## Warning in sqrt(diag(vc)[np]): NaNs produced

```

```
summary(nbmod4zip)
```

```
##
## Call:
## zeroinfl(formula = TARGET ~ . | 1, data = quine4, dist = "negbin")
##
## Pearson residuals:
##      Min      1Q   Median      3Q      Max
## -1.80479 -0.61159  0.06037  0.53998  5.74814
##
## Count model coefficients (negbin with log link):
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    2.474e+00  2.059e-01  12.014 < 2e-16 ***
## ResidualSugar_MISS  2.185e-02  2.400e-02   0.910  0.36263
## Chlorides_MISS    7.386e-03  2.394e-02   0.309  0.75767
## FreeSulfurDioxide_MISS  2.010e-02  2.421e-02   0.830  0.40654
## TotalSulfurDioxide_MISS  2.352e-02  2.306e-02   1.020  0.30768
## pH_MISS        -2.769e-02  3.076e-02  -0.900  0.36809
## Sulphates_MISS  -6.193e-03  1.805e-02  -0.343  0.73158
## Alcohol_MISS    1.700e-02  2.363e-02   0.719  0.47187
## STARS_MISS      -1.360e+00  2.627e-02 -51.785 < 2e-16 ***
## FixedAcidity_CAP -4.503e-04  9.424e-04  -0.478  0.63275
## VolatileAcidity_CAP -3.032e-02  7.460e-03  -4.064  4.82e-05 ***
## CitricAcid_CAP   5.586e-03  6.699e-03   0.834  0.40439
## ResidualSugar_CAP  7.936e-05  1.577e-04   0.503  0.61474
## Chlorides_CAP    -2.117e-02  1.659e-02  -1.276  0.20197
## FreeSulfurDioxide_CAP  1.529e-04  5.382e-05   2.841  0.00449 **
## TotalSulfurDioxide_CAP  5.926e-05  2.641e-05   2.244  0.02483 *
## Density_CAP      -2.951e-01  2.000e-01  -1.475  0.14012
## pH_CAP           -8.061e-03  8.914e-03  -0.904  0.36582
## Sulphates_CAP    -9.398e-03  6.070e-03  -1.548  0.12161
## Alcohol_CAP      4.776e-03  1.687e-03   2.831  0.00464 **
## AcidIndex_CAP    -6.702e-02  5.560e-03 -12.054 < 2e-16 ***
## LabelAppeal_Positive -2.721e-02  1.903e-02  -1.430  0.15274
## STARS_1          -6.211e-01  2.191e-02 -28.347 < 2e-16 ***
## STARS_2          -3.267e-01  1.948e-02 -16.773 < 2e-16 ***
## STARS_3          -1.730e-01  2.006e-02  -8.625 < 2e-16 ***
## Log(theta)       1.725e+01      NA      NA      NA
##
## Zero-inflation model coefficients (binomial with logit link):
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -2.80654    0.08463 -33.16 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Theta = 31167640.5163
## Number of iterations in BFGS optimization: 56
## Log-likelihood: -2.306e+04 on 27 Df
```

```
#####
```

```
nbzpr4 <- residuals(nbmod4zip,"pearson")
nbzphi4 <- sum(nbzpr4^2)/df.residual(nbmod4zip)
nbzphi4
```

```
## [1] 0.8386927
```

Again, Please note that the zip model output above does not indicate in any way if our zero-inflated model is an improvement over a standard Negative Binomial regression. We can determine this by running the corresponding standard Negative Binomial model and then performing a Vuong test of the two models against the transformed data.

```
vuong(nbmod4zip,nbmod4)
```

```
## Vuong Non-Nested Hypothesis Test-Statistic:
## (test-statistic is asymptotically distributed N(0,1) under the
## null that the models are indistinguishable)
## -----
##              Vuong z-statistic              H_A      p-value
## Raw              6.163416 model1 > model2 3.5596e-10
## AIC-corrected    6.163416 model1 > model2 3.5596e-10
## BIC-corrected    6.163416 model1 > model2 3.5596e-10
```

The Vuong test suggests that the zero-inflated Negative Binomial model is slight improvement over a standard Negative Binomial model using the transformed data. Please note that The model1 from the vuong() function output in this case refers to the first argument in our vuong(mod4zip,nbmod4) function which is the zero-inflation model Negative Binomial Model (4)

```
library(AICcmodavg)
```

```
## Warning: package 'AICcmodavg' was built under R version 3.3.1
```

```
#AICc(list(fm_qpois))
AICc(fm_qpois, return.K = FALSE, second.ord = TRUE,nobs = NULL, c.hat = 1)
```

```
## [1] NA
```

### 3.3 Linear Regression models

Although it is highly recommended for continuous variables instead of count variables, we will also create two linear regression models.

#### 3.3.1 Linear Regression Model 5

We will explore the Linear models Using original data with replacing all missing data with the means.

```
##
## Call:
## lm(formula = TARGET ~ ., data = winedata_orig)
```



```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.5476 -0.9475  0.0669  0.9047  5.9903
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    3.986e+00  4.487e-01   8.883  < 2e-16 ***
## FixedAcidity    1.608e-06  1.884e-03   0.001  0.999319
## VolatileAcidity -9.923e-02  1.498e-02  -6.625  3.61e-11 ***
## CitricAcid      2.085e-02  1.362e-02   1.531  0.125804
## ResidualSugar   2.012e-04  3.559e-04   0.565  0.571860
## Chlorides      -1.243e-01  3.777e-02  -3.290  0.001003 **
## FreeSulfurDioxide 3.153e-04  8.093e-05   3.897  9.80e-05 ***
## TotalSulfurDioxide 2.264e-04  5.201e-05   4.353  1.35e-05 ***
## Density        -8.012e-01  4.419e-01  -1.813  0.069829 .
## pH             -3.453e-02  1.754e-02  -1.969  0.049012 *
## Sulphates      -3.271e-02  1.322e-02  -2.475  0.013352 *
## Alcohol         1.094e-02  3.234e-03   3.384  0.000717 ***
## LabelAppeal     4.326e-01  1.367e-02  31.654  < 2e-16 ***
## AcidIndex      -2.084e-01  9.212e-03 -22.619  < 2e-16 ***
## STARS           9.767e-01  1.045e-02  93.433  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.324 on 12780 degrees of freedom
## Multiple R-squared:  0.528, Adjusted R-squared:  0.5275
## F-statistic: 1021 on 14 and 12780 DF, p-value: < 2.2e-16
```

\*\*\* Interpretation of Linear Model 5\*\*\*

Based on the summary for Linear Model 5, below are the characteristics :

- The Residual standard error is 1.3242
- Multiple R-squared: 0.528
- Adjusted R-squared: 0.5275
- F-statistic: 1021 on 14 and 12780 DF
- p-value: < 2.2e-16

Based on the available coefficients, we can make the following observations:

- Positive Impact - The following variables have a positive impact on TARGET, meaning an increase in the values of these variables leads to an increase in the number of cases sold: STARS, LabelAppeal, Alcohol, TotalSulfurDioxide, FreeSulfurDioxide, ResidualSugar, CitricAcid, FixedAcidity
- Negative Impact - The following variables have a negative impact on TARGET, meaning an increase in the values of these variables leads to a decrease in the number of cases sold: AcidIndex, Sulphates, pH, Density, Chlorides, VolatileAcidity
- The following variables have a 'significant' impact. These are the more important predictors for TARGET: STARS, AcidIndex, LabelAppeal, Alcohol, Sulphates, pH, TotalSulfurDioxide, FreeSulfurDioxide, Chlorides, VolatileAcidity
- Finally, the Linear Model equation is given by the following:

$3.9861 + 2e-06 * \text{FixedAcidity} - 0.099232 * \text{VolatileAcidity} + 0.020854 * \text{CitricAcid} + 0.000201 * \text{ResidualSugar}$   
 $- 0.124266 * \text{Chlorides} + 0.000315 * \text{FreeSulfurDioxide} + 0.000226 * \text{TotalSulfurDioxide} - 0.801199 * \text{Density}$   
 $- 0.034527 * \text{pH} - 0.032707 * \text{Sulphates} + 0.010942 * \text{Alcohol} + 0.432607 * \text{LabelAppeal} - 0.208371 * \text{AcidIndex}$   
 $+ 0.976721 * \text{STARS}$

### 3.3.1 Linear Regression Model 6

In this model we will be using the Linear Regression model; however using transformed data.

```
##
## Call:
## lm(formula = TARGET ~ ., data = winedata_trans)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.1375 -0.9450  0.0246  0.9372  6.7449
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      7.938e+00  4.779e-01  16.611 < 2e-16 ***
## ResidualSugar_MISS  6.293e-02  5.649e-02   1.114 0.265301
## Chlorides_MISS     6.208e-03  5.560e-02   0.112 0.911102
## FreeSulfurDioxide_MISS 6.444e-02  5.624e-02   1.146 0.251917
## TotalSulfurDioxide_MISS 4.998e-02  5.383e-02   0.929 0.353164
## pH_MISS           -8.562e-02  6.990e-02  -1.225 0.220634
## Sulphates_MISS    -2.325e-02  4.132e-02  -0.563 0.573704
## Alcohol_MISS       6.140e-02  5.494e-02   1.118 0.263709
## STARS_MISS         -4.092e+00  6.051e-02 -67.620 < 2e-16 ***
## FixedAcidity_CAP   -1.190e-03  2.175e-03  -0.547 0.584328
## VolatileAcidity_CAP -1.065e-01  1.720e-02  -6.193 6.07e-10 ***
## CitricAcid_CAP      2.220e-02  1.554e-02   1.429 0.153062
## ResidualSugar_CAP    3.782e-04  3.646e-04   1.038 0.299517
## Chlorides_CAP       -7.754e-02  3.840e-02  -2.020 0.043447 *
## FreeSulfurDioxide_CAP 4.803e-04  1.261e-04   3.809 0.000140 ***
## TotalSulfurDioxide_CAP 2.303e-04  6.162e-05   3.737 0.000187 ***
## Density_CAP        -9.171e-01  4.642e-01  -1.976 0.048215 *
## pH_CAP             -3.814e-02  2.061e-02  -1.850 0.064286 .
## Sulphates_CAP      -3.372e-02  1.406e-02  -2.397 0.016538 *
## Alcohol_CAP         1.311e-02  3.894e-03   3.367 0.000761 ***
## AcidIndex_CAP      -2.108e-01  1.177e-02 -17.916 < 2e-16 ***
## LabelAppeal_Positive -7.695e-02  4.394e-02  -1.751 0.079931 .
## STARS_1            -2.770e+00  6.074e-02 -45.611 < 2e-16 ***
## STARS_2            -1.586e+00  5.992e-02 -26.462 < 2e-16 ***
## STARS_3            -8.683e-01  6.248e-02 -13.898 < 2e-16 ***
## STARS_4              NA         NA      NA      NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.367 on 12770 degrees of freedom
## Multiple R-squared:  0.4976, Adjusted R-squared:  0.4967
## F-statistic:  527 on 24 and 12770 DF,  p-value: < 2.2e-16
```

### \*\*\* Interpretation of Linear Model 6\*\*\*

Based on the summary for Linear Model 6, below are the characteristics :

- The Residual standard error is 1.3667
- Multiple R-squared: 0.4976
- Adjusted R-squared: 0.4967
- F-statistic: 527 on 24 and 12770 DF
- p-value: < 2.2e-16

Based on the available coefficients, we can make the following observations:

- Positive Impact - The following variables have a positive impact on TARGET, meaning an increase in the values of these variables leads to an increase in the number of cases sold: Alcohol\_CAP, TotalSulfurDioxide\_CAP, FreeSulfurDioxide\_CAP, ResidualSugar\_CAP, CitricAcid\_CAP, Alcohol\_MISS, TotalSulfurDioxide\_MISS, FreeSulfurDioxide\_MISS, Chlorides\_MISS, ResidualSugar\_MISS
- Negative Impact - The following variables have a negative impact on TARGET, meaning an increase in the values of these variables leads to a decrease in the number of cases sold: STARS\_3, STARS\_2, STARS\_1, LabelAppeal\_Positive, AcidIndex\_CAP, Sulphates\_CAP, pH\_CAP, Density\_CAP, Chlorides\_CAP, VolatileAcidity\_CAP, FixedAcidity\_CAP, STARS\_MISS, Sulphates\_MISS, pH\_MISS
- The following variables have a 'significant' impact. These are the more important predictors for TARGET: STARS\_3, STARS\_2, STARS\_1, AcidIndex\_CAP, Alcohol\_CAP, Sulphates\_CAP, Density\_CAP, TotalSulfurDioxide\_CAP, FreeSulfurDioxide\_CAP, Chlorides\_CAP, VolatileAcidity\_CAP, STARS\_MISS
- Finally, the Linear Model equation is given by the following:

$$7.938 + 0.062931 * \text{ResidualSugar\_MISS} + 0.006208 * \text{Chlorides\_MISS} + 0.064437 * \text{FreeSulfurDioxide\_MISS} + 0.049984 * \text{TotalSulfurDioxide\_MISS} - 0.085625 * \text{pH\_MISS} - 0.023249 * \text{Sulphates\_MISS} + 0.061402 * \text{Alcohol\_MISS} - 4.092034 * \text{STARS\_MISS} - 0.00119 * \text{FixedAcidity\_CAP} - 0.106543 * \text{VolatileAcidity\_CAP} + 0.022202 * \text{CitricAcid\_CAP} + 0.000378 * \text{ResidualSugar\_CAP} - 0.077542 * \text{Chlorides\_CAP} + 0.00048 * \text{FreeSulfurDioxide\_CAP} + 0.00023 * \text{TotalSulfurDioxide\_CAP} - 0.917053 * \text{Density\_CAP} - 0.038139 * \text{pH\_CAP} - 0.033715 * \text{Sulphates\_CAP} + 0.013111 * \text{Alcohol\_CAP} - 0.210836 * \text{AcidIndex\_CAP} - 0.07695 * \text{LabelAppeal\_Positive} - 2.770326 * \text{STARS\_1} - 1.585505 * \text{STARS\_2} - 0.868345 * \text{STARS\_3}$$

```
AIC(step5)
```

```
## [1] 43508.94
```

```
AIC(step6)
```

```
## [1] 44321.76
```

```
lmpr5 <- residuals(lmod5,"pearson")
lmphi5 <- sum(lmpr5^2)/df.residual(lmod5)
lmphi5
```

```
## [1] 1.753383
```

```
lmpr6 <- residuals(lmod6,"pearson")
lmphi6 <- sum(lmpr6^2)/df.residual(lmod6)
lmphi6
```

```
## [1] 1.867863
```

## 4 Model Selection

Before we proceed with our model selection, let take a quick look at our models inventory. We have 12 models using a combination of three different type distributions. First we created our models using GLM distribution; then we created few using the zero Augmented distribution, and finally the Linear distribution. Hence our models selection will be based on the best AIC/  $\phi$  =Dispersion parameter for the GLM, AIC for Linear regression; and Vuong test for the zero Augmented distribution. Below is summary table of model selection strategy:

```
modselect<- read.csv("https://raw.githubusercontent.com/kishkp/data621-ctg5/master/HW5/modelselection2.csv")
kable(modselect, caption = "Model Selection Strategy")
```

Table 9: Model Selection Strategy

Distribution.Type	Model.Description	Comparaision.KPI
Classical Poisson	Poisson using original data	AIC
	Poisson using Transformed data	AIC
Quasi-Poisson	Quasi Poisson using original data	$\phi$ =Dispersion parameter
	Quasi Poisson using transformed data	$\phi$ =Dispersion parameter
Negative Binomial	NB using original data	AIC
	NB using transformed data	AIC
zero-inflation Poisson	zero inflated Pois using original data	Vuong test
	zero inflated Pois using Transforemed data	Vuong test
zero-inflation NB	zero inflated NB using original data	Vuong test
	zero inflated NB using transformed data	Vuong test
LM	linear regression using original data	AIC
	linear regression using transformed data	AIC

Below in the Model Selection KPI table is a summary of the major indicators use to select the best fit. To selefct the best model we will be using a combination of the AIC, Dispersion parameter, as well as the Vuong closeness test specifically for the zero inflation distributions.

However, since our data is count data and the problem of dispersion occurs more frequently in count data set, we will be using Dispersion parameter first in our process elimination, followed by AIC, and Voung test. Hence, the “Model Selection KPI” table nelow is sorted using the Dispersion parameter.

```
modmetrics<- read.csv("https://raw.githubusercontent.com/kishkp/data621-ctg5/master/HW5/modelmetrics2.csv")
kable(modmetrics, caption = "Model Selection KPI")
```

Table 10: Model Selection KPI

Model.Type	AIC	Dispersion.parameter	Vuong.Selected
Linear model with transformed data	44321.76	1.8678630	
Linear model with original data	43508.94	1.7533830	
Pois with transformed data	46368	0.9667917	
Quasi-Poisson with transformed data	Undefined	0.9667917	
Negative binomial /transformed data	46370	0.9667395	
Quasi-Poisson with Original data	Undefined	0.8515200	
Pois with original data	46700	0.8515130	
Negative binomial /original data	46703	0.8514770	
zero inflation NB with transformed data	Undefined	0.8386927	zero inflation NB transformed data
zero inflation Poisson with transformed data	Undefined	0.8386535	zero inflation Poisson transformed data
zero inflation NB with orig data	Undefined	0.4637071	zero inflation NB with orig data
zero inflation Poisson with orig data	Undefined	0.4636815	zero inflation Poisson with orig data

Therefore, from the above table, we can easily eliminate the Linear models both for in the original and transformed data as they respectively have a dispersion parameter of 1.867863 and 1.753383 which are much higher than 1.

Next we will eliminate the zero inflation Negative Binomial and Poisson for the original as they respectively have a dispersion parameter of 0.4637071 and 0.4636815 which are much lower than 1.

We will also eliminate the zero inflation Negative Binomial and Poisson for the transformed data as they respectively have a dispersion parameter of 0.8386927 and 0.8386535 which are not close to 1 compared to the rest of the models.

Also, based on dispersion parameter, we will eliminate the Poisson, Quasi-Poisson, and Negative binomial with original data as they respectively have a dispersion parameter of 0.851513, 0.85152, and 0.851477 which are not close to 1 compared to the rest of the models.

Finally we are left with the following 3 models:

Poisson with transformed data, with Dispersion parameter = 0.9667917 Quasi-Poisson with transformed data with Dispersion parameter = 0.9667917 Negative binomial /transformed data Dispersion parameter = 0.9667395

Since we have a virtual tie in the remaining 3 models from dispersion parameter perspective, we will use the second metric, AIC, as defining factor for our remaining 3 model selection. Hence, the Poisson model with transformed data as it has an AIC of 46368 compared to the Negative Binomial which is 46370.

## 5 Prediction Using Evaluation Data

Now that we have selected the final model, we will go ahead and use this model to predict the results for the evaluation dataset. After transforming the data to meet the needs of the trained model, we will apply the model.

## 5.1 Tranformation of Evaluation Data

First we need to transform the evaluation dataset to account for all the predictors that were used in the model.

## 5.2 Model Output

For ease of display we will display, in transposed format, only the first six rows as we have 42 variables.

### First six Records from output

Table 11: Model Output / Results

	1	2	3	4	5	6
IN	3.00000	21.000000	37.000000	39.0000	47.00000	62.00000
TARGET	1.00000	1.000000	1.000000	1.0000	1.00000	1.00000
FixedAcidity	5.40000	11.400000	15.900000	11.6000	3.80000	9.00000
VolatileAcidity	-0.86000	0.210000	1.190000	0.3200	0.22000	-0.21000
CitricAcid	0.27000	0.280000	1.140000	0.5500	0.31000	0.04000
ResidualSugar	-10.70000	1.200000	31.900000	-50.9000	-7.70000	51.40000
Chlorides	0.09200	0.038000	-0.299000	0.0760	0.03900	0.23700
FreeSulfurDioxide	23.00000	70.000000	115.000000	35.0000	40.00000	-213.00000
TotalSulfurDioxide	398.00000	53.000000	381.000000	83.0000	129.00000	-527.00000
Density	0.98527	1.028990	1.034160	1.0002	0.90610	0.99516
pH	5.02000	2.540000	2.990000	3.3200	4.72000	3.16000
Sulphates	0.64000	-0.070000	0.310000	2.1800	-0.64000	0.70000
Alcohol	12.30000	4.800000	11.400000	-0.5000	10.90000	14.70000
LabelAppeal	-1.00000	0.000000	1.000000	0.0000	0.00000	1.00000
AcidIndex	6.00000	10.000000	7.000000	12.0000	7.00000	10.00000
STARS	0.00000	0.000000	0.000000	0.0000	0.00000	0.00000
ResidualSugar_MISS	0.00000	0.000000	0.000000	0.0000	0.00000	0.00000
Chlorides_MISS	0.00000	0.000000	0.000000	0.0000	0.00000	0.00000
FreeSulfurDioxide_MISS	0.00000	0.000000	0.000000	0.0000	0.00000	0.00000
TotalSulfurDioxide_MISS	0.00000	0.000000	0.000000	0.0000	0.00000	0.00000
pH_MISS	0.00000	0.000000	0.000000	0.0000	0.00000	0.00000
Sulphates_MISS	0.00000	0.000000	0.000000	0.0000	0.00000	0.00000
Alcohol_MISS	0.00000	0.000000	0.000000	0.0000	0.00000	0.00000
STARS_MISS	1.00000	1.000000	1.000000	1.0000	1.00000	1.00000
FixedAcidity_CAP	5.40000	11.400000	17.500000	11.6000	3.80000	9.00000
VolatileAcidity_CAP	-1.04600	0.210000	1.190000	0.3200	0.22000	-0.21000
CitricAcid_CAP	0.27000	0.280000	1.140000	0.5500	0.31000	0.04000
ResidualSugar_CAP	-10.70000	1.200000	31.900000	-51.9000	-7.70000	61.56500
Chlorides_CAP	0.09200	0.038000	-0.479300	0.0760	0.03900	0.23700
FreeSulfurDioxide_CAP	23.00000	70.000000	115.000000	35.0000	40.00000	-216.30000
TotalSulfurDioxide_CAP	398.00000	53.000000	381.000000	83.0000	129.00000	-253.00000
Density_CAP	0.98527	1.040107	1.040107	1.0002	0.95028	0.99516
pH_CAP	4.37300	2.540000	2.990000	3.3200	4.37300	3.16000
Sulphates_CAP	0.64000	-0.070000	0.310000	2.0300	-0.99000	0.70000
Alcohol_CAP	12.30000	4.800000	11.400000	4.3000	10.90000	14.70000
AcidIndex_CAP	6.00000	10.000000	7.000000	10.0000	7.00000	10.00000
LabelAppeal_Positive	1.00000	1.000000	1.000000	1.0000	1.00000	0.00000
STARS_1	0.00000	0.000000	0.000000	0.0000	0.00000	0.00000

	1	2	3	4	5	6
STARS_2	0.00000	0.000000	0.000000	0.0000	0.00000	0.00000
STARS_3	0.00000	0.000000	0.000000	0.0000	0.00000	0.00000
STARS_4	0.00000	0.000000	0.000000	0.0000	0.00000	0.00000

### 5.3 Conclusion

After fitting multiple models using the classical Linear, classical Poisson, and the Binomial distributions using original data and transformed data, we think that the Poisson model has performed well once we have treated the outliers and missing data.

We also felt confident that the Negative Binomial would perform good as well as it has the same dispersion parameter as classical Poisson. However, the NB AIC was bit higher by .000043 which could be negligible.

In addition we felt confident that Quasi-Poisson would perform well as its dispersion parameter was .96 close to 1. However, we were not comfortable selecting the Quasi-Poisson as we could not generate the AIC value.

The zero inflation models for both Poisson and Negative yielded to promising results especially when using the Young test. However, lack of AIC and its lower dispersion parameter had made us reconsider our decision in favor of the Poisson.

Over all, we were little bit overwhelmed with analyzing about 12 models. However, we are very satisfied with our Poisson model selection especially that it had leveraged our data preparation and transformation efforts.

## Appendix A: DATA621 Homework 05 R Code