Home Work Assignment - 01

Critical Thinking Group 5
Arindam Barman
Mohamed Elmoudni
Shazia Khan
Kishore Prasad

Contents

Overview	2
Objective	2
1 Data Exploration Analysis	2
1.1 Variable identification	2
1.2 Variable Relationships	
	3
1.3 Data summary analysis	
	4
1.4 Outliers Identification	
	6

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
${\bf Free Sulfur Dioxide}$	num	predictor
${\bf Total Sulfur Dioxide}$	num	predictor
Density	num	predictor
рН	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	None
TARGET	Number of Cases Purchased None	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	
${\bf Free Sulfur Dioxide}$	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 customes 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	e Total Sulfur Dioxide of Wine	
VolatileAcidity pH	Volatile Acid content of wine 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 $SO2 = free SO_2 + bound SO_2$.

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 breifly 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.

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
${\bf Free Sulfur Dioxide}$	30.8455713	148.7145577	30.00000	30.9334877
${\bf Total Sulfur Dioxide}$	120.7142326	231.9132105	123.00000	120.8895367
Density	0.9942027	0.0265376	0.99449	0.9942130
pН	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
рН	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:

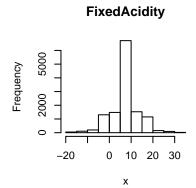
- Variable ResidualSugar has 616 and 0.0164913 correlation. Given the low correlation we will try try
 some imputation techniques to handle the missing the values and replace missing values with their
 respective value.
- variable Chlorides 638 -0.0382631 correlation. . Given the low negative correlation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

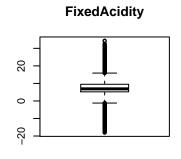
Please note that ResidualSugar, Chlorides, FreeSulfurDioxide, Alcohol, and TotalSulfurDioxide variables have similar number of missing values. They are chemically related. However, we don't think they are arithmetically related.

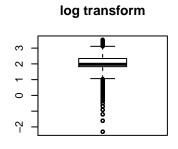
- In addition, variable pH has 395 missing values with negative correlation of -0.0094448. Again we may just ignore these missing values especially that it has very low negative correlation to the target variable.
- Variable Sulphates has much higher missing values of 1210 with low negative correlation 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

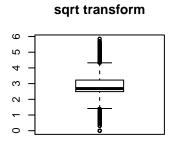
1.4 Outliers Identification

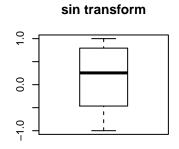
In this section we look at boxplots to determine the outliers in variables and decide on whether to act on the outliers. Lets do some univariate analysis. We will look at the Histogram and Boxplot for each variable to detect outliers if any and treat it accordingly.

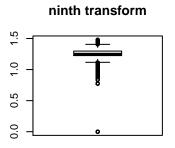


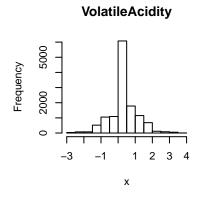


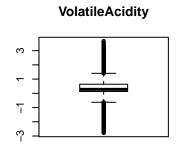


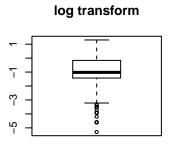


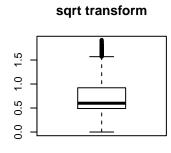


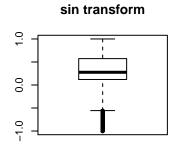


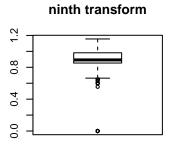




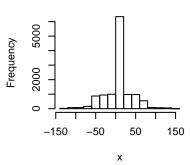




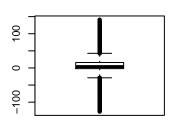




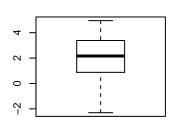
ResidualSugar



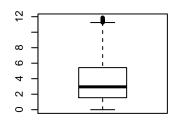
ResidualSugar



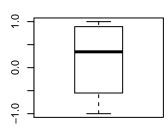
log transform



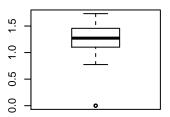
sqrt transform

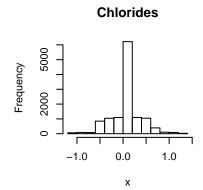


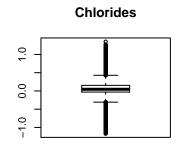
sin transform

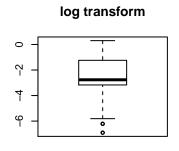


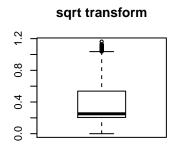
ninth transform

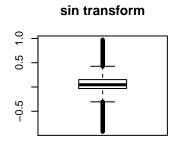


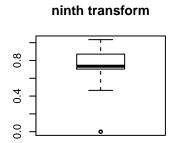




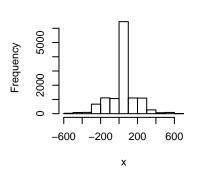




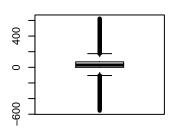




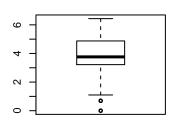
FreeSulfurDioxide



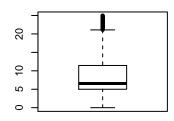
FreeSulfurDioxide



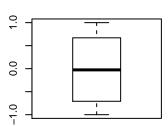
log transform



sqrt transform



sin transform



ninth transform

