

Regional Liquor Sales in Des Moines, Iowa

Christophe Hunt^{*,a}, Senthil Dhanapal^a, Yadu Chittampalli^a

^aCUNY School of Professional Studies, Data Analytics, New York, NY

Abstract

This is the abstract.

It consists of two paragraphs

Keywords: Liquor Sales, Naive Forecast*

Problem

The objective of this report is to create a statistical model for the volume sold of liquor in gallons and the profit dollars in the City of Des Moines which is within the state of Iowa. This can help us make informed decisions on inventory production, sales, and assist wholesale distributors to plan for the predicted volume of distribution.

Introduction

In February, the Distilled Spirits Council (DISCUS), announced that spirits had an estimated retail sales of nearly \$72 billion in 2015. Additionally, DISCUS credits the continuous growth of the distilled spirits industry to several key factors - continuous fascination with American Whiskeys in the United States and abroad, innovations in flavors, permutation across all spirits categories leading to consumer interest, improved regulatory and tax environment resulting in expanded market access and a relatively low number of state tax threats, and the growth of small distillers, which expanded grassroots and overall interest in the spirits category (Del Buono (2016)).

This establishes that spirit sales in the United States is a valuable market worth exploring for a more detailed and statistical understanding of sales and volume. We hope to more thoroughly understand what impact specific store sights may have accounting for the seasonal impact in November that might effect liquor sales. We will limit the analysis to the City of Des Moines for only whiskey sales in the month of November. In 2000 the State of Iowa reported sales at a record pace during the last half of 2000 (Boshart (2001)). The later part of the year has an increase in sales so planning to meet capacity is a suitable goal for any company. Our years of interest for this analysis will be the month of November for 2015 and 2016.

Research Background

The main goal that has to be achieved in inventory prediction is increasing the efficiency without decreasing the service value offered to the customers. When managing the levels of inventory, it is important to maintain moderate level(s) - not too high and not too low. If the inventory level is excessive, business funds can get wasted. These funds would not be able to be used for any other purpose, thus involving an opportunity cost.

* Authors

Email addresses: christophe.hunt@spsmail.cuny.edu (Christophe Hunt), senthil.dhanapal@spsmail.cuny.edu (Senthil Dhanapal), yadu.chittampalli@spsmail.cuny.edu (Yadu Chittampalli)

The costs of shortage, handling insurance, recording and inspection would proportionately increase along with inventory volume, thus impairing profitability.

On the other hand, low level(s) of inventory may result in frequent interruptions in the production schedule resulting in under-utilization of capacity and lower sales. When making predictions about orders that should be placed, assumptions are made as follows - uncertainty always exists regardless of the method(s) used, new technologies cannot always be forecasted for which paradigms do not exist, and social policy will be formulated where the future would be affected, changing the accuracy of the forecast (David S. Walonick (1993)).

One useful method for predicting inventories is the extrapolation of trends. In this method, trends and cycles in the historical data are examined and mathematical techniques are used to extrapolate to the future. The model chosen for forecasting would depend on the historical data (David S. Walonick (1993)).

One of the most common models used in this method is decomposition, where historical data is separated into trend, seasonal, and random components. As a result, forecasts are produced using "turning point analysis". Other examples of models used are adaptive filtering, Box-Jenkins analysis, simple linear regression, curve fitting, and weighted smoothing (David S. Walonick (1993)).

According to Makridakis, "Judgmental forecasting is superior to mathematical models, although there are several forecasting applications where computer-generated models would be more feasible." When inventory levels for bulk-quantity items would need to be forecasted monthly by large manufacturing companies, generating models through computer software would be more efficient (David S. Walonick (1993)).

Forecasting the demand of a product is very essential in predicting the order quantity. As a result, a data bank is created, helping the decision makers settle targets, create plans, and demonstrate changes in the business setting.

Two different methods are utilized in the investigation of the future demand - quantitative and qualitative. In quantitative methods, mathematical consistencies in the history are searched for. Two subcategories exist in quantitative methods - time series models and correlation models. On the other hand, qualitative methods are based on the opinions that people have had about the product in the past based on their experiences, premonitions, and emotions (Kumar (2012)).

However, when the most suitable forecast model gets selected, it is not necessarily based solely on quantitative or qualitative variables. The forecast model can even combine several models (Kumar (2012)).

Methodology

Our initial data set is sufficiently large in that it includes sales by individual stores and the invoices for each store. The reason for the large size of the initial data set is due to it including every liquor transaction from 2012 to present in Iowa, so it approaches 2.68 GB. For the purposes of this analysis, to analyze a data set this large is not feasible. Therefore, we reduced the number of variables and summarized to a regional aggregate.

Additionally, we looked into the top 10 liquor categories for each year by number of bottles sold. In 2015, the top categories were American Cocktails, Blended Whiskies, Canadian Whiskies, Imported Vodka, Puerto Rico & Virgin Islands Rum, Spiced Rum, Straight Bourbon Whiskies, Tequila, Vodka 80 Proof, and Whiskey Liqueur. Interestingly straight bourbon appears to have more sales in 2015 than 2014 which coincides with the literature of strong growing whiskey sales for every whiskey segment (Anonymous (2016)). We decided to focus on whiskey due to its strong sales and growing interest in the US.

We initially attempted to model for a dependent variable of `Volume Sold in Gallons` by the top Counties and top liquor categories, however, the distribution of this variable becomes over-dispersed and negatively skewed when aggregating the data set and therefore we were unsuccessful in modeling this variable across regions. Our hope was that if we could model for the gallons sold by entirety of Iowa we could more accurately predict our planned inventory and anticipate production goals. We were more successful in limiting our regional analysis to one city and the stores within that city.

We accomplished this by looking into volume of sales by the largest County for Iowa which is Polk county. The City of Des Moines has the largest volume of whiskey sales in Polk county so we limited our analysis to this city.

Therefore, our final evaluation data set is the following subset of variables for largest city in Iowa of Des Moines as follows; Vendor Name, Pack (pack size of bottles sold) Bottle Volume, State Bottle Cost, State Bottle Retail, Sales Dollars (Total sales), and our dependent variables are Volume sold in Gallons and Profit Dollars.

By modeling Volume sold in Gallons and Profit Dollars we aim to predict the volume of production needed and the possible profit dollars when producing at the predicted volume. We first began by using linear regression to model the Volume Sold in Gallons, however, the distribution was initially non-normal so we used the BoxCox method to transform our dependent variable to a more normal distribution. We then used forward selection method to determine the final form for the model of Volume Sold in Gallons. We further removed points that had unusually high influence on the model.

We then modeled our second dependent variable of Profit Dollars, which is the difference between the sale price and retail price, by using linear regression. After using forward selection and removing highly influence points in the second model we were able more accurately model the data set.

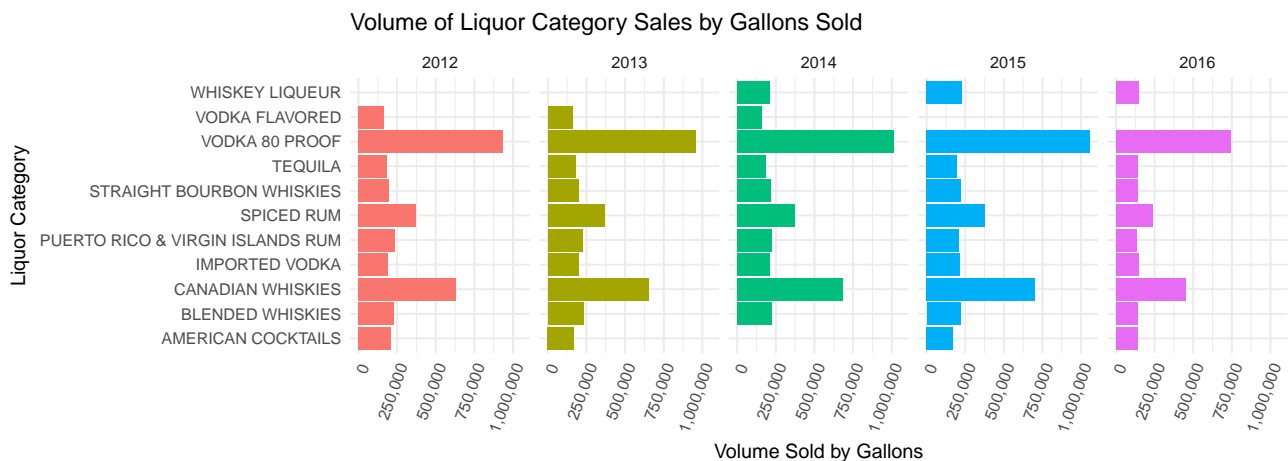
Experimentation and Results

Data Acquisition

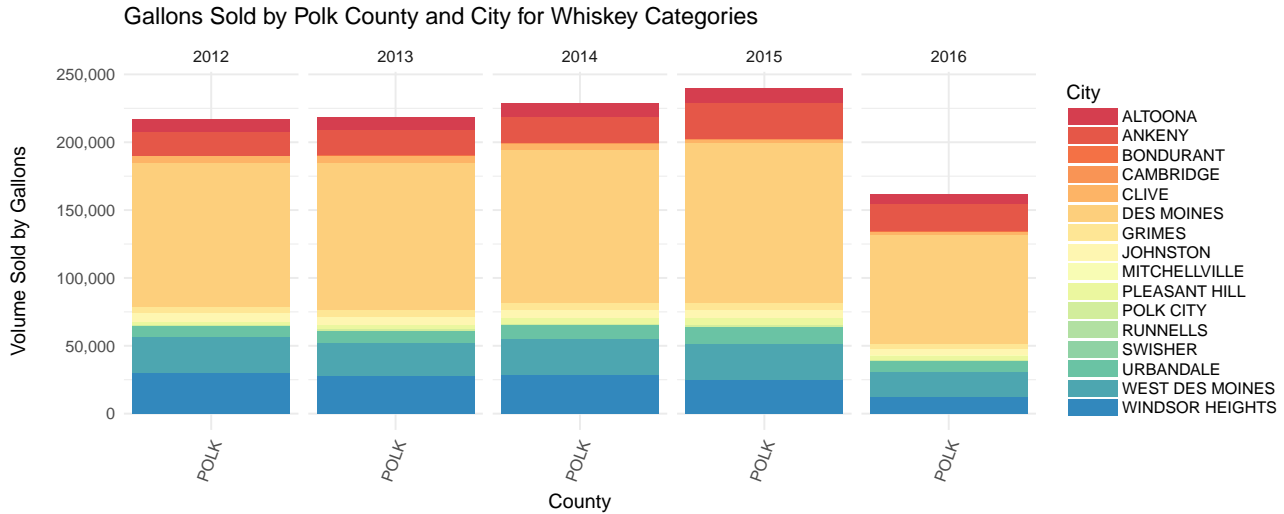
The data set contains the spirits purchase information of Iowa Class “E” liquor licensees by product and date of purchase from January 1, 2012 to current. The data set is provided by the Iowa Department of Commerce, Alcoholic Beverages Division, [click here](#) to view the data set at Data.Iowa.Gov.

As previously discussed, the data set is 2.68 GB in total size and much to large to use in a meaningful model.

We reviewed the liquor sales by gallons sold per year by Liquor Category. Initially, we viewed the top 5 Liquor Categories by volume sold but there were large disparities between years, suggesting that the top 5 change often and is likely due to changing consumer tastes. We do see a more stable set of liquor categories for the top 10 category which suggests that while tastes may change we don't see large movements in liquor categories at this level. We focused our attention on the whiskey categories.



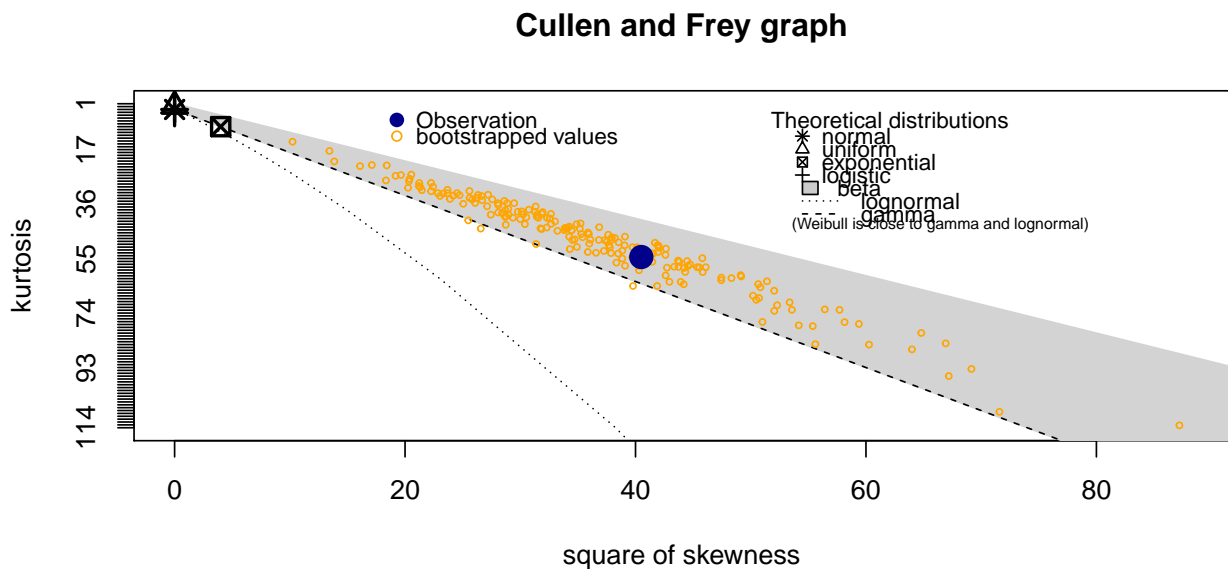
We can further see that Des Moines accounts for a significant portion of the liquor sales in Polk County. Polk County is the most populous county in Iowa so we will limit our analysis to this city.



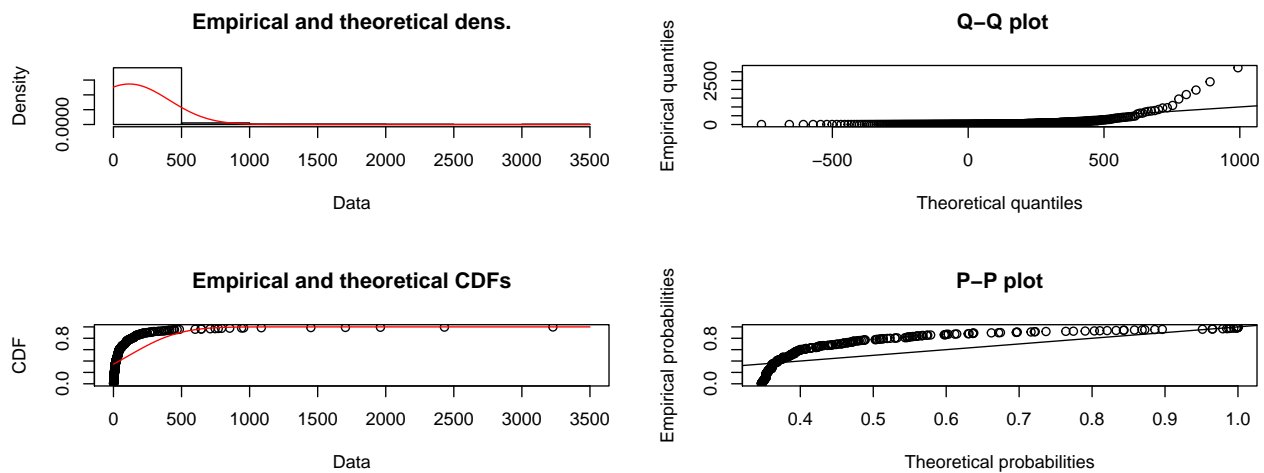
Our first attempt was to use a Poisson regression due to the over-dispersion created by aggregating the data set. However, the distribution was far to negatively skewed to fit a Poisson distribution.

Distributions of Dependent Variables

We see that the dependent variable of Bottles Sold has a non-normal distribution as illustrated in the Cullen and Frey graph AC and HC (1999) from the `fitdistrplus` package and the non normal distribution is further highlighted in the diagnostic plots.

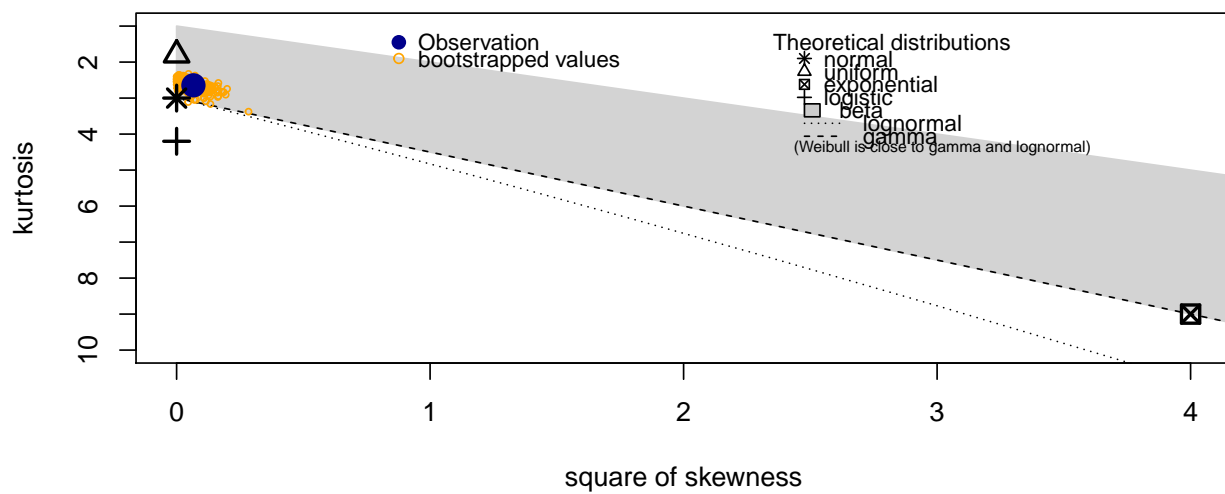


min	max	median	mean	sd	skewness	kurtosis	method
1	3228	32	116.0505	292.3069	6.364168	54.45724	unbiased



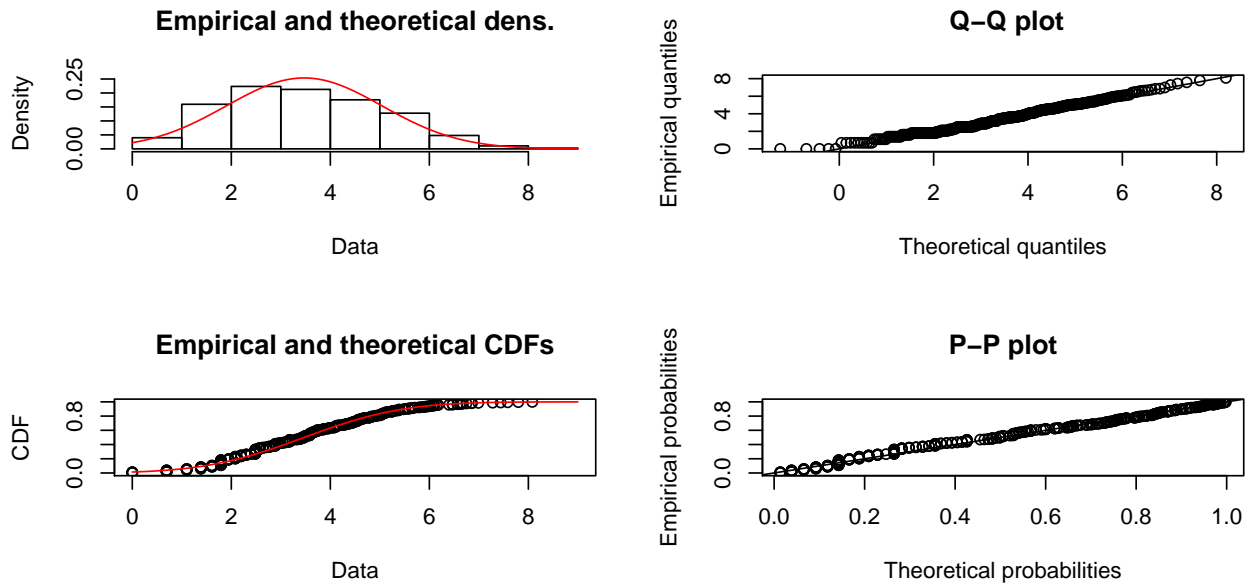
The lambda value of the BoxCox transformation of the Bottles Sold is 0. We therefore apply the log transformation of the dependent variable of Bottles Sold and see a more normal distribution.

Cullen and Frey graph



min	max	median	mean	sd	skewness	kurtosis	method
0	8.079618	3.465736	3.471787	1.57412	0.2576122	2.645664	unbiased

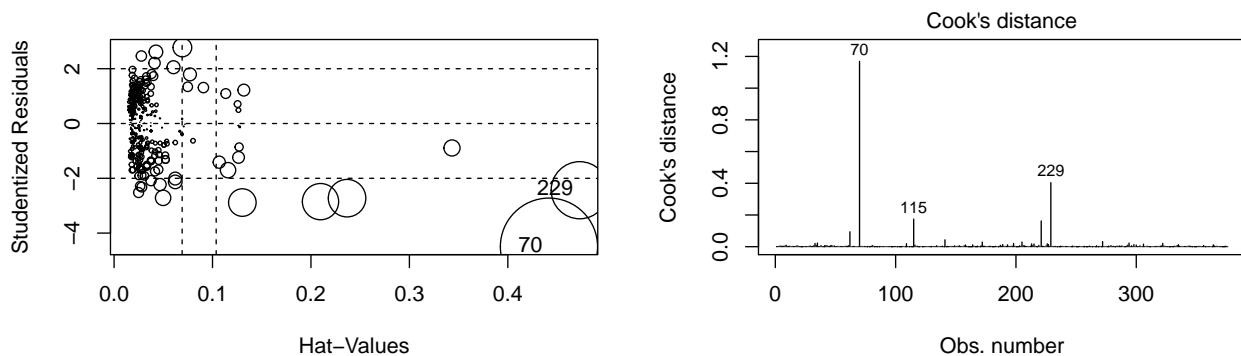
The diagnostic plots show a much more normal distribution after the log transformation.

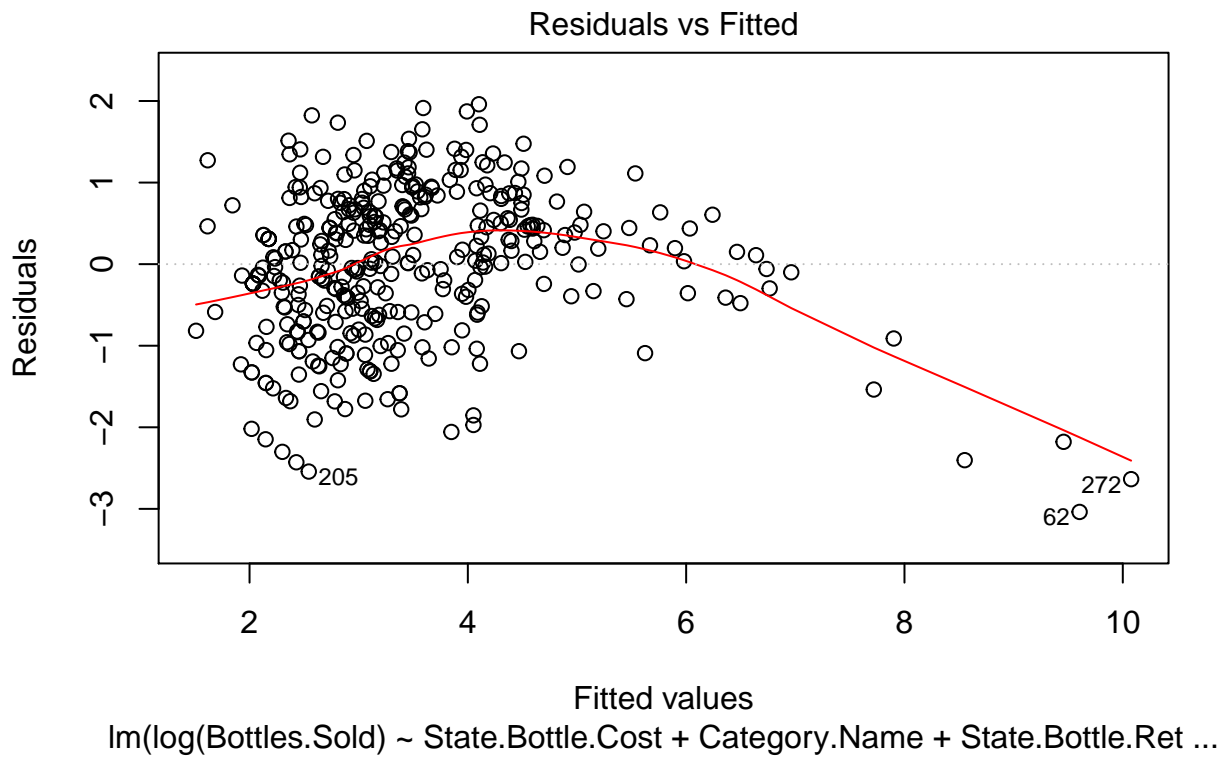


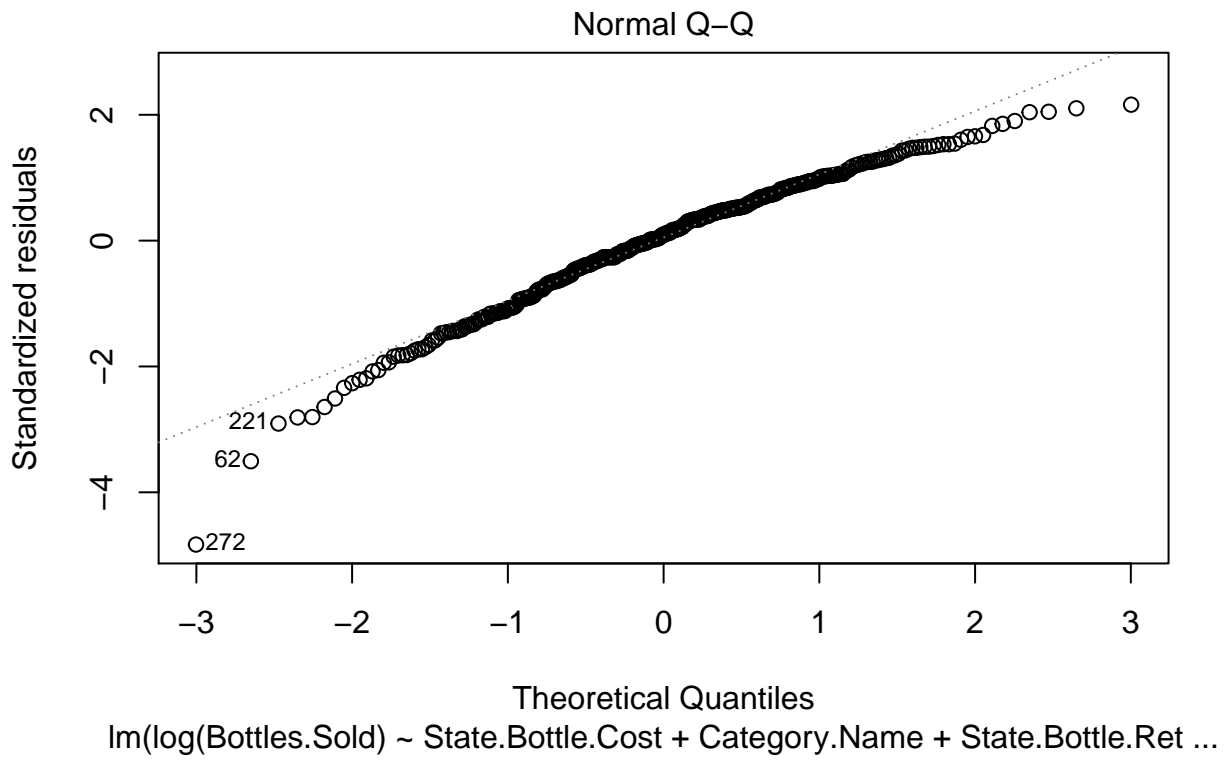
Model Development

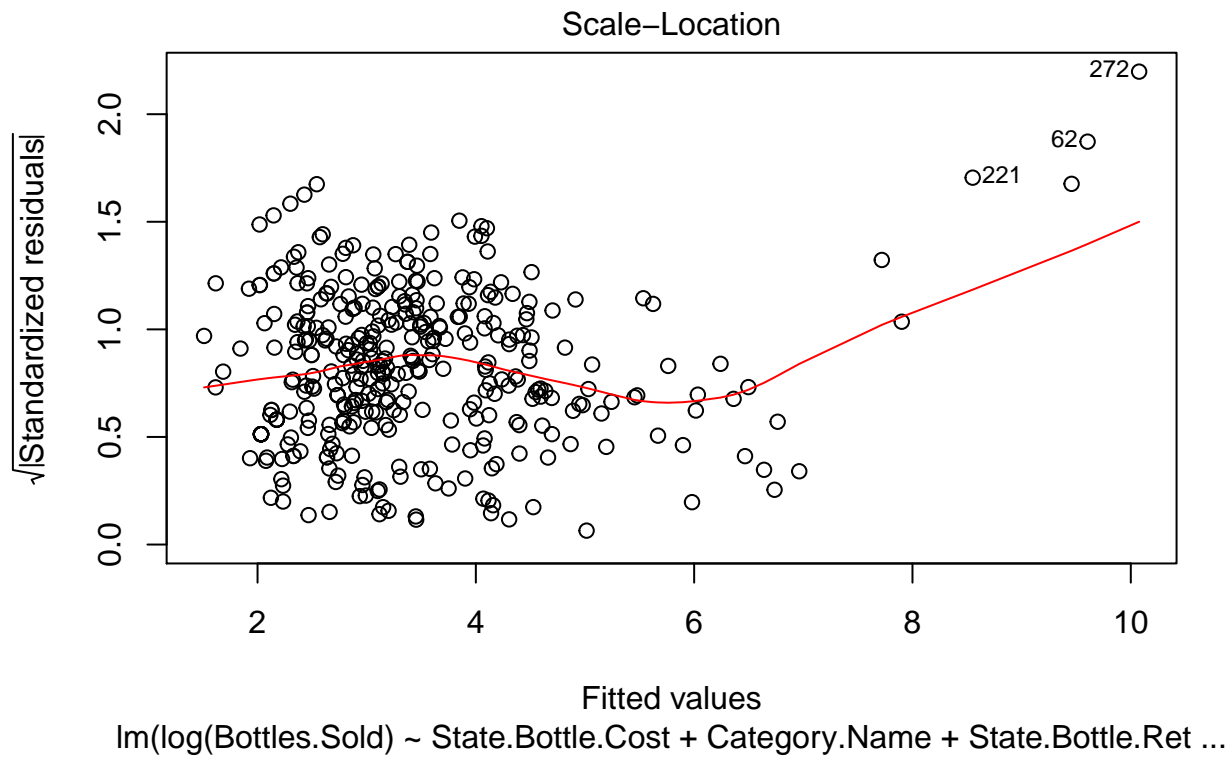
Our first model uses the forward selection method and Bottles Sold as our dependent variable. By using this selection method, the final model excludes the Sales Dollars variable.

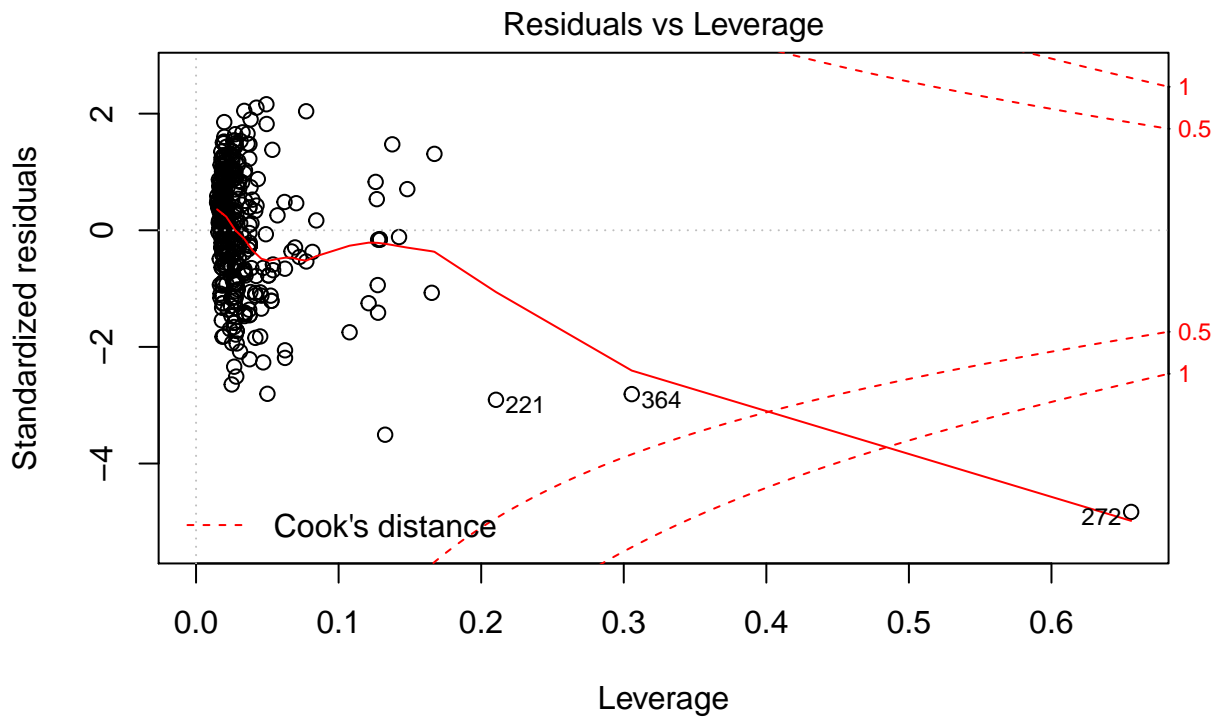
However, several values may have undue influence on the final form of our model. Using the `influencePlot` function from the `car` package and Cook's Distance, we can see which values that have the greatest impact on our model.



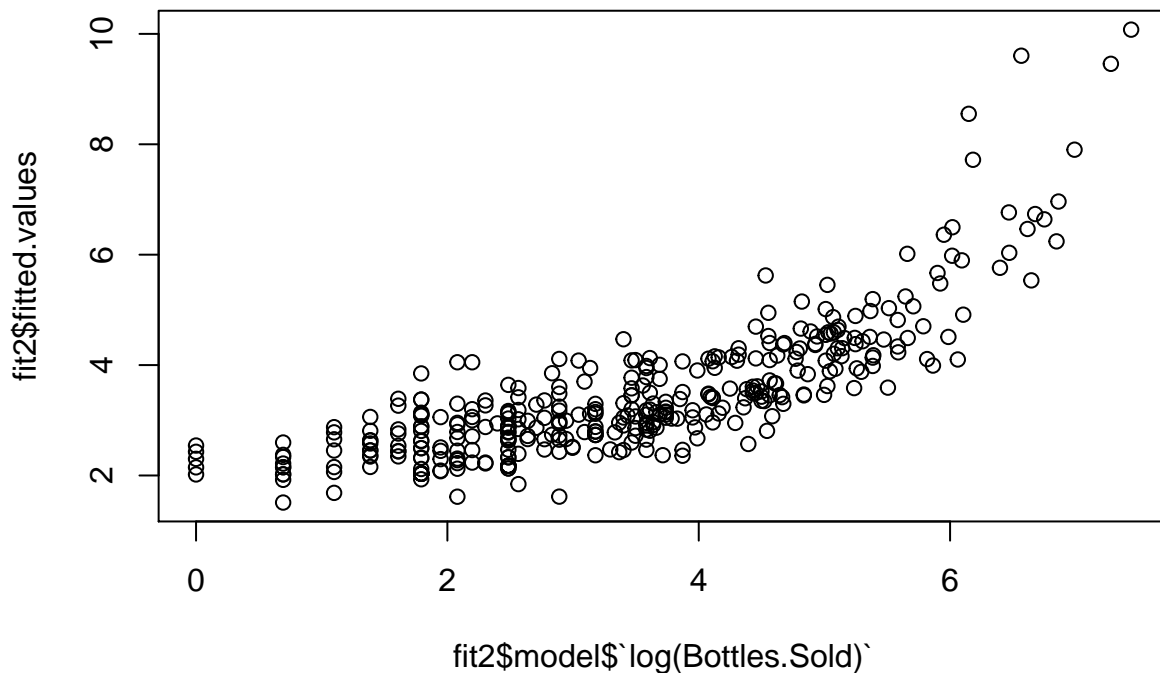








Im(log(Bottles.Sold) ~ State.Bottle.Cost + Category.Name + State.Bottle.Ret ...



Volume.Sold..Gallons=One unit of increase in volume sold will increase the bottles sold by 4.10 units

Category.Name

CANADIAN WHISKIES is the only one that is significant.

CANDIAN WHISKIES vs BLENDED WHISKIES coefficient is 5.492

State.Bottle.Retail=One unit of increase in State.Bottle.Retail will increase the bottles sold by 7.133 units

Pack=One unit of increase in Pack will increase the bottles sold by 7.133 units

Sale..Dollars.=It is slightly significant. One unit of increase in Sale..Dollars. will decrease the bottles sold by 2.189 units

Formula = $-4.196 + 25.314 \text{Volume.Sold..Gallons.} + 5.492 (\text{Category.Name}=\text{"CANADIAN WHISKIES"}) +$

$7.133 * (\text{State.Bottle.Retail}) + 7.521 * (\text{Pack}) - 2.189 * \text{Sale..Dollars.}$

Discussion and Conclusions

In two studies that were conducted, the statistical data from the previous time interval was used to predict what would happen in the next time interval. One study was for pharmaceutical distribution companies. The other study involved data mining models.

In one study involving pharmaceutical distribution companies, the purpose was to propose a novel method to forecast the sales of the companies. Network-based analysis was conducted to find clique sets and group members and to use the sales data of comembers. The reason for this was the lack of sufficient historic sales records of each drug.

Three methods were used to build time series models forecasting sales - ARIMA methodology, neural network, and an advanced hybrid neural network approach. The performance of the proposed method was evaluated using a real dataset provided by one of the leading pharmacy distribution companies in Iran. The results of the evaluation indicated that the proposed method can cope with the low number of past records while accurately forecasting medicine sales.

After exploratory analysis was done on the data, it was concluded that most medicines had different and specific characteristics and sales behavior, it was impractical to make a single prediction model for all medicines, and most sales records had nonlinear relationships.

The reason why the hybrid neural network method was carried out was due to the fact that it is not acceptable to apply a fully linear or nonlinear model on sales data.

The two forecast error measures that were used to evaluate and compare model performance were mean squared error and mean absolute error. The performance of the predicted data was significantly improved when the past records of co members were used.

http://onsearch.cuny.edu/primo_library/libweb/action/display.do?tabs=detailsTab&ct=display&fn=search&doc=TN_gale_ofa

The other study that examined prediction-based inventory optimization using data mining models, the Back propagation neural network was used for training the prediction model. The idea that gave rise to this method of inventory prediction was the idea that the demand of marketing is viewed as the foundation of inventory management.

On the basis of the prediction result, a simple and concise inventory policy was established. Following this, the historic sales data was used to estimate a normal distribution of demand and to calculate the inventory cost with inventory strategy.

Two models (Back Propagation Neural Network and Support Vector Regression) were established using three input variables (historical sales data, the frequency of searching the commodity, and the click volume of the commodity page).

When the back-propagation neural network method was used, there was more accuracy shown in the performance because the predicted values almost matched the actual values in the graphs.

The mean absolute percentage error calculated for this model is 0.06.

http://onsearch.cuny.edu/primo_library/libweb/action/display.do?frbrVersion=2&tabs=detailsTab&ct=display&fn=search&doc

In another study conducted in 2012 in Idaho, the monthly revenue generated was examined rather than the yearly revenue generated. The continued growth was rather owed to the number of weekends a month has (five instead of four) and to the higher prices in neighboring states. In Washington, the voters approved an initiative that led the state to sell its liquor stores and add new distributor and retail fees, making prices in the neighboring states (Idaho and Oregon) look better. There were no changes made in marketing or pricing in response to the regulatory shift in Washington (???). Further research into the proximity of our counties to states and towns with higher prices and regulation may provide more insight into sales and volume of liquor sold. Additionally, reviewing the data by identifying months that has 5 weekends instead of four could provide further insights.

conclude your findings, limitations, and suggest areas for future work.

Appendices

Supplemental tables and/or figures.

dfLiquorSales 14 Variables 376 Observations

X

n	missing	distinct	Info	Mean	Gmd	.05	.10	.25	.50	.75	.90	.95
376	0	376	1	188.5	125.7	19.75	38.50	94.75	188.50	282.25	338.50	357.25

lowest : 1 2 3 4 5, highest: 372 373 374 375 376

Month

n	missing	distinct	Info	Mean	Gmd
376	0	1	0	11	0

Value 11
Frequency 376
Proportion 1

Year

n	missing	distinct	Info	Mean	Gmd
376	0	1	0	2015	0

Value 2015
Frequency 376
Proportion 1

City

n	missing	distinct	value
376	0	1	DES MOINES

Value DES MOINES
Frequency 376
Proportion 1

Category.Name

n	missing	distinct
376	0	8

BLENDED WHISKIES (61, 0.162), CANADIAN WHISKIES (71, 0.189), IRISH WHISKIES (39, 0.104), SCOTCH WHISKIES (41, 0.109), SINGLE BARREL BOURBON WHISKIES (8, 0.021), STRAIGHT BOURBON WHISKIES (66, 0.176), STRAIGHT RYE WHISKIES (27, 0.072), TENNESSEE WHISKIES (63, 0.168)

Store.Number

n	missing	distinct
376	0	73

lowest : 2190 2248 2527 2528 2532, highest: 5131 5132 5137 5145 5169

Store.Name

n	missing	distinct
376	0	72

lowest : AV Superstop Best Food Mart / Des Moines C Fresh Market Cash Saver / E Euclid Ave C
highest: Walgreens #05852 / Des Moines Walgreens #07452 / Des Moines Walgreens #07453 / Des Moines Walgreens #07833 / Des Moines W

Pack

n	missing	distinct	Info	Mean	Gmd	.05	.10	.25	.50	.75	.90	.95
376	0	141	0.99	13.38	6.201	6.00	6.00	9.75	12.00	16.00	21.37	24.00

lowest : 3.00000 4.50000 5.00000 5.25000 5.40000, highest: 30.00000 32.00000 32.57143 33.81818 36.00000

Bottles.Sold

n	missing	distinct	Info	Mean	Gmd	.05	.10	.25	.50	.75	.90	.95
376	0	153	0.999	116.1	173.3	3.00	5.00	9.75	32.00	96.50	247.00	444.25

lowest : 1 2 3 4 5, highest: 1450 1704 1961 2431 3228

Sale..Dollars.

n	missing	distinct	Info	Mean	Gmd	.05	.10	.25	.50	.75	.90	.95
376	0	343	1	1891	3024	44.98	69.84	162.63	371.25	1078.08	2937.44	6963.48

lowest : 7.20 21.74 22.49 26.25 27.14, highest: 35818.26 38945.46 40135.98 54923.22 71157.84

Bottle.Volume..ml.

n	missing	distinct	Info	Mean	Gmd	.05	.10	.25	.50	.75	.90	.95
376	0	180	0.986	803.8	345.4	375.0	435.4	573.8	750.0	953.8	1250.0	1564.3

lowest : 200.000 270.000 287.500 300.000 310.000, highest: 1416.667 1500.000 1550.000 1607.143 1750.000

State.Bottle.Cost

n	missing	distinct	Info	Mean	Gmd	.05	.10	.25	.50
376	0	319	1	101.7	127.4	7.262	10.670	18.495	47.375
.75	.90	.95							
103.537	225.940	338.075							

lowest : 3.21 3.46 3.50 4.10 4.40, highest: 887.95 918.77 1045.11 1362.12 1640.44

State.Bottle.Retail

n	missing	distinct	Info	Mean	Gmd	.05	.10	.25	.50	.75	.90	.95
376	0	322	1	152.6	191.3	10.90	16.01	27.75	71.09	155.34	338.93	507.33

lowest : 4.82 5.19 5.25 6.15 6.60, highest: 1332.12 1378.66 1568.12 2049.35 2467.91

Volume.Sold..Gallons.

n	missing	distinct	Info	Mean	Gmd	.05	.10	.25	.50	.75	.90	.95
376	0	254	1	24.4	38.04	0.4225	0.7350	1.7900	5.1500	18.2400	47.0350	88.4975

lowest : 0.13 0.20 0.30 0.32 0.39, highest: 383.96 417.74 483.05 623.52 830.55

```
describe(dfLiquorSales$Store.Number)
```

```
## dfLiquorSales$Store.Number
##      n missing distinct
##    376      0        73
##
## lowest : 2190 2248 2527 2528 2532, highest: 5131 5132 5137 5145 5169
```

```
describe(dfLiquorSales$Category.Name)
```

```
## dfLiquorSales$Category.Name
##      n missing distinct
##    376      0         8
##
## BLENDED WHISKIES (61, 0.162), CANADIAN WHISKIES (71, 0.189), IRISH
## WHISKIES (39, 0.104), SCOTCH WHISKIES (41, 0.109), SINGLE BARREL BOURBON
## WHISKIES (8, 0.021), STRAIGHT BOURBON WHISKIES (66, 0.176), STRAIGHT RYE
## WHISKIES (27, 0.072), TENNESSEE WHISKIES (63, 0.168)
```

Session Info

- R version 3.3.2 (2016-10-31), x86_64-mingw32
- Locale: LC_COLLATE=English_United States.1252, LC_CTYPE=English_United States.1252, LC_MONETARY=English_United States.1252, LC_NUMERIC=C, LC_TIME=English_United States.1252
- Base packages: base, datasets, graphics, grDevices, methods, stats, utils
- Other packages: car 2.1-4, data.table 1.10.0, dplyr 0.5.0, fitdistrplus 1.0-7, forecast 7.3, Formula 1.2-1, ggplot2 2.2.0, Hmisc 4.0-1, lattice 0.20-34, logspline 2.1.9, MASS 7.3-45, pacman 0.4.1, pander 0.6.0, purrr 0.2.2, RColorBrewer 1.1-2, readr 1.0.0, stargazer 5.2, survival 2.40-1, tibble 1.2, tidyr 0.6.0, tidyverse 1.0.0, timeDate 3012.100, zoo 1.7-13
- Loaded via a namespace (and not attached): acepack 1.4.1, assertthat 0.1, backports 1.0.4, base64 2.0, cluster 2.0.5, colorspace 1.3-1, DBI 0.5-1, digest 0.6.10, evaluate 0.10, foreign 0.8-67, fracdiff 1.4-2, grid 3.3.2, gridExtra 2.2.1, gtable 0.2.0, htmlTable 1.7, htmltools 0.3.5, knitr 1.15.1, latticeExtra 0.6-28, lazyeval 0.2.0, lme4 1.1-12, magrittr 1.5, Matrix 1.2-7.1, MatrixModels 0.4-1, mgcv 1.8-16, minqa 1.2.4, munsell 0.4.3, nlme 3.1-128, nloptr 1.0.4, nnet 7.3-12, openssl 0.9.5, parallel 3.3.2, pbkrtest 0.4-6, plyr 1.8.4, quadprog 1.5-5, quantreg 5.29, R6 2.2.0, Rcpp 0.12.8, rmarkdown 1.2, rpart 4.1-10, rprojroot 1.1, rticles 0.2, scales 0.4.1, SparseM 1.74, splines 3.3.2, stringi 1.1.2, stringr 1.1.0, tools 3.3.2, tseries 0.10-35, yaml 2.1.14

R statistical programming code.

Please see [Final Project.rmd](#) on GitHub for source code.

<https://github.com/ChristopheHunt/DATA-621-Group-1/blob/master/Final%20Project/Final%20Project.Rmd>

References

AC, Cullen, and Frey HC. 1999. "Probabilistic Techniques in Exposure Assessment," 181–241.

Anonymous. 2016. "Specialty Products Grow in Wine, Spirits." *Beverage Industry* 107(7).

Boshart, Rod. 2001. "Liquor Sales in Iowa Set Record." *Gazette*.

David S. Walonick, Ph.D. 1993. "An Overview of Forecasting Methodology." *Survival Statistics*. doi:<http://www.statpac.org/research/library/forecasting.htm>.

Del Buono, Amanda. 2016. "Keeping Spirits High." *Beverage Industry* 107.4: 14–16, 18.

Kumar, Abhishek. 2012. "Demand Forecast Process and Inventory Management."

List of Tables

3	Forward Selection Linear Model 1 with Influencial Points	17
4	Forward Selection Linear Model 1	18

Table 3: Forward Selection Linear Model 1 with Influential Points

	<i>Dependent variable:</i>
	log(Bottles.Sold)
Constant	1.744*** (0.438)
State.Bottle.Cost	0.864*** (0.212)
Category.NameCANADIAN WHISKIES	0.677*** (0.187)
Category.NameIRISH WHISKIES	−0.774*** (0.221)
Category.NameSCOTCH WHISKIES	−0.740*** (0.216)
Category.NameSINGLE BARREL BOURBON WHISKIES	−0.881** (0.423)
Category.NameSTRAIGHT BOURBON WHISKIES	−0.218 (0.200)
Category.NameSTRAIGHT RYE WHISKIES	−0.634** (0.280)
Category.NameTENNESSEE WHISKIES	−0.328* (0.194)
State.Bottle.Retail	−0.572*** (0.141)
Pack	0.061*** (0.014)
Volume.Sold..Gallons.	0.003*** (0.001)
Bottle.Volume..ml.	0.001*** (0.0002)
Observations	376
R ²	0.601
Adjusted R ²	0.588
Residual Std. Error	1.011 (df = 363)
F Statistic	45.537*** (df = 12; 363) (p = 0.000)
Note:	*p<0.1; **p<0.05; ***p<0.01

Table 4: Forward Selection Linear Model 1

	Dependent variable:
	Bottles.Sold
Constant	-31.762*** (7.570)
Volume.Sold..Gallons.	4.093*** (0.162)
Category.NameCANADIAN WHISKIES	36.254*** (6.601)
Category.NameIRISH WHISKIES	-4.874 (8.077)
Category.NameSCOTCH WHISKIES	-8.116 (8.189)
Category.NameSINGLE BARREL BOURBON WHISKIES	12.384 (14.612)
Category.NameSTRAIGHT BOURBON WHISKIES	-3.416 (6.985)
Category.NameSTRAIGHT RYE WHISKIES	17.005* (9.423)
Category.NameTENNESSEE WHISKIES	-1.311 (7.125)
State.Bottle.Retail	0.079*** (0.011)
Pack	2.760*** (0.367)
Sale..Dollars.	-0.004** (0.002)
Observations	376
R ²	0.984
Adjusted R ²	0.984
Residual Std. Error	37.162 (df = 364)
F Statistic	2,076.152*** (df = 11; 364) (p = 0.000)
Note:	*p<0.1; **p<0.05; ***p<0.01