# Modeling Restaurant-Goer's Behavior in The Great Recession: An Econometric Case Study

*Lester Pi*

*4/12/2017*

## Abstract

This is a case study into Yelp restaurant-goers' consumer behavior during The Great Recession based off data from the Yelp Dataset Challenge 9. By manipulating and transforming the dataset, I rebuilt the relevant data into an econometric framework. Combined with a very light semantic analysis, we are able to see how consumer behavior changed during The Great Recession. With this information and econometric models, we can effectively determine how restaurant-goers' behavior will change in the event of a future recession.

## Introduction

Yelp is a platform where users can review businesses based off a star system, with one star being the lowest and five stars being the highest. Along with the review, users write their thoughts about it which typically include why they feel the business earned the score they gave it. Another key feature of Yelp is business information that includes attributes such as price, business type, and location.

The Great Recession hit in December 2007 and lasted until June 2009 and was related to the financial crisis of 2007-08 and subprime mortgage crisis of 2007-09. One of the key aspects from a consumer standpoint is the Consumer Price Index (CPI), which is an indexed measure of prices and purchasing power. According to the Federal Reserve Economic Data of St. Louis (FRED), the CPI for urban food and beverages increases during the most of the recession, but declines towards the end. In other words, purchasing power was weaker during most of the recession.

## Motivation

I wanted to focus on the restaurant industry for three main reasons: Yelp is very well-known for their restaurant reviews, I majoritively use Yelp for restaurant reviews, and I have a personal interest in the food and restaurant world. The Great Recession sets the stage for a great case study in that it was recent enough to occur after Yelp's conception and this also allows to account for systematic time differences. It was also a very interesting recession since it had huge ramifications both domestically and worldwide.

If restaurant Yelpers' behavior during and around The Great Recession period can be modeled, then we can apply this model to a future recession. Depending on the results, this can have important insights for restaurants who are looking to survive, or perhaps even take advantage of, a recession.

## Datasets

The following are the datasets which I used in this case study.

Yelp Dataset Challenge 9: Contains a selective subset of Yelp data covering reviews, users, businesses, tips, and check-ins. Core datasets that will be examined.

FRED GDP: U.S. Real GDP data pulled from FRED. Used to examine The Great Recession.

Yahoo Finance Yelp Stock: Yelp's adjusted closure stock price. Used to measure Yelp company success and performance.

BEA Restaurant Expenditures: Seasonally adjusted real restaurant expenditures from The BEA. Used to connect Yelp restaurant data with a generalized restaurant industry.

## Data Work Oveview

The Yelp data came in large Json files that needed to be converted into R-workable dataframes. Once in a workable format, I explored the data and extracted the relevant information. Regular expressions were used to analyze large groups of text when only a small specific portion was needed. By using SQL queries, I subsetted the interesting data into the needed date ranges and the associated categories, such as by geography and business type.

The other data sets were pulled from their sources, either manually or through R, and formatted and transformed accordingly.

Most of the data is either transformed into growth rates, detrended and/or seasonally adjusted, or kept in original levels.
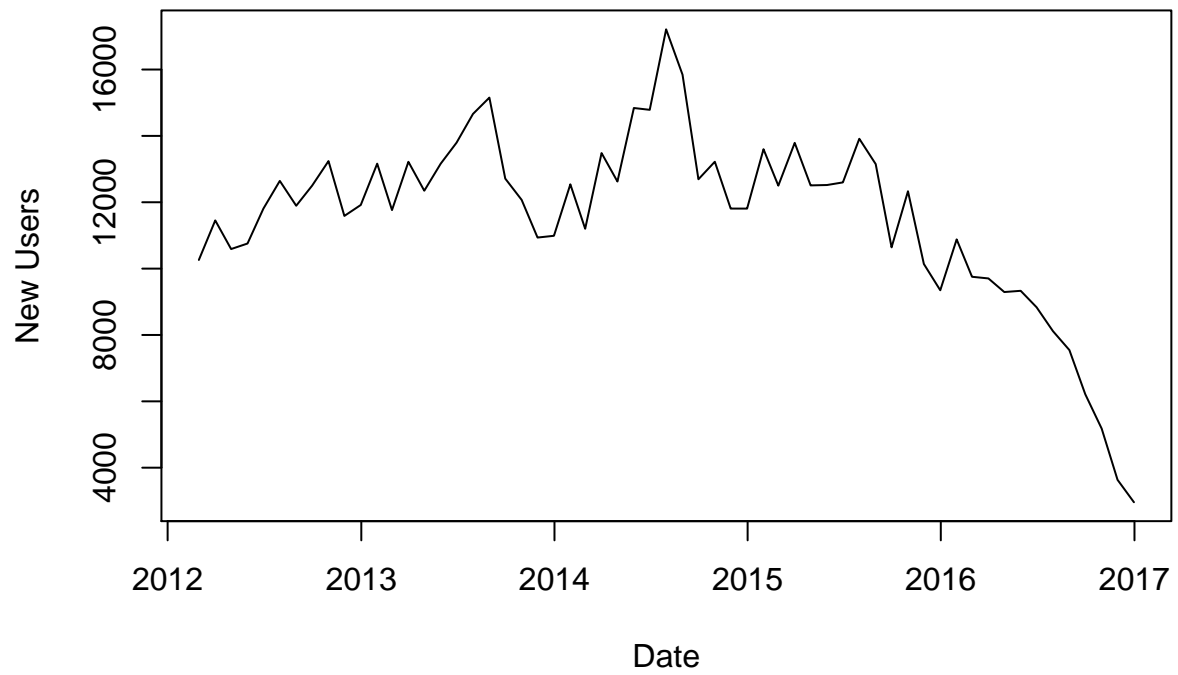
## Challenges and Data Issues Addressed

1. Yelp Data has a disproportionate amount of observations during its early days as well as the latest month due to not having a complete month's worth of data. This can cause statistical insignificance and heteroskedasticity. In order to prevent this, I omitted some of the earliest and latest data.

2. Almost everything is a affected by endogeneity. I tried to prove or disprove what I thought could be a potential instrumental variable that could be used to reduce endogeneity.

3. The Yelp data is a Yelp-decided subset of their data. This can cause large selection bias. I examined the data and saw that it includes small, medium, and large cities alike. There is no way to obtain the unreleased portion of data.

4. Level sets vs. growth rates. Growth rates allow our data to be transformed into stationary (or nearly stationary) data. However, they do not always make intuitive sense for this case study. Therefore, I used the levels for creating linear models, but used growth rates for determining Granger causality and VAR models.

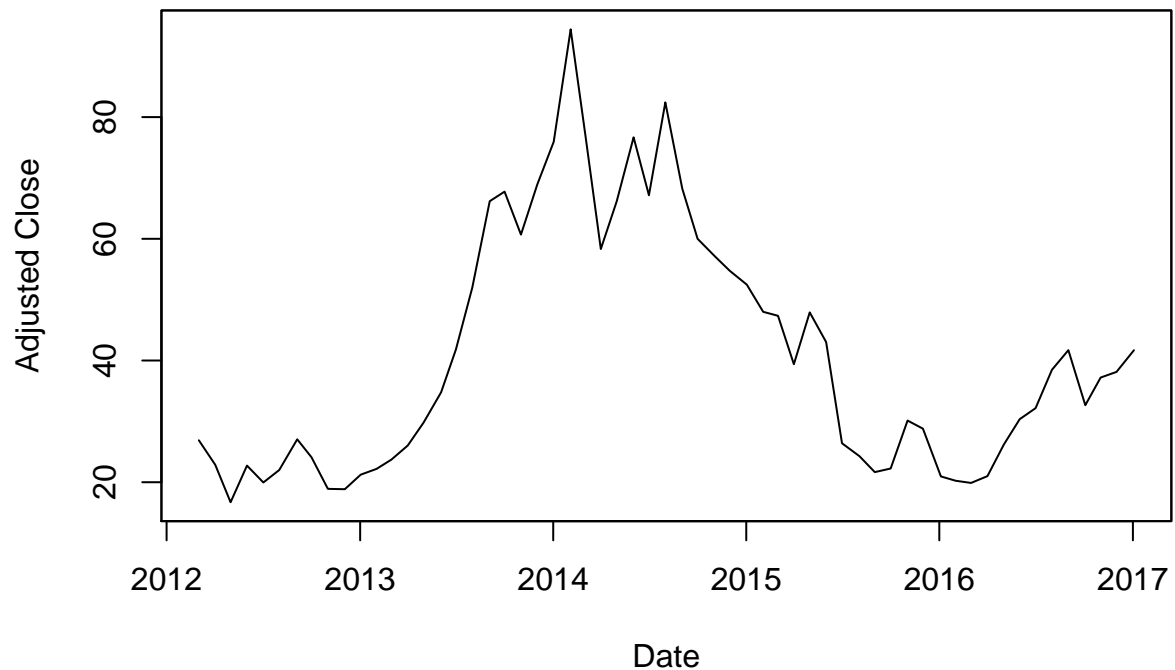## Does Yelp Performance Affect the Userbase?

First, I will explore the Yelp stock data to see if it is of any relevance in an attempt to handle any cases of endogeneity. In this case, the stock data will be a representation of the company performance. Intuitively, there is a chance it may be an instrumental variable such that the better the company is doing, the more they can advertise to, accumulate, and support a larger userbase.

## New User Acounts by Month



```
## time series starts 2012-03-02
## time series ends   2017-01-03
```

# Yelp Stock



```
## 
##  KPSS Test for Level Stationarity
## 
## data:  t
## KPSS Level = 1.1995, Truncation lag parameter = 1, p-value = 0.01
## 
## 
##  Augmented Dickey-Fuller Test
## 
## data:  t
## Dickey-Fuller = -1.0302, Lag order = 3, p-value = 0.9257
## alternative hypothesis: stationary

## 
##  KPSS Test for Level Stationarity
## 
## data:  t
## KPSS Level = 0.53728, Truncation lag parameter = 1, p-value =
## 0.03327
## 
## 
##  Augmented Dickey-Fuller Test
## 
## data:  t
## Dickey-Fuller = -1.4945, Lag order = 3, p-value = 0.7789
## alternative hypothesis: stationary
```

```
## 
##  KPSS Test for Level Stationarity
## 
## data:  t
## KPSS Level = 0.8162, Truncation lag parameter = 1, p-value = 0.01
## 
## 
##  Augmented Dickey-Fuller Test
## 
## data:  t
## Dickey-Fuller = -2.5657, Lag order = 3, p-value = 0.3463
## alternative hypothesis: stationary

## 
##  KPSS Test for Level Stationarity
## 
## data:  t
## KPSS Level = 0.13073, Truncation lag parameter = 1, p-value = 0.1
## 
## 
##  Augmented Dickey-Fuller Test
## 
## data:  t
## Dickey-Fuller = -3.9873, Lag order = 3, p-value = 0.01641
## alternative hypothesis: stationary

## 
## VAR Estimation Results:
## ========================= 
## Endogenous variables: ts_yelp, ts_users
## Deterministic variables: const
## Sample size: 46
## Log Likelihood: 104.197
## Roots of the characteristic polynomial:
## 1.121 0.9968 0.9968 0.9859 0.9859 0.9739 0.9739 0.9716 0.9716 0.9556 0.9556 0.9471 0.9423 0.9423 0.9:
## Call:
## VAR(y = combined, p = select$select[1])
## 
## 
## Estimation results for equation ts_yelp:
## ======================================= 
## ts_yelp = ts_yelp.l1 + ts_users.l1 + ts_yelp.l2 + ts_users.l2 + ts_yelp.l3 + ts_users.l3 + ts_yelp.l
## 
##              Estimate Std. Error t value Pr(>|t|)
## ts_yelp.l1    0.22443    0.24350   0.922    0.367
## ts_users.l1   0.27950    0.43228   0.647    0.525
## ts_yelp.l2   -0.02723    0.26692  -0.102    0.920
## ts_users.l2  -0.12729    0.44781  -0.284    0.779
## ts_yelp.l3    0.06446    0.25133   0.256    0.800
## ts_users.l3  -0.11399    0.47105  -0.242    0.811
## ts_yelp.l4   -0.15175    0.26027  -0.583    0.566
## ts_users.l4   0.41380    0.39727   1.042    0.309
## ts_yelp.l5    0.34786    0.24875   1.398    0.177
## ts_users.l5  -0.04790    0.41595  -0.115    0.909
## ts_yelp.l6    0.23889    0.29288   0.816    0.424
```

```
## ts_users.l6   -0.01589    0.38874   -0.041     0.968
## ts_yelp.l7    -0.07162    0.29937   -0.239     0.813
## ts_users.l7   -0.33711    0.38189   -0.883     0.387
## ts_yelp.l8     0.01820    0.29149    0.062     0.951
## ts_users.l8   -0.22491    0.40855   -0.551     0.588
## ts_yelp.l9     0.07345    0.27828    0.264     0.794
## ts_users.l9    0.40564    0.37089    1.094     0.286
## ts_yelp.l10   -0.15597    0.24322   -0.641     0.528
## ts_users.l10   0.31981    0.38672    0.827     0.418
## ts_yelp.l11   -0.16642    0.22890   -0.727     0.475
## ts_users.l11  -0.26574    0.41547   -0.640     0.529
## ts_yelp.l12   -0.13732    0.23143   -0.593     0.559
## ts_users.l12   0.05773    0.44119    0.131     0.897
## const          0.01366    0.03194    0.428     0.673
##
##
## Residual standard error: 0.1904 on 21 degrees of freedom
## Multiple R-Squared: 0.411,   Adjusted R-squared: -0.2621
## F-statistic: 0.6106 on 24 and 21 DF,  p-value: 0.8779
##
##
## Estimation results for equation ts_users:
## =========================================
## ts_users = ts_yelp.l1 + ts_users.l1 + ts_yelp.l2 + ts_users.l2 + ts_yelp.l3 + ts_users.l3 + ts_yelp.l
##
##              Estimate Std. Error t value Pr(>|t|)
## ts_yelp.l1   -0.158567   0.103539  -1.531  0.14058
## ts_users.l1   0.062517   0.183814   0.340  0.73715
## ts_yelp.l2   -0.034134   0.113501  -0.301  0.76657
## ts_users.l2   0.353135   0.190419   1.855  0.07776 .
## ts_yelp.l3    0.004447   0.106872   0.042  0.96720
## ts_users.l3   0.185054   0.200299   0.924  0.36604
## ts_yelp.l4   -0.072499   0.110674  -0.655  0.51954
## ts_users.l4  -0.073652   0.168926  -0.436  0.66728
## ts_yelp.l5   -0.163585   0.105773  -1.547  0.13691
## ts_users.l5   0.041078   0.176871   0.232  0.81859
## ts_yelp.l6    0.235357   0.124538   1.890  0.07266 .
## ts_users.l6   0.140393   0.165300   0.849  0.40528
## ts_yelp.l7   -0.105272   0.127298  -0.827  0.41755
## ts_users.l7   0.428628   0.162388   2.640  0.01533 *
## ts_yelp.l8    0.053978   0.123949   0.435  0.66765
## ts_users.l8   0.031861   0.173722   0.183  0.85624
## ts_yelp.l9   -0.015866   0.118330  -0.134  0.89462
## ts_users.l9  -0.097292   0.157710  -0.617  0.54393
## ts_yelp.l10   0.053363   0.103420   0.516  0.61126
## ts_users.l10  0.264388   0.164440   1.608  0.12281
## ts_yelp.l11   0.092298   0.097331   0.948  0.35377
## ts_users.l11  0.439468   0.176665   2.488  0.02134 *
## ts_yelp.l12  -0.029449   0.098410  -0.299  0.76769
## ts_users.l12  0.662858   0.187602   3.533  0.00197 **
## const        -0.009345   0.013580  -0.688  0.49886
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```
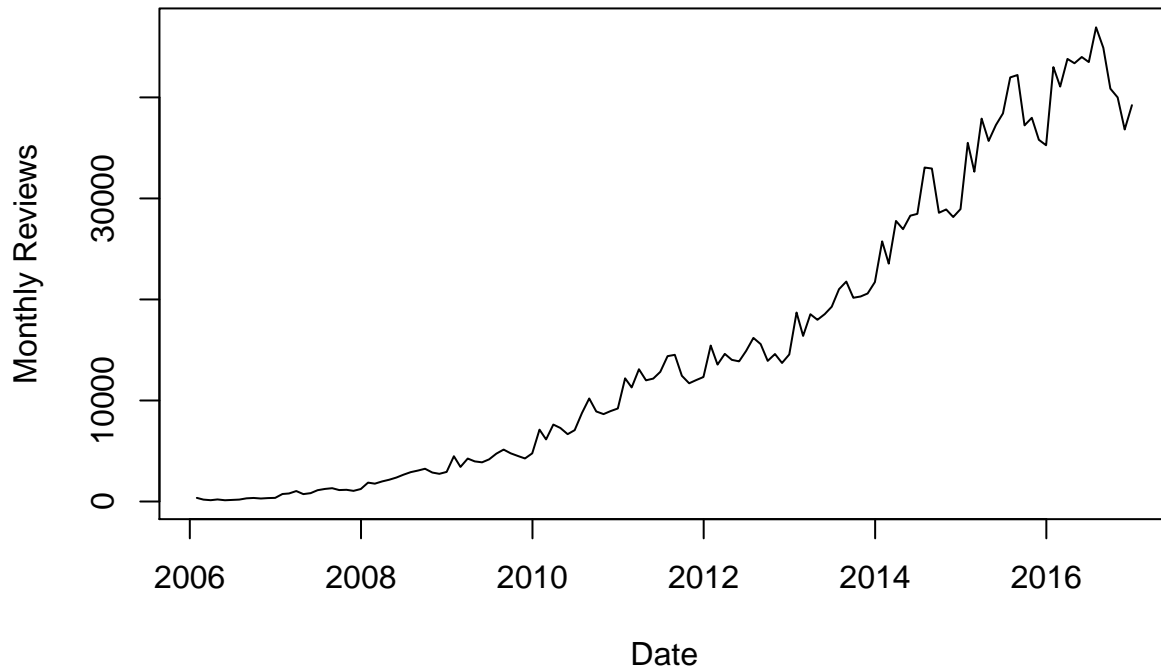
6

```
## 
## Residual standard error: 0.08095 on 21 degrees of freedom
## Multiple R-Squared: 0.798,   Adjusted R-squared: 0.5672
## F-statistic: 3.457 on 24 and 21 DF,  p-value: 0.002768
## 
## 
## 
## Covariance matrix of residuals:
##          ts_yelp ts_users
## ts_yelp  0.036241 0.007759
## ts_users 0.007759 0.006553
## 
## Correlation matrix of residuals:
##         ts_yelp ts_users
## ts_yelp  1.0000   0.5035
## ts_users 0.5035   1.0000

## Granger causality test
## 
## Model 1: ts_users ~ Lags(ts_users, 1:12) + Lags(ts_yelp, 1:12)
## Model 2: ts_users ~ Lags(ts_users, 1:12)
##   Res.Df  Df     F Pr(>F)
## 1     21
## 2     33 -12 1.071 0.4288

## Granger causality test
## 
## Model 1: ts_yelp ~ Lags(ts_yelp, 1:12) + Lags(ts_users, 1:12)
## Model 2: ts_yelp ~ Lags(ts_yelp, 1:12)
##   Res.Df  Df     F Pr(>F)
## 1     21
## 2     33 -12 0.5091 0.8856
```

The number of new users do not have an effect on the stock value of Yelp and vice versa. This allows to eliminate the stock value as an instrumental variable.

# Connecting Users and Reviews

We will attempt to connect new users with new reviews.

## New Reviews by Month



```
##
## VAR Estimation Results:
## =========================
## Endogenous variables: log_rev_count, log_user_count
## Deterministic variables: const
## Sample size: 119
## Log Likelihood: 297.714
## Roots of the characteristic polynomial:
## 1.048 0.9737 0.9737 0.9716 0.9716  0.97  0.97 0.9687 0.9581 0.9581 0.9101 0.9101 0.8859 0.8859 0.857
## Call:
## VAR(y = rates_combined, p = select$select[1])
##
##
## Estimation results for equation log_rev_count:
## ===============================================
## log_rev_count = log_rev_count.l1 + log_user_count.l1 + log_rev_count.l2 + log_user_count.l2 + log_re
##
##                   Estimate Std. Error t value Pr(>|t|)
## log_rev_count.l1  -0.447591   0.102482  -4.367 3.23e-05 ***
## log_user_count.l1  0.073934   0.092299   0.801 0.425138
## log_rev_count.l2  -0.140889   0.109751  -1.284 0.202398
## log_user_count.l2  0.208054   0.097639   2.131 0.035711 *
## log_rev_count.l3   0.024924   0.105031   0.237 0.812936
## log_user_count.l3 -0.053311   0.095380  -0.559 0.577540
## log_rev_count.l4   0.210029   0.098349   2.136 0.035316 *
## log_user_count.l4 -0.405138   0.096096  -4.216 5.72e-05 ***
```

```
## log_rev_count.l5     0.005835   0.095989    0.061 0.951655
## log_user_count.l5   -0.071718   0.102943   -0.697 0.487721
## log_rev_count.l6     0.113966   0.094900    1.201 0.232807
## log_user_count.l6    0.108670   0.101916    1.066 0.289033
## log_rev_count.l7     0.115462   0.094805    1.218 0.226318
## log_user_count.l7    0.124041   0.103038    1.204 0.231678
## log_rev_count.l8    -0.084966   0.095834   -0.887 0.377561
## log_user_count.l8    0.003598   0.103881    0.035 0.972444
## log_rev_count.l9    -0.243305   0.090552   -2.687 0.008528 **
## log_user_count.l9    0.122566   0.101162    1.212 0.228713
## log_rev_count.l10   -0.184751   0.095714   -1.930 0.056592 .
## log_user_count.l10   0.299530   0.100755    2.973 0.003748 **
## log_rev_count.l11   -0.132124   0.086128   -1.534 0.128376
## log_user_count.l11   0.214219   0.101571    2.109 0.037598 *
## log_rev_count.l12    0.092019   0.066103    1.392 0.167189
## log_user_count.l12   0.373805   0.096160    3.887 0.000189 ***
## const                0.032366   0.013296    2.434 0.016811 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##
## Residual standard error: 0.07862 on 94 degrees of freedom
## Multiple R-Squared: 0.6877,  Adjusted R-squared: 0.6079
## F-statistic: 8.623 on 24 and 94 DF,  p-value: 7.305e-15
##
##
## Estimation results for equation log_user_count:
## ================================================
## log_user_count = log_rev_count.l1 + log_user_count.l1 + log_rev_count.l2 + log_user_count.l2 + log_r
##
##                    Estimate Std. Error t value Pr(>|t|)
## log_rev_count.l1   -0.228340   0.108927   -2.096  0.03874 *
## log_user_count.l1   0.085515   0.098103    0.872  0.38560
## log_rev_count.l2    0.109457   0.116653    0.938  0.35049
## log_user_count.l2   0.032154   0.103779    0.310  0.75738
## log_rev_count.l3    0.155007   0.111635    1.389  0.16826
## log_user_count.l3   0.086143   0.101378    0.850  0.39764
## log_rev_count.l4    0.212001   0.104533    2.028  0.04538 *
## log_user_count.l4  -0.139470   0.102139   -1.365  0.17536
## log_rev_count.l5   -0.003095   0.102025   -0.030  0.97586
## log_user_count.l5   0.064438   0.109416    0.589  0.55732
## log_rev_count.l6   -0.001041   0.100868   -0.010  0.99179
## log_user_count.l6   0.187274   0.108324    1.729  0.08712 .
## log_rev_count.l7   -0.085882   0.100767   -0.852  0.39622
## log_user_count.l7   0.260097   0.109517    2.375  0.01958 *
## log_rev_count.l8   -0.212682   0.101861   -2.088  0.03951 *
## log_user_count.l8   0.012833   0.110413    0.116  0.90772
## log_rev_count.l9   -0.310263   0.096246   -3.224  0.00174 **
## log_user_count.l9   0.237115   0.107524    2.205  0.02988 *
## log_rev_count.l10  -0.282063   0.101732   -2.773  0.00671 **
## log_user_count.l10  0.349787   0.107091    3.266  0.00152 **
## log_rev_count.l11  -0.161645   0.091543   -1.766  0.08068 .
## log_user_count.l11  0.258147   0.107957    2.391  0.01879 *
## log_rev_count.l12   0.256980   0.070259    3.658  0.00042 ***
```

```
## log_user_count.l12   0.257381    0.102207    2.518   0.01348 *
## const              -0.008062    0.014132   -0.570   0.56971
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##
## Residual standard error: 0.08356 on 94 degrees of freedom
## Multiple R-Squared: 0.618,   Adjusted R-squared: 0.5204
## F-statistic: 6.336 on 24 and 94 DF,  p-value: 3.029e-11
##
##
##
## Covariance matrix of residuals:
##              log_rev_count log_user_count
## log_rev_count       0.006181        0.002504
## log_user_count      0.002504        0.006982
##
## Correlation matrix of residuals:
##              log_rev_count log_user_count
## log_rev_count        1.0000          0.3811
## log_user_count       0.3811          1.0000
## Granger causality test
##
## Model 1: log_rev_count ~ Lags(log_rev_count, 1:12) + Lags(log_user_count, 1:12)
## Model 2: log_rev_count ~ Lags(log_rev_count, 1:12)
##   Res.Df  Df      F    Pr(>F)
## 1     94
## 2    106 -12 5.9487 1.318e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## Granger causality test
##
## Model 1: log_user_count ~ Lags(log_user_count, 1:12) + Lags(log_rev_count, 1:12)
## Model 2: log_user_count ~ Lags(log_user_count, 1:12)
##   Res.Df  Df      F    Pr(>F)
## 1     94
## 2    106 -12 3.4167 0.000354 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Call:
## lm(formula = log_rev_count ~ log_user_count)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.68573 -0.04008  0.00958  0.04946  0.41635
##
## Coefficients:
##                Estimate Std. Error t value Pr(>|t|)
## (Intercept)     0.02131    0.01183   1.801   0.0741 .
## log_user_count  0.77640    0.08763   8.860 5.44e-15 ***
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1341 on 129 degrees of freedom
## Multiple R-squared:  0.3783, Adjusted R-squared:  0.3735
## F-statistic: 78.51 on 1 and 129 DF,  p-value: 5.444e-15

## NULL

##
##  KPSS Test for Level Stationarity
##
## data:  t
## KPSS Level = 0.27805, Truncation lag parameter = 2, p-value = 0.1
##
##
##  Augmented Dickey-Fuller Test
##
## data:  t
## Dickey-Fuller = -4.91, Lag order = 5, p-value = 0.01
## alternative hypothesis: stationary
##
##
##  Augmented Dickey-Fuller Test
##
## data:  t
## Dickey-Fuller = -4.91, Lag order = 5, p-value = 0.01
## alternative hypothesis: stationary
```
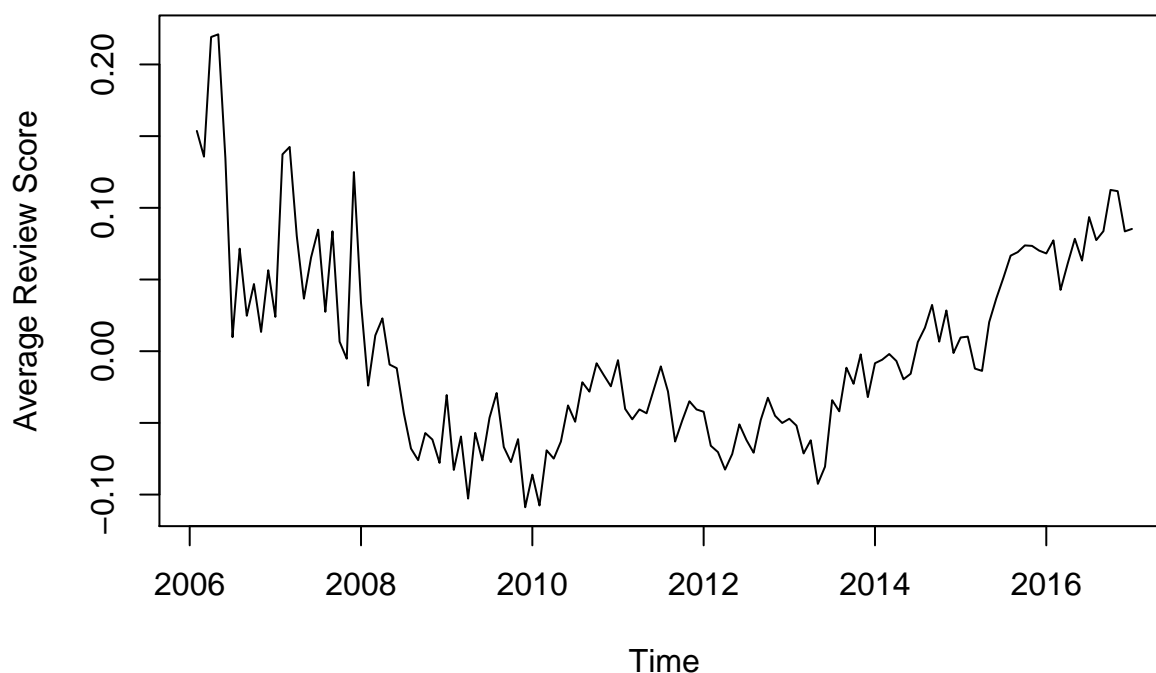
Through the VAR model results, Granger causality test, and cointegration test, we can conclude that users and reviews can be used interchangably. The error descriptive statistics (Note: error descriptive statistics surpressed in writeup) do not look perfect, but we just want to see if reviews can be used in place of users. From here on out, we will focus on reviews since reviews contain more valuable information than the user data.

# Do Review Scores Change?

By examining review scores, we can see if people become more or less critical during the recession, and many other insights.
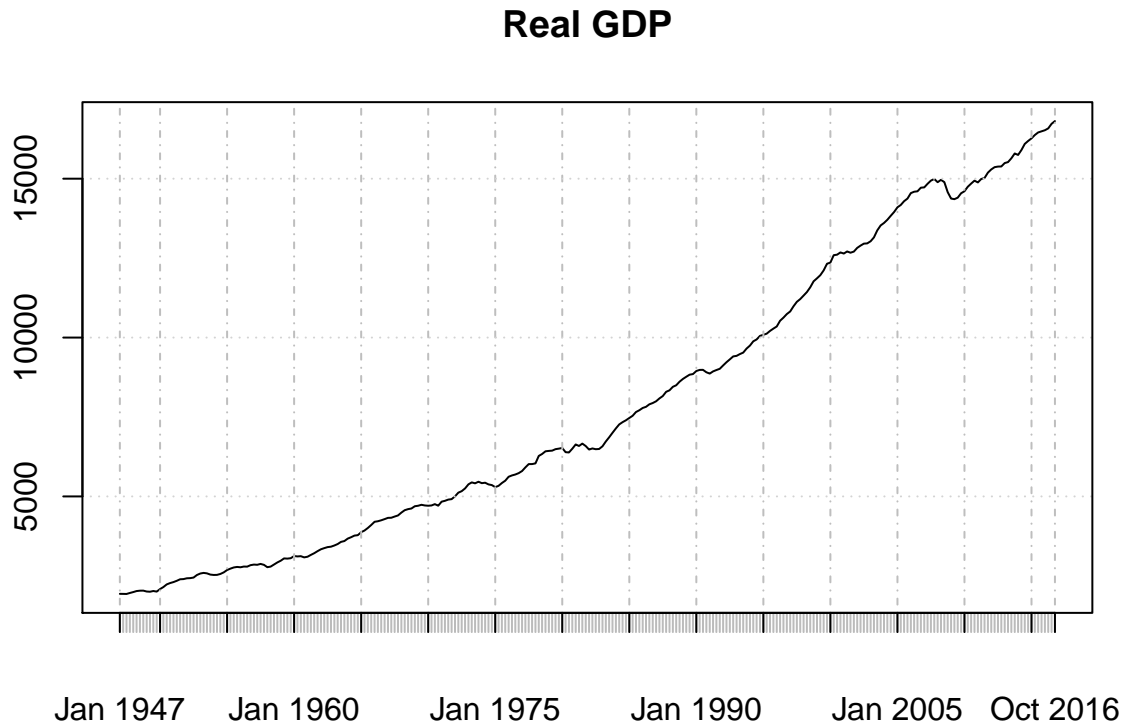
## Average Review Scores by Month



```
## 
## Call:
## lm(formula = stars_avg ~ stars_recession_dummy)
## 
## Residuals:
##      Min       1Q   Median       3Q      Max 
## -0.12291 -0.04541 -0.01606  0.04614  0.22330 
## 
## Coefficients:
##                        Estimate Std. Error t value Pr(>|t|)    
## (Intercept)            3.712971   0.006261 593.052  < 2e-16 ***
## stars_recession_dummy -0.048566   0.016502  -2.943  0.00385 ** 
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 0.06655 on 130 degrees of freedom
## Multiple R-squared:  0.06246,    Adjusted R-squared:  0.05525 
## F-statistic: 8.661 on 1 and 130 DF,  p-value: 0.003851
## 
## NULL
## 
## 
##  KPSS Test for Level Stationarity
## 
## data:  t
## KPSS Level = 0.84648, Truncation lag parameter = 2, p-value = 0.01
## 
```
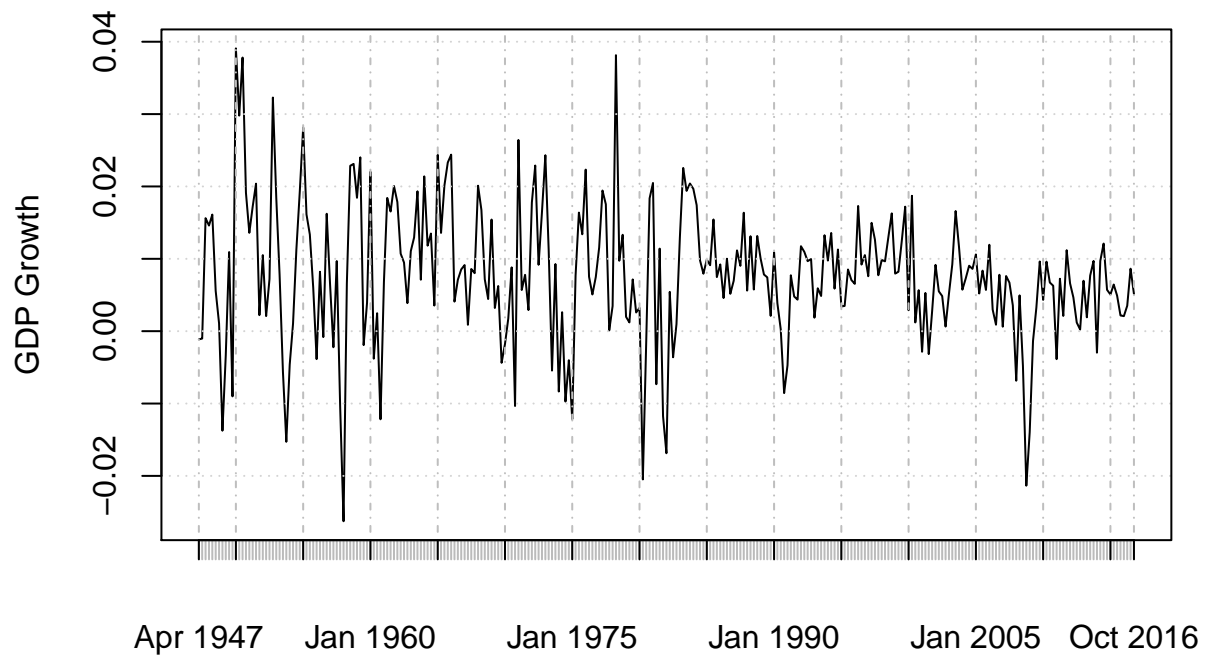
```
##
##  Augmented Dickey-Fuller Test
##
## data:  t
## Dickey-Fuller = -1.9674, Lag order = 5, p-value = 0.5901
## alternative hypothesis: stationary

##
##  KPSS Test for Level Stationarity
##
## data:  t
## KPSS Level = 0.20563, Truncation lag parameter = 2, p-value = 0.1
##
##
##  Augmented Dickey-Fuller Test
##
## data:  t
## Dickey-Fuller = -5.5779, Lag order = 5, p-value = 0.01
## alternative hypothesis: stationary

##
## Call:
## tslm(formula = ts_stars ~ trend + season)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.10883 -0.04940 -0.01198  0.04378  0.22099
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3.677e+00  2.426e-02 151.563   <2e-16 ***
## trend       -2.068e-05  1.619e-04  -0.128    0.899
## season2      4.736e-02  3.015e-02   1.571    0.119
## season3      4.357e-02  3.014e-02   1.446    0.151
## season4      3.505e-02  3.013e-02   1.163    0.247
## season5      3.822e-02  3.012e-02   1.269    0.207
## season6      2.959e-02  3.011e-02   0.983    0.328
## season7      3.114e-02  3.011e-02   1.034    0.303
## season8      4.315e-02  3.010e-02   1.433    0.154
## season9      2.931e-02  3.010e-02   0.974    0.332
## season10     1.768e-02  3.010e-02   0.587    0.558
## season11     1.906e-02  3.009e-02   0.633    0.528
## season12     2.838e-02  3.009e-02   0.943    0.347
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.07057 on 119 degrees of freedom
## Multiple R-squared:  0.03497,    Adjusted R-squared:  -0.06235
## F-statistic: 0.3593 on 12 and 119 DF,  p-value: 0.9748
```

## Detrended & Seasonally Adjusted Review Scores by Month



```
##
## Call:
## lm(formula = detrend_stars ~ stars_recession_dummy)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.11605 -0.04819 -0.01401  0.05040  0.21376
##
## Coefficients:
##                        Estimate Std. Error t value Pr(>|t|)
## (Intercept)            0.007226   0.006128   1.179  0.24049
## stars_recession_dummy -0.050203   0.016153  -3.108  0.00231 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.06514 on 130 degrees of freedom
## Multiple R-squared:  0.06917,    Adjusted R-squared:  0.06201
## F-statistic:  9.66 on 1 and 130 DF,  p-value: 0.002314

## NULL
```

The regression results cannot be trusted because the error statistics are all over the place. No conclusion can be made yet.

```
##     As of 0.4-0, 'getSymbols' uses env=parent.frame() and
##  auto.assign=TRUE by default.
##
##  This  behavior  will be  phased out in 0.5-0  when the call  will
```

```
##  default to use auto.assign=FALSE. getOption("getSymbols.env") and
##  getOptions("getSymbols.auto.assign") are now checked for alternate defaults
##
##  This message is shown once per session and may be disabled by setting
##  options("getSymbols.warning4.0"=FALSE). See ?getSymbols for more details.

## [1] "GDPC96"
```

**Real GDP**

# Real GDP Growth Rate

# Real GDP Growth 2006Q2 to 2016Q4

## Review Growth Rate by Quarter



```
##
## VAR Estimation Results:
## =========================
## Endogenous variables: log_rev_quarter, gdp_growth_subset
## Deterministic variables: const
## Sample size: 31
## Log Likelihood: 239.132
## Roots of the characteristic polynomial:
## 0.9961 0.9961 0.9922 0.9922 0.986 0.9774 0.9774 0.9695 0.9612 0.9612 0.9571 0.9571 0.9324 0.9324 0.9(
## Call:
## VAR(y = gdp_combined, p = 12)
##
##
## Estimation results for equation log_rev_quarter:
## ================================================
## log_rev_quarter = log_rev_quarter.l1 + gdp_growth_subset.l1 + log_rev_quarter.l2 + gdp_growth_subset
##
##                     Estimate Std. Error t value Pr(>|t|)
## log_rev_quarter.l1    0.41495    0.36251   1.145   0.2960
## gdp_growth_subset.l1  -0.69319    3.48625  -0.199   0.8490
## log_rev_quarter.l2    0.28843    0.37771   0.764   0.4740
## gdp_growth_subset.l2   4.10864    3.15933   1.300   0.2412
## log_rev_quarter.l3    0.27193    0.32528   0.836   0.4352
## gdp_growth_subset.l3   1.55990    3.26731   0.477   0.6500
## log_rev_quarter.l4    0.63886    0.20809   3.070   0.0219 *
## gdp_growth_subset.l4   1.46361    2.96026   0.494   0.6386
```

```
## log_rev_quarter.l5       -0.21532    0.23826  -0.904    0.4010
## gdp_growth_subset.l5    -0.43820    3.45264  -0.127    0.9032
## log_rev_quarter.l6       -0.16586    0.22996  -0.721    0.4979
## gdp_growth_subset.l6     4.69612    3.45845   1.358    0.2233
## log_rev_quarter.l7       -0.13165    0.24321  -0.541    0.6078
## gdp_growth_subset.l7    -4.04228    2.95912  -1.366    0.2209
## log_rev_quarter.l8        0.05838    0.24925   0.234    0.8226
## gdp_growth_subset.l8     5.04951    3.49119   1.446    0.1982
## log_rev_quarter.l9        0.11847    0.20277   0.584    0.5803
## gdp_growth_subset.l9     0.05217    3.36605   0.016    0.9881
## log_rev_quarter.l10       0.11208    0.13776   0.814    0.4470
## gdp_growth_subset.l10    6.34267    3.22213   1.968    0.0966 .
## log_rev_quarter.l11      -0.08093    0.15120  -0.535    0.6117
## gdp_growth_subset.l11   -1.75022    2.11397  -0.828    0.4394
## log_rev_quarter.l12       0.43781    0.17667   2.478    0.0479 *
## gdp_growth_subset.l12    2.53313    2.12131   1.194    0.2775
## const                    -0.16460    0.20382  -0.808    0.4502
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##
## Residual standard error: 0.04107 on 6 degrees of freedom
## Multiple R-Squared: 0.9823,  Adjusted R-squared: 0.9117
## F-statistic:  13.9 on 24 and 6 DF,  p-value: 0.001734
##
##
## Estimation results for equation gdp_growth_subset:
## ====================================================
## gdp_growth_subset = log_rev_quarter.l1 + gdp_growth_subset.l1 + log_rev_quarter.l2 + gdp_growth_subs
##
##                     Estimate Std. Error t value Pr(>|t|)
## log_rev_quarter.l1    -0.038721   0.037669  -1.028    0.3436
## gdp_growth_subset.l1  -0.467584   0.362259  -1.291    0.2443
## log_rev_quarter.l2    -0.043579   0.039248  -1.110    0.3094
## gdp_growth_subset.l2  -0.419066   0.328289  -1.277    0.2490
## log_rev_quarter.l3    -0.006052   0.033800  -0.179    0.8638
## gdp_growth_subset.l3  -0.409673   0.339509  -1.207    0.2730
## log_rev_quarter.l4    -0.035484   0.021623  -1.641    0.1519
## gdp_growth_subset.l4  -0.620561   0.307603  -2.017    0.0902 .
## log_rev_quarter.l5    -0.003939   0.024758  -0.159    0.8788
## gdp_growth_subset.l5  -0.263539   0.358766  -0.735    0.4903
## log_rev_quarter.l6     0.025456   0.023895   1.065    0.3277
## gdp_growth_subset.l6  -0.486785   0.359369  -1.355    0.2243
## log_rev_quarter.l7    -0.015218   0.025272  -0.602    0.5691
## gdp_growth_subset.l7  -0.176916   0.307484  -0.575    0.5860
## log_rev_quarter.l8    -0.029032   0.025900  -1.121    0.3052
## gdp_growth_subset.l8  -0.585704   0.362772  -1.615    0.1575
## log_rev_quarter.l9    -0.002672   0.021070  -0.127    0.9032
## gdp_growth_subset.l9  -0.410525   0.349768  -1.174    0.2850
## log_rev_quarter.l10   -0.022118   0.014315  -1.545    0.1733
## gdp_growth_subset.l10 -0.197934   0.334813  -0.591    0.5760
## log_rev_quarter.l11   -0.034505   0.015711  -2.196    0.0705 .
## gdp_growth_subset.l11  0.088467   0.219664   0.403    0.7011
## log_rev_quarter.l12   -0.014660   0.018358  -0.799    0.4550
```

```
## gdp_growth_subset.l12 -0.410368   0.220426  -1.862   0.1120
## const                   0.044382   0.021179   2.096   0.0810 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##
## Residual standard error: 0.004267 on 6 degrees of freedom
## Multiple R-Squared: 0.7701,  Adjusted R-squared: -0.1493
## F-statistic: 0.8376 on 24 and 6 DF,  p-value: 0.657
##
##
##
## Covariance matrix of residuals:
##                 log_rev_quarter gdp_growth_subset
## log_rev_quarter       0.0016865          1.116e-04
## gdp_growth_subset     0.0001116          1.821e-05
##
## Correlation matrix of residuals:
##                 log_rev_quarter gdp_growth_subset
## log_rev_quarter           1.000             0.637
## gdp_growth_subset         0.637             1.000
##
## Call:
## lm(formula = df_rev_quarter$coredata.df_rev_m_quarter. ~ recession_dummy_reviews_q)
##
## Residuals:
##    Min    1Q Median    3Q    Max
## -53417 -28547  -3776  24826  78800
##
## Coefficients:
##                            Estimate Std. Error t value Pr(>|t|)
## (Intercept)                   53903       6553   8.225 2.72e-10 ***
## recession_dummy_reviews_q    -45589      16430  -2.775  0.00821 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 39860 on 42 degrees of freedom
## Multiple R-squared:  0.1549, Adjusted R-squared:  0.1348
## F-statistic: 7.699 on 1 and 42 DF,  p-value: 0.008211

## NULL

##
## Call:
## lm(formula = log_rev_quarter ~ rec_dummy_rev_growth_q)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.44588 -0.13871 -0.02718  0.09270  0.83847
##
## Coefficients:
##                     Estimate Std. Error t value Pr(>|t|)
## (Intercept)          0.10999    0.03739   2.941  0.00535 **
## rec_dummy_rev_growth_q 0.05857    0.09268   0.632  0.53093
```
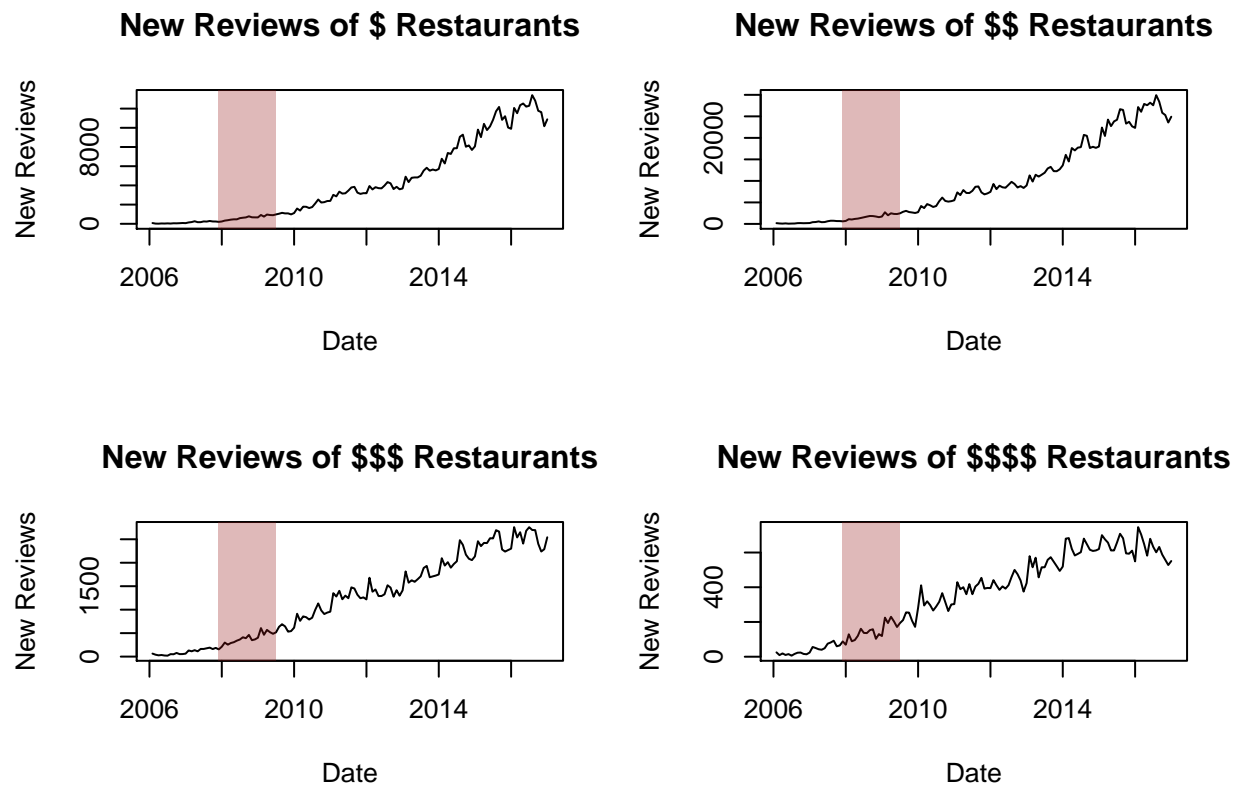
20

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2244 on 41 degrees of freedom
## Multiple R-squared:  0.009647,   Adjusted R-squared:  -0.01451
## F-statistic: 0.3994 on 1 and 41 DF,  p-value: 0.5309

## NULL
```

It looks that using the growth rates are better suited based off the error descriptive statistics, but we want to examine the social dynamics and not reviews as a whole.

Let's continue by splitting the restaurants by dollar signs ($).

# Examining Restaurants by Prices

How do people react to prices during a recession? To examine this, we will subset our data into 4 categories by price, one for each price level. This will be refered to as how many dollar signs ($) a business is.

### New Reviews of $ Restaurants

### New Reviews of $$ Restaurants

### New Reviews of $$$ Restaurants

### New Reviews of $$$$ Restaurants

There is a strong and obvious trend along with seasonality. This is something that should be wiped out as much as possible to see the true effect of the recession.

```
##
## Call:
## lm(formula = tslm_d1_resid ~ recession_dummy_dollars_m)
##
## Residuals:
```

```
##     Min     1Q  Median      3Q     Max
## -1932.8 -1098.0  -162.2   975.3  2613.1
##
## Coefficients:
##                          Estimate Std. Error t value Pr(>|t|)
## (Intercept)                 33.99     120.76   0.281    0.779
## recession_dummy_dollars_m  -236.11     318.30  -0.742    0.460
##
## Residual standard error: 1284 on 130 degrees of freedom
## Multiple R-squared:  0.004215,   Adjusted R-squared:  -0.003445
## F-statistic: 0.5502 on 1 and 130 DF,  p-value: 0.4596

##
## Call:
## lm(formula = tslm_d2_resid ~ recession_dummy_dollars_m)
##
## Residuals:
##     Min     1Q  Median      3Q     Max
## -4262.1 -2293.1  -201.7  2170.2  5499.9
##
## Coefficients:
##                          Estimate Std. Error t value Pr(>|t|)
## (Intercept)                 71.32     259.47   0.275    0.784
## recession_dummy_dollars_m  -495.52     683.92  -0.725    0.470
##
## Residual standard error: 2758 on 130 degrees of freedom
## Multiple R-squared:  0.004022,   Adjusted R-squared:  -0.00364
## F-statistic: 0.5249 on 1 and 130 DF,  p-value: 0.47

##
## Call:
## lm(formula = tslm_d3_resid ~ recession_dummy_dollars_m)
##
## Residuals:
##     Min     1Q  Median      3Q     Max
## -301.93  -84.56  -10.65   76.69  290.26
##
## Coefficients:
##                          Estimate Std. Error t value Pr(>|t|)
## (Intercept)                 14.40      10.99   1.310  0.19253
## recession_dummy_dollars_m  -100.01      28.97  -3.453  0.00075 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 116.8 on 130 degrees of freedom
## Multiple R-squared:  0.084,   Adjusted R-squared:  0.07695
## F-statistic: 11.92 on 1 and 130 DF,  p-value: 0.0007496
```

## Adjusted Reviews, $



## Adjusted Reviews, $$



## Adjusted Reviews, $$$



## Adjusted Reviews, $$$$



```
##
## Call:
## lm(formula = tslm_d4_resid ~ recession_dummy_dollars_m)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -158.13  -21.63    5.18   29.74  107.62
##
## Coefficients:
##                           Estimate Std. Error t value Pr(>|t|)
## (Intercept)                  4.168      4.353   0.958   0.3400
## recession_dummy_dollars_m  -28.960     11.472  -2.524   0.0128 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 46.27 on 130 degrees of freedom
## Multiple R-squared:  0.04673,    Adjusted R-squared:  0.03939
## F-statistic: 6.372 on 1 and 130 DF,  p-value: 0.0128

## NULL

## NULL

## NULL

## NULL
```

The regression results and error descriptive statistics do not look too promising, but could be leading in the right direction.

Let's try to make it better.

# Making a Better Model

By controlling for the aggegated number of reviews between dollar signs, we can eliminate the effect that a drop in aggregated reviews (again, by dollar signs) can potentially have.

```
##
## Call:
## lm(formula = tslm_d1_resid ~ recession_dummy_dollars_m + df_rev_count$coredata.df_rev_m.)
##
## Residuals:
##     Min     1Q  Median     3Q     Max
## -1829.6  -917.4  -307.4  1123.4  2837.2
##
## Coefficients:
##                                  Estimate Std. Error t value Pr(>|t|)
## (Intercept)                    -4.521e+02  1.852e+02  -2.441 0.016016 *
## recession_dummy_dollars_m       1.689e+02  3.291e+02   0.513 0.608659
## df_rev_count$coredata.df_rev_m.  2.751e-02  8.162e-03   3.370 0.000991 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1235 on 129 degrees of freedom
## Multiple R-squared:  0.0848, Adjusted R-squared:  0.07061
## F-statistic: 5.976 on 2 and 129 DF,  p-value: 0.003295

## NULL

##
## Call:
## lm(formula = tslm_d2_resid ~ recession_dummy_dollars_m + df_rev_count$coredata.df_rev_m.)
##
## Residuals:
##     Min      1Q  Median     3Q     Max
## -4138.1 -1943.8  -620.9  2367.0  6147.7
##
## Coefficients:
##                                  Estimate Std. Error t value Pr(>|t|)
## (Intercept)                    -977.68395  397.81376  -2.458 0.015313 *
## recession_dummy_dollars_m       378.55976  706.74989   0.536 0.593133
## df_rev_count$coredata.df_rev_m.   0.05937    0.01753   3.387 0.000938 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2653 on 129 degrees of freedom
## Multiple R-squared:  0.08534,    Adjusted R-squared:  0.07116
## F-statistic: 6.018 on 2 and 129 DF,  p-value: 0.003171

## NULL

##
## Call:
## lm(formula = tslm_d3_resid ~ recession_dummy_dollars_m + df_rev_count$coredata.df_rev_m.)
##
```
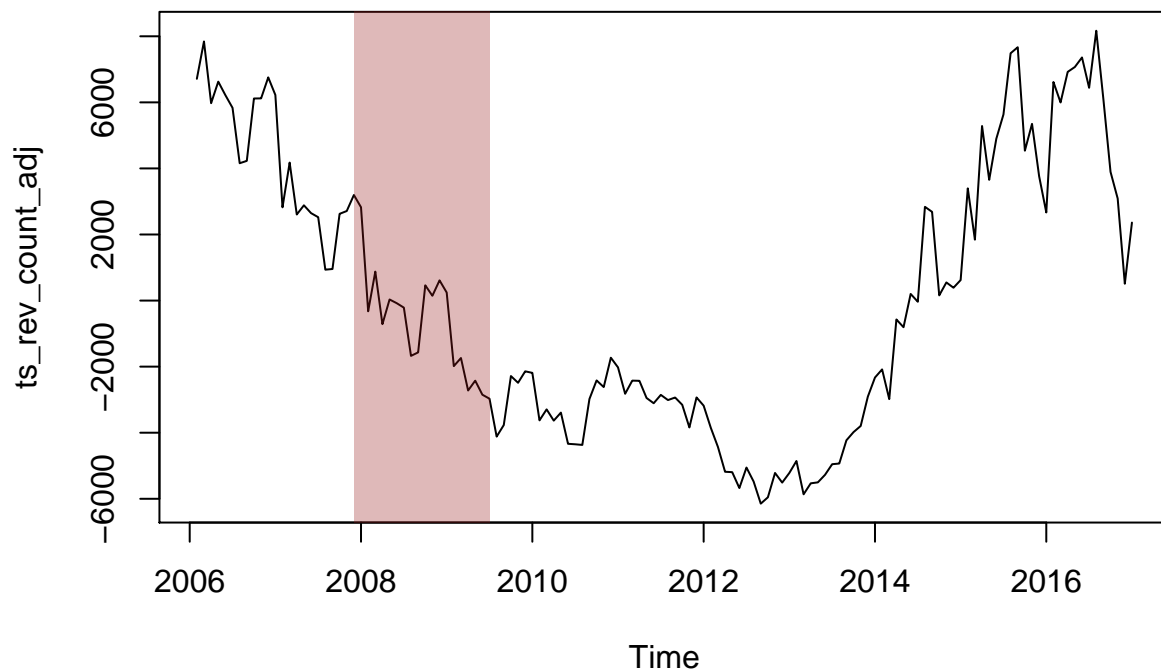
```
## Residuals:
##      Min       1Q   Median       3Q      Max
## -315.824  -83.622   -8.919   77.962  280.690
##
## Coefficients:
##                                 Estimate Std. Error t value Pr(>|t|)
## (Intercept)                    3.403e+00  1.754e+01    0.194  0.84646
## recession_dummy_dollars_m     -9.085e+01  3.116e+01   -2.916  0.00418 **
## df_rev_count$coredata.df_rev_m. 6.221e-04  7.729e-04    0.805  0.42231
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 117 on 129 degrees of freedom
## Multiple R-squared:  0.08857,    Adjusted R-squared:  0.07444
## F-statistic: 6.268 on 2 and 129 DF,  p-value: 0.002523

## NULL

##
## Call:
## lm(formula = tslm_d4_resid ~ recession_dummy_dollars_m + df_rev_count$coredata.df_rev_m.)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -145.400  -28.213    4.314   30.616  111.525
##
## Coefficients:
##                                  Estimate Std. Error t value Pr(>|t|)
## (Intercept)                     1.591e+01  6.836e+00    2.327  0.02152 *
## recession_dummy_dollars_m      -3.874e+01  1.214e+01   -3.190  0.00179 **
## df_rev_count$coredata.df_rev_m. -6.644e-04  3.012e-04   -2.206  0.02919 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 45.6 on 129 degrees of freedom
## Multiple R-squared:  0.08137,    Adjusted R-squared:  0.06712
## F-statistic: 5.713 on 2 and 129 DF,  p-value: 0.004195

## NULL
```

Adding new reviews in doesn't help much, but what about if we detrend and seasonally adjust?

```
##
## Call:
## tslm(formula = ts_rev_count ~ trend + season)
##
## Residuals:
##     Min       1Q  Median       3Q      Max
## -6146.8 -3161.7  -642.9  3115.8  8167.5
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -10135.277   1480.253  -6.847 3.52e-10 ***
## trend          356.093      9.876  36.056  < 2e-16 ***
## season2       3427.756   1839.239   1.864   0.0648 .
## season3       1767.844   1838.682   0.961   0.3383
```

```
## season4          3223.751     1838.178     1.754      0.0820 .
## season5          2287.839     1837.727     1.245      0.2156
## season6          2263.018     1837.329     1.232      0.2205
## season7          2320.834     1836.984     1.263      0.2089
## season8          3676.922     1836.692     2.002      0.0476 *
## season9          3374.556     1836.453     1.838      0.0686 .
## season10         1166.008     1836.267     0.635      0.5267
## season11          750.823     1836.134     0.409      0.6833
## season12         -197.361     1836.054    -0.107      0.9146
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4306 on 119 degrees of freedom
## Multiple R-squared:  0.9165,	Adjusted R-squared:  0.9081
## F-statistic: 108.9 on 12 and 119 DF,  p-value: < 2.2e-16
```



```
##
## Call:
## lm(formula = tslm_d1_resid ~ recession_dummy_dollars_m + ts_rev_count_adj)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -303.05  -50.07   -3.11   56.55  359.65
##
## Coefficients:
##                             Estimate Std. Error t value Pr(>|t|)
```

```
## (Intercept)                   -4.796724   9.627037  -0.498    0.619
## recession_dummy_dollars_m 33.324607  25.434935   1.310    0.192
## ts_rev_count_adj            0.311485   0.002184 142.630   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 102.3 on 129 degrees of freedom
## Multiple R-squared:  0.9937, Adjusted R-squared:  0.9936
## F-statistic: 1.021e+04 on 2 and 129 DF,  p-value: < 2.2e-16

## NULL

##
## Call:
## lm(formula = tslm_d2_resid ~ recession_dummy_dollars_m + ts_rev_count_adj)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -257.51  -74.84   -1.27   80.55  428.91
##
## Coefficients:
##                            Estimate Std. Error t value Pr(>|t|)
## (Intercept)              -12.196577  10.777970  -1.132  0.25989
## recession_dummy_dollars_m 84.734117  28.475734   2.976  0.00349 **
## ts_rev_count_adj           0.670808   0.002445 274.365  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 114.5 on 129 degrees of freedom
## Multiple R-squared:  0.9983, Adjusted R-squared:  0.9983
## F-statistic: 3.779e+04 on 2 and 129 DF,  p-value: < 2.2e-16

## NULL

##
## Call:
## lm(formula = tslm_d3_resid ~ recession_dummy_dollars_m + ts_rev_count_adj)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -351.65  -55.54   -0.17   58.72  244.76
##
## Coefficients:
##                            Estimate Std. Error t value Pr(>|t|)
## (Intercept)               12.307561   8.912974   1.381 0.169710
## recession_dummy_dollars_m -85.505163  23.548357  -3.631 0.000406 ***
## ts_rev_count_adj           0.016770   0.002022   8.294 1.25e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 94.71 on 129 degrees of freedom
## Multiple R-squared:  0.4026, Adjusted R-squared:  0.3933
## F-statistic: 43.47 on 2 and 129 DF,  p-value: 3.712e-15

## NULL

##
```

```
## Call:
## lm(formula = tslm_d4_resid ~ recession_dummy_dollars_m + ts_rev_count_adj)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -156.602  -22.164    6.943   27.480  101.189
##
## Coefficients:
##                             Estimate Std. Error t value Pr(>|t|)
## (Intercept)                4.662e+00  4.092e+00   1.139  0.25668
## recession_dummy_dollars_m -3.239e+01  1.081e+01  -2.996  0.00328 **
## ts_rev_count_adj          -3.962e-03  9.282e-04  -4.269 3.78e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 43.48 on 129 degrees of freedom
## Multiple R-squared:  0.1647, Adjusted R-squared:  0.1518
## F-statistic: 12.72 on 2 and 129 DF,  p-value: 9.082e-06

## NULL

## recession_dummy_dollars_m           ts_rev_count_adj
##                  1.005547                   1.005547

## recession_dummy_dollars_m           ts_rev_count_adj
##                  1.005547                   1.005547

## recession_dummy_dollars_m           ts_rev_count_adj
##                  1.005547                   1.005547

## recession_dummy_dollars_m           ts_rev_count_adj
##                  1.005547                   1.005547
```

Although not perfect, using detrended and seasonally adjusted data for the new reviews improves our results drastically. We finally have results that can be trustworthy.

So what do our new results tell us?

We see a decrease in ($$$$) and ($$$) reviews, but an increase in ($$) reviews. Not only do people prefer ($$) restaurants over the more expensive ($$$) and ($$$$) restaurants, it seems that customers who would originally have dined at the more expensive eateries are now choosing the less expensive ($$) restaurants. A surprising results is that the least expensive ($) restaurants do not see a significant change. This could be because inexpensive ($) restauarants are not substitutes for the others while the ($$) restaurants can be substitutes for ($$$) and ($$$$) restaurants.

Since reviews can be modeled when broken down by prices, this gives more promise to the previous review score analysis as long as it is also broken down in the same way.

## Breaking Down Review Scores by Price

After splitting our reviews by prices, a re-examination of review scores is due.
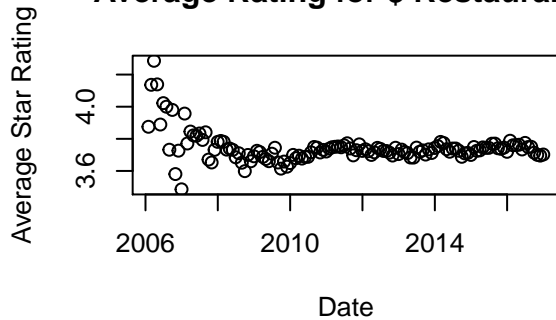
```
##
## Call:
## lm(formula = df_dollars_1_star$avg ~ recession_dummy_dollars_m)
##
## Residuals:
```

```
##     Min      1Q  Median      3Q     Max
## -0.26498 -0.03895 -0.01530  0.01052  0.53502
##
## Coefficients:
##                          Estimate Std. Error t value Pr(>|t|)
## (Intercept)              3.750694   0.009036 415.069   <2e-16 ***
## recession_dummy_dollars_m -0.044273   0.023818  -1.859   0.0653 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.09606 on 130 degrees of freedom
## Multiple R-squared:  0.02589,   Adjusted R-squared:  0.0184
## F-statistic: 3.455 on 1 and 130 DF,  p-value: 0.06532

##
## Call:
## lm(formula = df_dollars_2_star$avg ~ recession_dummy_dollars_m)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -0.16317 -0.06408 -0.01288  0.06248  0.17477
##
## Coefficients:
##                          Estimate Std. Error t value Pr(>|t|)
## (Intercept)              3.671420   0.007485 490.516  < 2e-16 ***
## recession_dummy_dollars_m -0.076331   0.019728  -3.869 0.000172 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.07956 on 130 degrees of freedom
## Multiple R-squared:  0.1033, Adjusted R-squared:  0.09637
## F-statistic: 14.97 on 1 and 130 DF,  p-value: 0.000172

##
## Call:
## lm(formula = df_dollars_3_star$avg ~ recession_dummy_dollars_m)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -0.25080 -0.05478 -0.00522  0.03954  0.40509
##
## Coefficients:
##                          Estimate Std. Error t value Pr(>|t|)
## (Intercept)              3.800796   0.008157 465.945   <2e-16 ***
## recession_dummy_dollars_m -0.020898   0.021501  -0.972    0.333
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.08671 on 130 degrees of freedom
## Multiple R-squared:  0.007215,   Adjusted R-squared:  -0.000422
## F-statistic: 0.9447 on 1 and 130 DF,  p-value: 0.3329
```
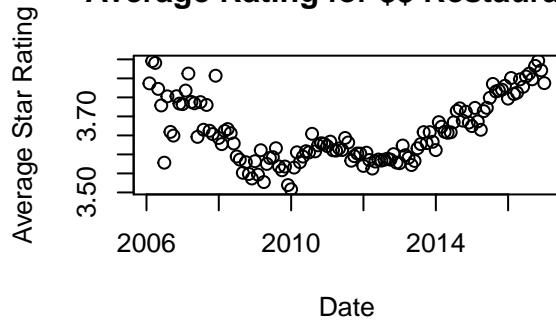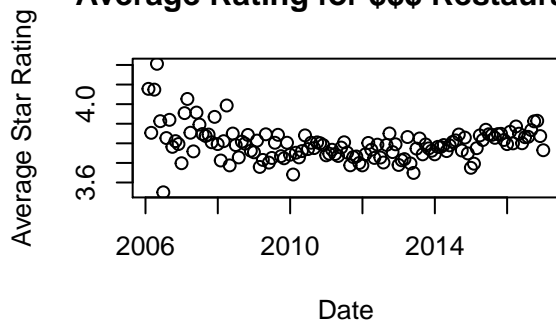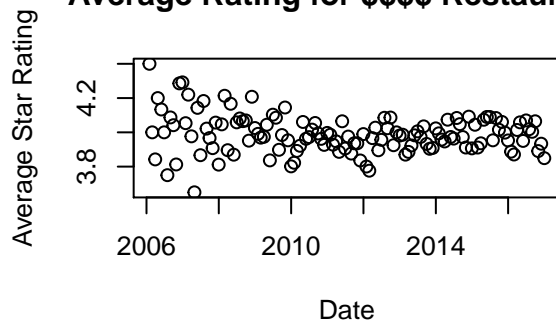
**Average Rating for $ Restaurants**

**Average Rating for $$ Restaurants**

**Average Rating for $$$ Restaurants**

**Average Rating for $$$$ Restaurants**

```
## 
## Call:
## lm(formula = df_dollars_4_star$avg ~ recession_dummy_dollars_m)
## 
## Residuals:
##      Min       1Q   Median       3Q      Max 
## -0.33494 -0.06462 -0.00311  0.06433  0.41506 
## 
## Coefficients:
##                           Estimate Std. Error t value Pr(>|t|)    
## (Intercept)                3.98494    0.01040  383.33   <2e-16 ***
## recession_dummy_dollars_m  0.03506    0.02740    1.28    0.203    
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 0.1105 on 130 degrees of freedom
## Multiple R-squared:  0.01244,    Adjusted R-squared:  0.004841 
## F-statistic: 1.637 on 1 and 130 DF,  p-value: 0.203

## NULL

## NULL

## NULL

## NULL
```
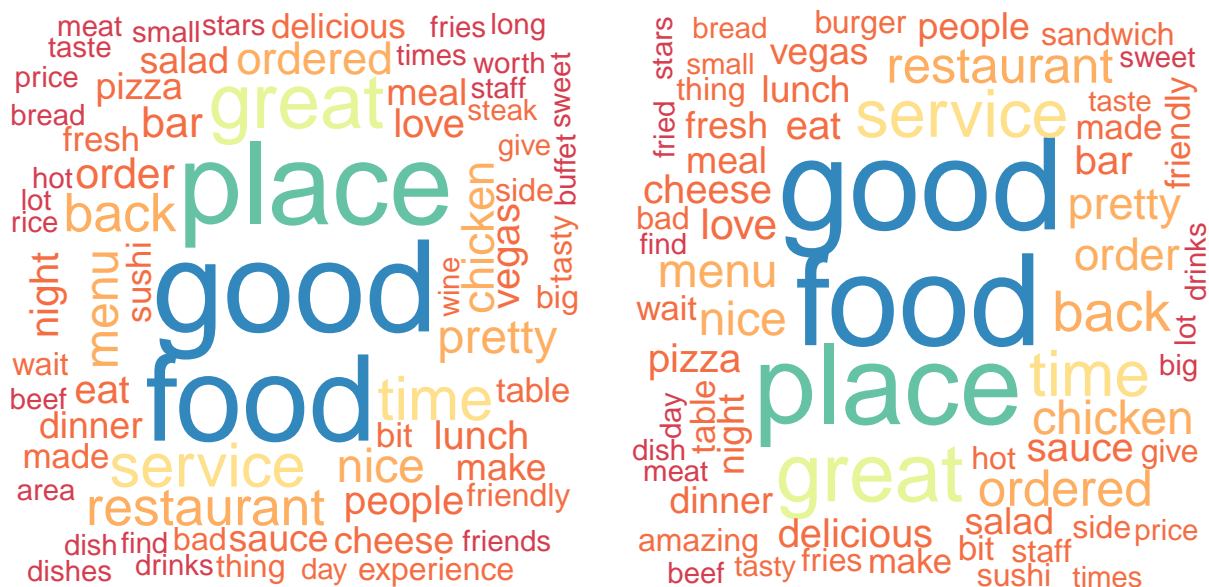
The results look better, but still are not that great. If we are to interpret the regression anyways, however, it looks as if the lower priced restaurants have lower review scores. This could be a sign of many things, such as

the substituters from ($$$) and ($$$$) restaurants having higher expectations or people just want more bang for their buck. The higher priced restaurants do not see a change, which could be due to the type of people who still dine there during recessions. They could be part of the group unaffected by the recession.

## Sentiment Analysis

By doing a sentiment analysis on the text in the reviews, we can see the association with words during the recession and the period of time (of equal length) directly after the recession.

Word Clouds:



It is hard to tell which word cloud is from the recession and which is from the period after, but the one on the left is from the recession.

Sentiment Visualization:

```
## # A tibble: 436,171 × 4
##        document          term count sentiment
##          <chr>          <chr> <dbl>     <chr>
## 1  character(0)          died     1  negative
## 2  character(0)     enthusiasm     1  positive
## 3  character(0)      fantastic     2  positive
## 4  character(0)           good     1  positive
## 5  character(0)       horrible     1  negative
## 6  character(0)           love     1  positive
## 7  character(0) recommendations     1  positive
## 8  character(0)           good     1  positive
```
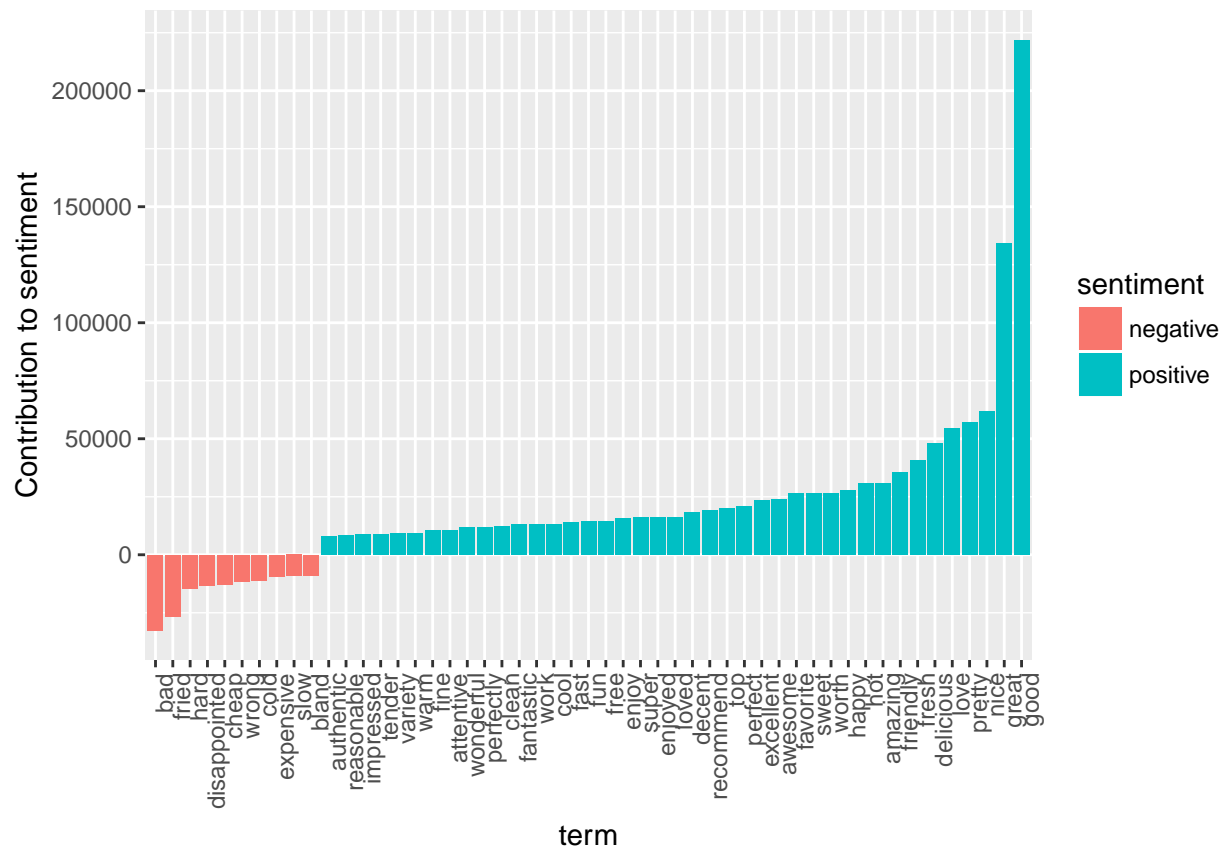
```
## 9  character(0)          nice     1  positive
## 10 character(0)          pure     1  positive
## # ... with 436,161 more rows

## # A tibble: 1 × 4
##      document negative positive sentiment
##         <chr>    <dbl>    <dbl>     <dbl>
## 1 character(0)   158714   346991    188277
```



```
## # A tibble: 2,121,102 × 4
##      document      term count sentiment
##         <chr>     <chr> <dbl>     <chr>
## 1  character(0) amazingly     1  positive
## 2  character(0)   awesome     1  positive
## 3  character(0)      fast     2  positive
## 4  character(0)    fucking     1  negative
## 5  character(0)      great     1  positive
## 6  character(0)      holy     1  positive
## 7  character(0)      nice     1  positive
## 8  character(0)      shit     1  negative
## 9  character(0)      weak     1  negative
## 10 character(0)      work     1  positive
## # ... with 2,121,092 more rows

## # A tibble: 1 × 4
##      document negative positive sentiment
##         <chr>    <dbl>    <dbl>     <dbl>
## 1 character(0)   733329  1731938    998609
```

Again, it is difficult to distinguish.

The Sentiment visualization on the top is from the recession. On closer examination of the sentiments, the following statistics are extracted.

Recession:
negative: 158714
positive: 346991
percent negative: 31.4%
cheap negative word rank: 3rd
expensive negative word rank: 7th

Post-Recession:
negative: 733329
positive: 1731938
percent negative: 29.75%
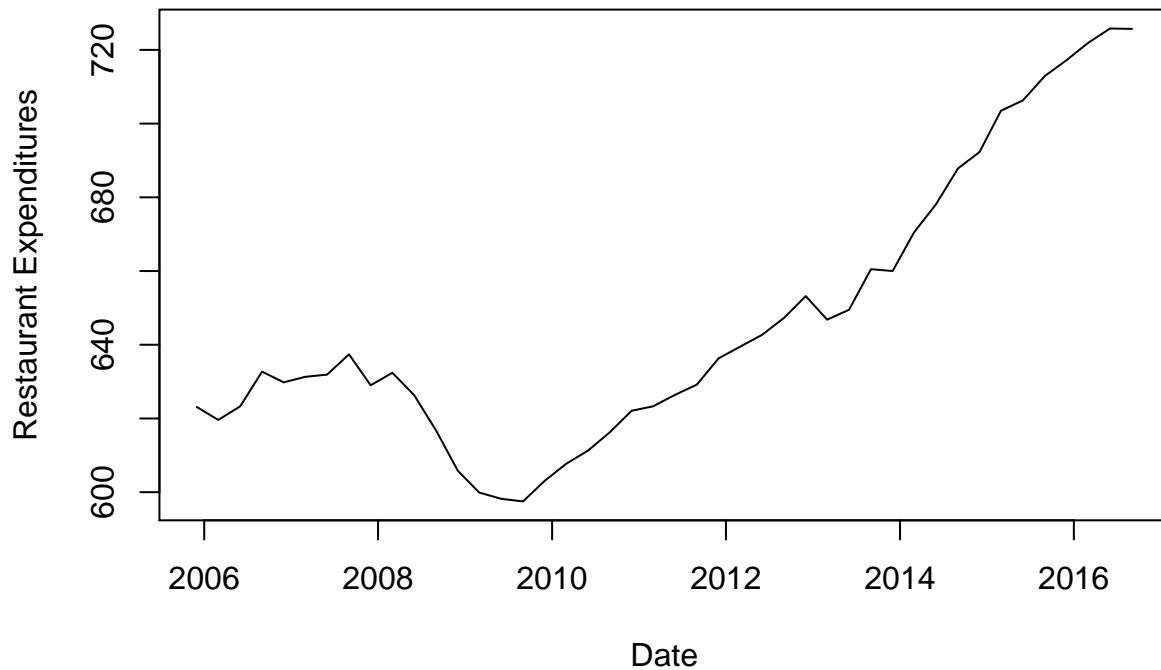cheap negative word rank: 5th
expensive negative word rank: 8th

There is a higher percentage of negative sentiment during the recession as well as having the words "cheap" and "expensive" rank higher for negative words. The word "cheap", however is not always used in a negative way, but it is still related to price. This shows that reviews are more concerned with price during the recession compared to the period directly following it. It should be noted that there is a large difference in the sample size of words from the two periods. This leads to a soft conclusion that reviewers are more concerned with prices during the recession, as it follows recessionary thinking.

# Connecting Yelp Reviews with the Restaurant Industry

By using restaurant expenditures, there can be a connection made between the review data and the actual restaurant industry.

But first, lets confirm that GDP and the recession can be linked to restaurant expenditures.

## Real Restaurant Expenditures, Quarterly



```
##
##   KPSS Test for Level Stationarity
##
## data:  t
## KPSS Level = 1.7047, Truncation lag parameter = 1, p-value = 0.01
##
##
##   Augmented Dickey-Fuller Test
##
## data:  t
## Dickey-Fuller = -2.4338, Lag order = 3, p-value = 0.4021
## alternative hypothesis: stationary

##
##   KPSS Test for Level Stationarity
##
## data:  t
## KPSS Level = 0.69776, Truncation lag parameter = 1, p-value =
## 0.01375
##
```

```
##
##  Augmented Dickey-Fuller Test
##
## data:  t
## Dickey-Fuller = -2.2003, Lag order = 3, p-value = 0.4946
## alternative hypothesis: stationary
```

## Real Restaurant Expenditures Growth Rate, Quarterly



```
##
## VAR Estimation Results:
## =========================
## Endogenous variables: rest_real_exp_diff_log, gdp_growth_subset
## Deterministic variables: const
## Sample size: 40
## Log Likelihood: 309.175
## Roots of the characteristic polynomial:
## 0.8295 0.7503 0.7503 0.7004 0.5468 0.5468
## Call:
## VAR(y = gdp_exp_combined, p = select$select[1])
##
##
## Estimation results for equation rest_real_exp_diff_log:
## ========================================================
## rest_real_exp_diff_log = rest_real_exp_diff_log.l1 + gdp_growth_subset.l1 + rest_real_exp_diff_log.l:
##
##                            Estimate Std. Error t value Pr(>|t|)
## rest_real_exp_diff_log.l1 -0.3061642  0.2190439  -1.398  0.17152
```

```
## gdp_growth_subset.l1        0.7168516  0.2564126   2.796  0.00857 **
## rest_real_exp_diff_log.l2  0.1645218  0.1996566   0.824  0.41584
## gdp_growth_subset.l2        0.1951121  0.2662260   0.733  0.46880
## rest_real_exp_diff_log.l3  0.1923951  0.1868020   1.030  0.31053
## gdp_growth_subset.l3        0.1265239  0.2690216   0.470  0.64123
## const                     -0.0001735  0.0013498  -0.129  0.89852
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##
## Residual standard error: 0.007011 on 33 degrees of freedom
## Multiple R-Squared: 0.4144,  Adjusted R-squared: 0.3079
## F-statistic: 3.892 on 6 and 33 DF,  p-value: 0.004777
##
##
## Estimation results for equation gdp_growth_subset:
## ===================================================
## gdp_growth_subset = rest_real_exp_diff_log.l1 + gdp_growth_subset.l1 + rest_real_exp_diff_log.l2 + g
##
##                          Estimate Std. Error t value Pr(>|t|)
## rest_real_exp_diff_log.l1 -0.096552   0.190925  -0.506   0.6164
## gdp_growth_subset.l1       0.472165   0.223497   2.113   0.0423 *
## rest_real_exp_diff_log.l2  0.071231   0.174027   0.409   0.6850
## gdp_growth_subset.l2      -0.009650   0.232051  -0.042   0.9671
## rest_real_exp_diff_log.l3  0.324927   0.162822   1.996   0.0543 .
## gdp_growth_subset.l3      -0.294587   0.234487  -1.256   0.2178
## const                      0.001649   0.001177   1.402   0.1703
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##
## Residual standard error: 0.006111 on 33 degrees of freedom
## Multiple R-Squared: 0.2688,  Adjusted R-squared: 0.1358
## F-statistic: 2.022 on 6 and 33 DF,  p-value: 0.09062
##
##
##
## Covariance matrix of residuals:
##                     rest_real_exp_diff_log gdp_growth_subset
## rest_real_exp_diff_log          4.915e-05          2.935e-05
## gdp_growth_subset               2.935e-05          3.734e-05
##
## Correlation matrix of residuals:
##                     rest_real_exp_diff_log gdp_growth_subset
## rest_real_exp_diff_log             1.0000            0.6851
## gdp_growth_subset                  0.6851            1.0000

## Granger causality test
##
## Model 1: rest_real_exp_diff_log ~ Lags(rest_real_exp_diff_log, 1:3) + Lags(gdp_growth_subset[1:lengt
## Model 2: rest_real_exp_diff_log ~ Lags(rest_real_exp_diff_log, 1:3)
##   Res.Df Df      F  Pr(>F)
## 1     33
## 2     36 -3 2.8912 0.05002 .
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

## Granger causality test
##
## Model 1: gdp_growth_subset[1:length(gdp_growth_subset)] ~ Lags(gdp_growth_subset[1:length(gdp_growth
## Model 2: gdp_growth_subset[1:length(gdp_growth_subset)] ~ Lags(gdp_growth_subset[1:length(gdp_growth
##   Res.Df Df      F Pr(>F)
## 1     33
## 2     36 -3 1.4038 0.2591

##
## Call:
## lm(formula = rest_real_exp_diff_log ~ rec_exp_diff_log_dummy)
##
## Residuals:
##        Min         1Q     Median         3Q        Max
## -0.0158435 -0.0038443 -0.0000908  0.0033176  0.0144583
##
## Coefficients:
##                         Estimate Std. Error t value Pr(>|t|)
## (Intercept)             0.005997   0.001061   5.654 1.34e-06 ***
## rec_exp_diff_log_dummy -0.015065   0.002629  -5.730 1.05e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.006364 on 41 degrees of freedom
## Multiple R-squared:  0.4447, Adjusted R-squared:  0.4312
## F-statistic: 32.84 on 1 and 41 DF,  p-value: 1.046e-06

## NULL
```

The assumption that restaurant expenditures are related to the recession and GDP is confirmed. GDP will not be needed in constructing a full model as restaurant expenditures will be used in its place.

By seasonally adjusting the growth rate in new reviews, a Granger causality test can be run between restaurant expenditures and new reviews.

```
##
## Call:
## tslm(formula = log_rev_quarter ~ season)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.34782 -0.07594 -0.00878  0.04565  0.59583
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.35264    0.05057   6.973 2.34e-08 ***
## season2     -0.34070    0.06988  -4.876 1.85e-05 ***
## season3     -0.15807    0.06988  -2.262   0.0293 *
## season4     -0.41248    0.06988  -5.903 7.06e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1599 on 39 degrees of freedom
## Multiple R-squared:  0.5214, Adjusted R-squared:  0.4846
```

```
## F-statistic: 14.16 on 3 and 39 DF,  p-value: 2.155e-06

##
## VAR Estimation Results:
## =========================
## Endogenous variables: rest_real_exp_diff_log, log_rev_quarter_adj
## Deterministic variables: const
## Sample size: 39
## Log Likelihood: 196.578
## Roots of the characteristic polynomial:
## 0.9055 0.828 0.721 0.721 0.5383 0.5383 0.4846 0.4846
## Call:
## VAR(y = rev_exp_adj_combined, p = select$select[1])
##
##
## Estimation results for equation rest_real_exp_diff_log:
## ========================================================
## rest_real_exp_diff_log = rest_real_exp_diff_log.l1 + log_rev_quarter_adj.l1 + rest_real_exp_diff_log
##
##                           Estimate Std. Error t value Pr(>|t|)
## rest_real_exp_diff_log.l1  0.180897   0.180663   1.001    0.325
## log_rev_quarter_adj.l1    -0.010709   0.014668  -0.730    0.471
## rest_real_exp_diff_log.l2  0.239154   0.178751   1.338    0.191
## log_rev_quarter_adj.l2    -0.004026   0.011875  -0.339    0.737
## rest_real_exp_diff_log.l3  0.099466   0.168370   0.591    0.559
## log_rev_quarter_adj.l3     0.003002   0.012594   0.238    0.813
## rest_real_exp_diff_log.l4  0.026576   0.166231   0.160    0.874
## log_rev_quarter_adj.l4    -0.011961   0.011398  -1.049    0.302
## const                      0.001700   0.001486   1.144    0.262
##
##
## Residual standard error: 0.007706 on 30 degrees of freedom
## Multiple R-Squared: 0.3411,  Adjusted R-squared: 0.1654
## F-statistic: 1.941 on 8 and 30 DF,  p-value: 0.09013
##
##
## Estimation results for equation log_rev_quarter_adj:
## =====================================================
## log_rev_quarter_adj = rest_real_exp_diff_log.l1 + log_rev_quarter_adj.l1 + rest_real_exp_diff_log.l2
##
##                           Estimate Std. Error t value Pr(>|t|)
## rest_real_exp_diff_log.l1 -2.47221    1.57285  -1.572  0.12649
## log_rev_quarter_adj.l1     0.15930    0.12770   1.247  0.22186
## rest_real_exp_diff_log.l2  2.79746    1.55620   1.798  0.08231 .
## log_rev_quarter_adj.l2     0.22813    0.10338   2.207  0.03513 *
## rest_real_exp_diff_log.l3 -0.09956    1.46582  -0.068  0.94630
## log_rev_quarter_adj.l3    -0.15938    0.10964  -1.454  0.15643
## rest_real_exp_diff_log.l4 -2.85185    1.44720  -1.971  0.05806 .
## log_rev_quarter_adj.l4     0.28534    0.09923   2.875  0.00735 **
## const                     -0.01439    0.01293  -1.112  0.27477
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##
```
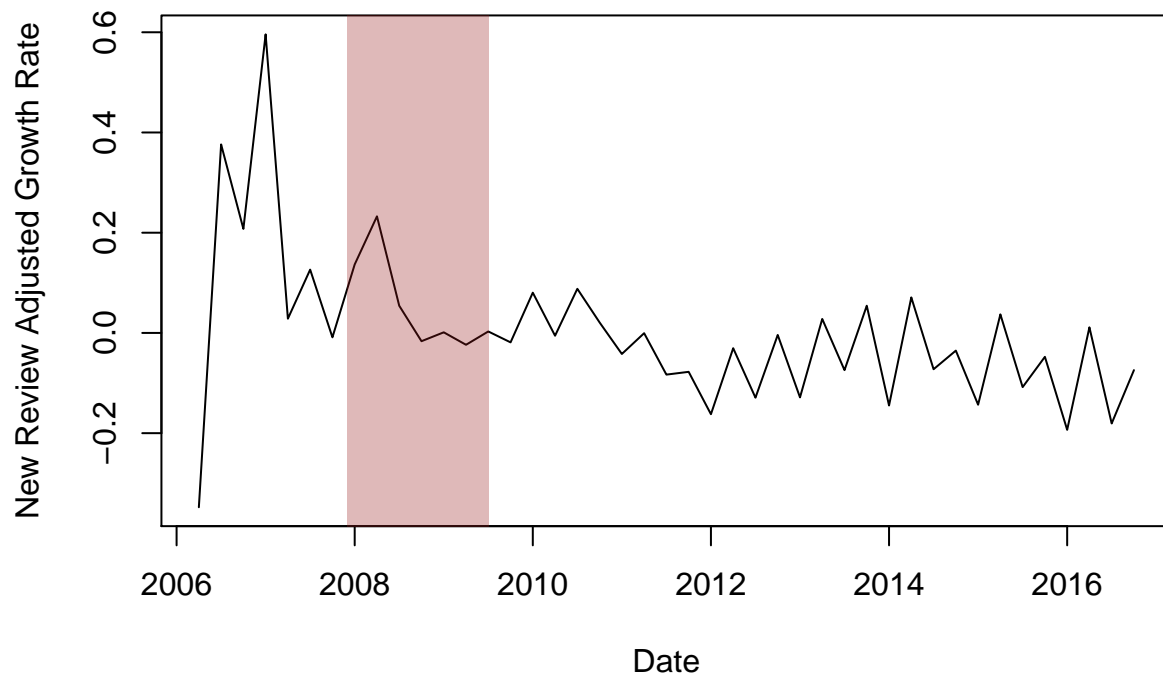
```
## Residual standard error: 0.06709 on 30 degrees of freedom
## Multiple R-Squared: 0.5769,  Adjusted R-squared: 0.4641
## F-statistic: 5.114 on 8 and 30 DF,  p-value: 0.0004486
##
##
##
## Covariance matrix of residuals:
##                     rest_real_exp_diff_log log_rev_quarter_adj
## rest_real_exp_diff_log              5.938e-05            0.000157
## log_rev_quarter_adj                 1.570e-04            0.004501
##
## Correlation matrix of residuals:
##                     rest_real_exp_diff_log log_rev_quarter_adj
## rest_real_exp_diff_log              1.0000            0.3038
## log_rev_quarter_adj                 0.3038            1.0000

## Granger causality test
##
## Model 1: rest_real_exp_diff_log ~ Lags(rest_real_exp_diff_log, 1:4) + Lags(log_rev_quarter_adj, 1:4)
## Model 2: rest_real_exp_diff_log ~ Lags(rest_real_exp_diff_log, 1:4)
##   Res.Df Df      F Pr(>F)
## 1     30
## 2     34 -4 0.6949 0.6014

## Granger causality test
##
## Model 1: log_rev_quarter_adj ~ Lags(log_rev_quarter_adj, 1:4) + Lags(rest_real_exp_diff_log, 1:4)
## Model 2: log_rev_quarter_adj ~ Lags(log_rev_quarter_adj, 1:4)
##   Res.Df Df      F  Pr(>F)
## 1     30
## 2     34 -4 2.2997 0.08183 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

## New Review Adjusted Growth Rate, Quarterly



New reviews is Granger caused by restaurant expenditures. This makes sense as reviews are typically made after people visit a restaurant. This shows that the number of reviews can in fact be related to the previous restaurant expenditures.

## Building Full Models

A full model that incorporates all the previous results will be useful to show the real effect of the recession and determine if there is a causal inference that can be made.

```
##
## Call:
## lm(formula = ts_d1_quarterly ~ rec_q + ts_rev_count_adj_q + ts_d1_star_quarterly +
##     rest_exp_q_adj)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -569.29  -97.37  -17.66  121.28  408.08
##
## Coefficients:
##                      Estimate Std. Error t value Pr(>|t|)
## (Intercept)         -8.777e+02  1.536e+03  -0.572    0.571
## rec_q                5.358e+01  8.628e+01   0.621    0.538
## ts_rev_count_adj_q   3.034e-01  5.399e-03  56.207   <2e-16 ***
## ts_d1_star_quarterly 7.738e+01  1.364e+02   0.567    0.574
## rest_exp_q_adj       5.007e+00  3.017e+00   1.660    0.105
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 201.6 on 39 degrees of freedom
## Multiple R-squared:  0.9975, Adjusted R-squared:  0.9972
## F-statistic:  3851 on 4 and 39 DF,  p-value: < 2.2e-16

## NULL

##
## Call:
## lm(formula = ts_d2_quarterly ~ rec_q + ts_rev_count_adj_q + ts_d2_star_quarterly +
##     rest_exp_q_adj)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -502.56 -211.31    6.58  205.47  656.43
##
## Coefficients:
##                       Estimate Std. Error t value Pr(>|t|)
## (Intercept)          -2.595e+03  4.914e+03  -0.528   0.6005
## rec_q                 2.731e+02  1.440e+02   1.896   0.0654 .
## ts_rev_count_adj_q    6.637e-01  8.055e-03  82.387   <2e-16 ***
## ts_d2_star_quarterly  2.323e+02  4.462e+02   0.521   0.6056
## rest_exp_q_adj        2.197e+00  4.883e+00   0.450   0.6552
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 277.9 on 39 degrees of freedom
## Multiple R-squared:  0.999,  Adjusted R-squared:  0.9989
## F-statistic:  9367 on 4 and 39 DF,  p-value: < 2.2e-16

## NULL

##
## Call:
## lm(formula = ts_d3_quarterly ~ rec_q + ts_rev_count_adj_q + ts_d3_star_quarterly +
##     rest_exp_q_adj)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -568.78 -164.74    9.84  174.17  465.42
##
## Coefficients:
##                       Estimate Std. Error t value Pr(>|t|)
## (Intercept)           4.254e+03  3.553e+03   1.197   0.2384
## rec_q                -2.089e+02  1.039e+02  -2.011   0.0513 .
## ts_rev_count_adj_q    2.963e-02  6.739e-03   4.397 8.22e-05 ***
## ts_d3_star_quarterly -3.705e+02  3.116e+02  -1.189   0.2418
## rest_exp_q_adj       -6.025e+00  3.829e+00  -1.573   0.1237
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 247.2 on 39 degrees of freedom
## Multiple R-squared:  0.5027, Adjusted R-squared:  0.4517
## F-statistic: 9.856 on 4 and 39 DF,  p-value: 1.311e-05
```

```
## NULL

##
## Call:
## lm(formula = ts_d4_quarterly ~ rec_q + ts_rev_count_adj_q + ts_d4_star_quarterly +
##     rest_exp_q_adj)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -316.46  -67.73   15.81   65.89  215.63
##
## Coefficients:
##                       Estimate Std. Error t value Pr(>|t|)
## (Intercept)          -9.689e+02  1.208e+03  -0.802   0.4273
## rec_q                -7.340e+01  4.909e+01  -1.495   0.1429
## ts_rev_count_adj_q    1.116e-03  3.115e-03   0.358   0.7221
## ts_d4_star_quarterly  8.192e+01  1.010e+02   0.811   0.4224
## rest_exp_q_adj       -3.740e+00  1.710e+00  -2.187   0.0348 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 114.6 on 39 degrees of freedom
## Multiple R-squared:  0.2986, Adjusted R-squared:  0.2267
## F-statistic: 4.152 on 4 and 39 DF,  p-value: 0.006758

## NULL

##              rec_q   ts_rev_count_adj_q ts_d1_star_quarterly
##           1.078443             4.614419             1.318413
##     rest_exp_q_adj
##           4.514264

##              rec_q   ts_rev_count_adj_q ts_d2_star_quarterly
##           1.581756             5.405495             6.411727
##     rest_exp_q_adj
##           6.224056

##              rec_q   ts_rev_count_adj_q ts_d3_star_quarterly
##           1.039758             4.779719             2.219843
##     rest_exp_q_adj
##           4.836175

##              rec_q   ts_rev_count_adj_q ts_d4_star_quarterly
##           1.079835             4.752153             1.268115
##     rest_exp_q_adj
##           4.487196
```
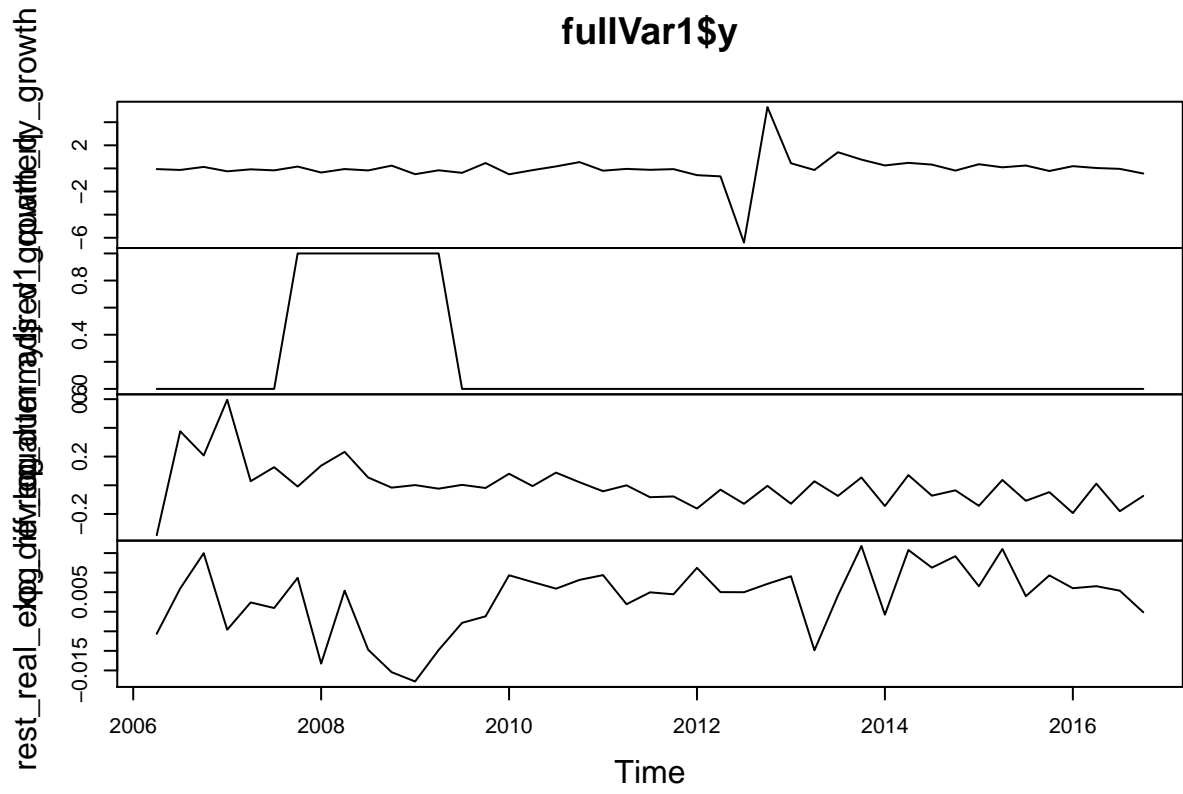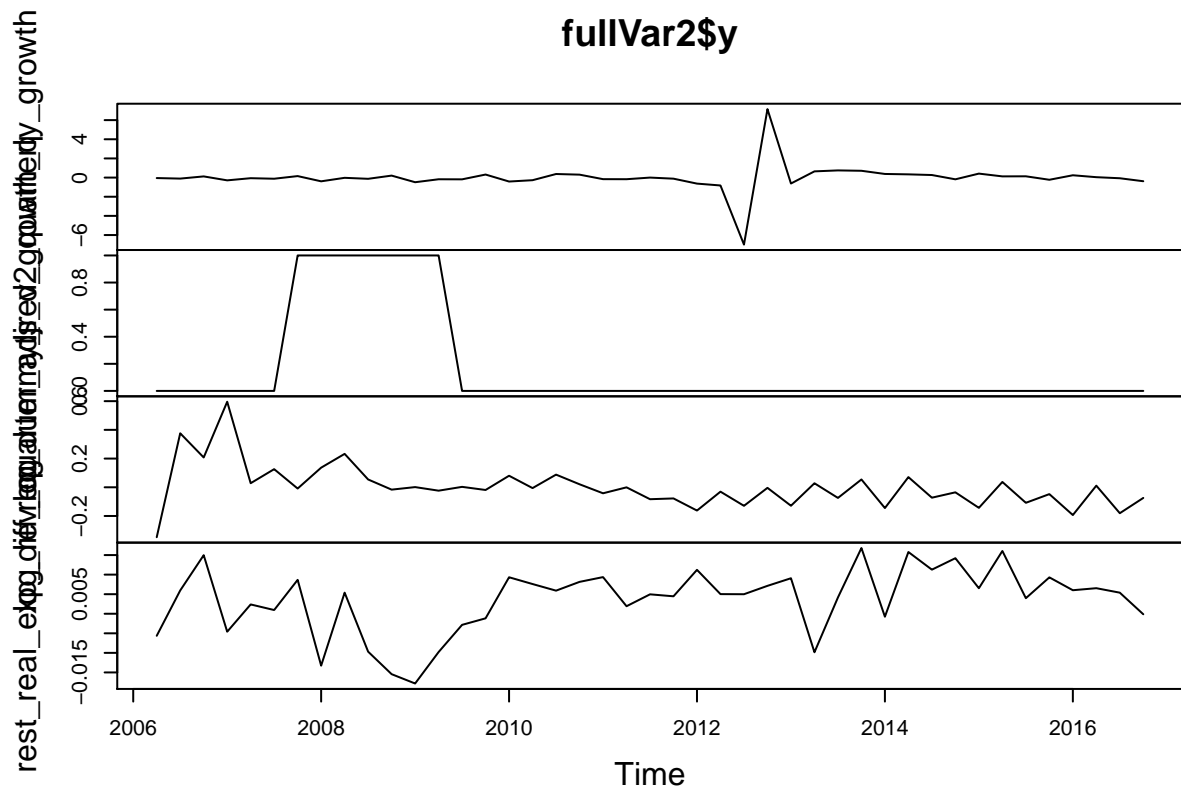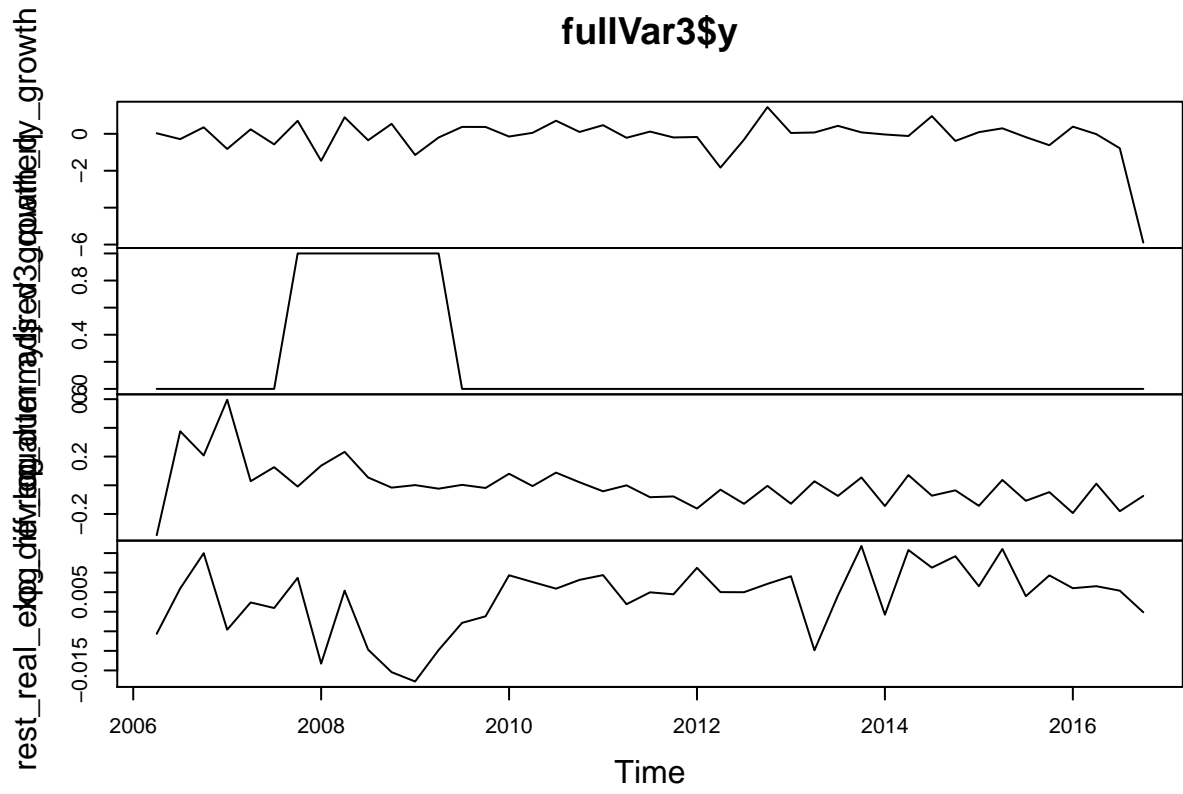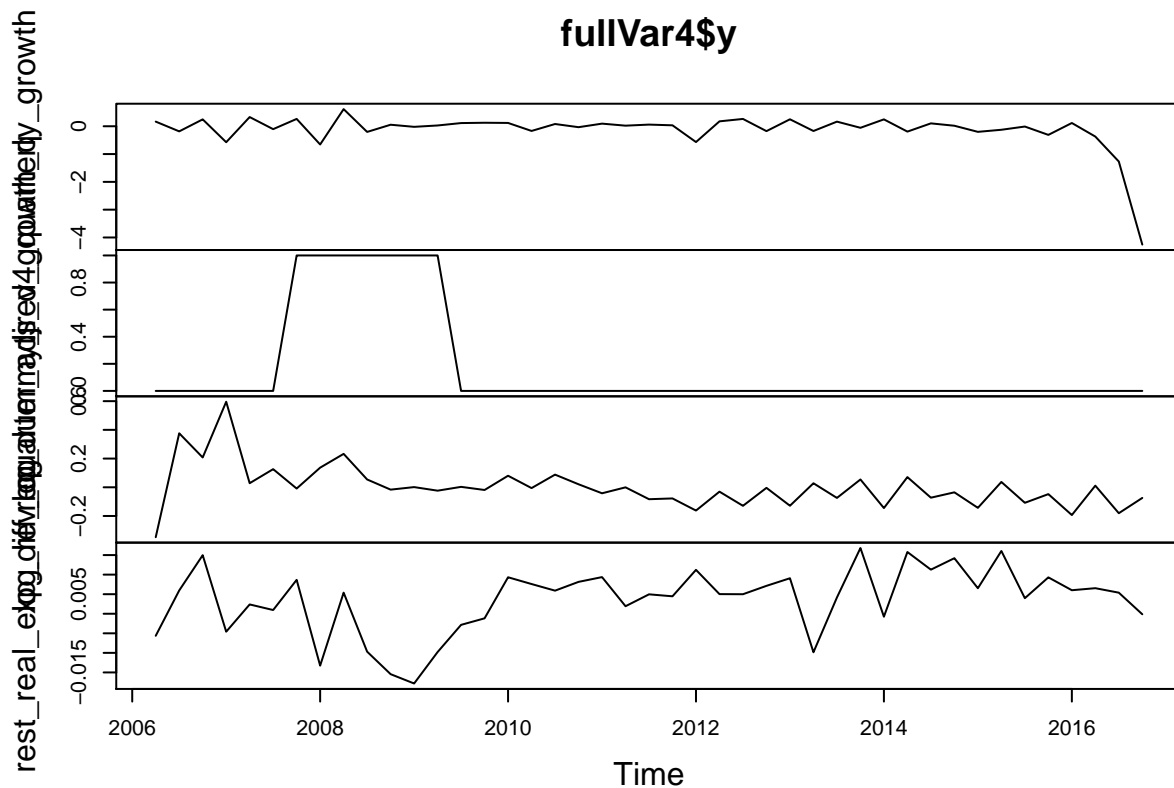
The regression results and error descriptive statistics are actually worse in this model. There is even multicolinearity occuring. The better model is actually in the form of Restaurant Review (dollars) regressed on Recession and Adjusted New Reviews.

Since new reviews is Granger caused by restaurant expenditures, can we find a causal inference through comprehensive VAR models that include the "Better" model, with restaurant expenditures added in?

**fullVar1$y**

**fullVar2$y**

**fullVar3$y**

**fullVar4$y**

The results of the VAR models (surpressed) are not very assuring. They don't point to any real causal inferences. More research and work needs to be done to create a truley causal model.

## Conclusion

The most effective model has been using the amount of reviews by dollar signs and regressing them on the recession and seasonally adjusted new reviews. These models showed that there is a drop in reviews for the higher priced ($$$) and ($$$$) restaurants while there is an increase in ($$) restaurants, with ($) restaurants being unaffected. This means that the ($$) restaurants are a substitute for the more expensive restaurants and that the lowest priced ($) restaurants are not interchangable with the other three categories.

Sentiment analysis on the review text showed that there was a possibility that reviews became more negative and price focused during the recession. However, the differences in sample size and potentially statistically insignificance of the results leave the sentiment analysis without a definite answer.

It was determined that the VAR model was not very effective in showing potential causality.

## Sources

YELP: https://www.yelp.com/dataset_challenge
FRED: https://fred.stlouisfed.org/
BEA: https://www.bea.gov/iTable/iTable.cfm?reqid=9&step=1&acrdn=2#reqid=9&step=1&isuri=1&904=2004&903=64&906=q&905=2016&910=x&911=0

Yahoo Finance: https://finance.yahoo.com/quote/YELP?p=YELP