

PSTAT 174 - Time Series Final Project: Total Generation of Electricity by US Electric Industry

Emeric Szaboky

12/6/2017

Abstract

This study focuses on the analysis and forecasting of monthly total electricity generation by the US electric industry, using data as provided by Makridakis, Wheelwright and Hyndman (1998). This time series spans from January 1985 to October 1996. The last two months of the final year 1996, are not included within the data. The primary concept of a study like this one is to forecast the future pattern of electric energy generation in the US, with the accompanying belief that it will continue to steadily increase due to demand caused by new technology. The project is designed, however, to forecast values for 1996 which were documented already in the data. This will allow the reader to verify the efficacy of this data in terms of forecast. The data will be split into a test set containing all 142 observations and training set with the last 10 observations (of 1996) missing. The training set of 132 observations will be used to predict and forecast the next ten months of data. The forecast is successful, offering a powerful glimpse into the future of the electric industry. This preliminary research is a helpful beginning for deeper exploration into the future trends of this data.

The Box-Jenkins method is used to perform analysis and forecast the total electricity generation ten months into the future. Multiple transformations are explored, with the Box-Cox transformation being chosen in order to stabilize variance. After this preliminary transformation, the time series is differenced once at lag 1 in order to remove its additive linear trend and a second time at lag 12 in order to remove its seasonal component. The data is confirmed to be stationary, and two models are chosen, one by analysis and observation of the ACF and PACF plots, and another with a computer-programmed function.

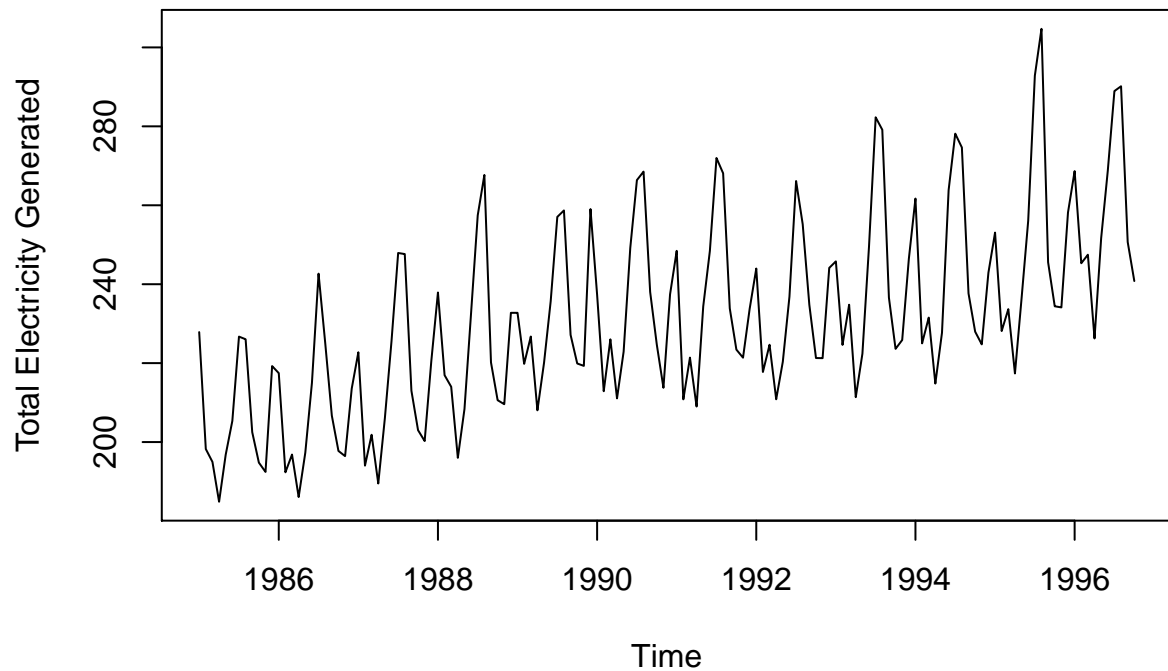
The two models are compared by AICc and diagnostics are explored in order to decide whether the models are of good enough quality to use for forecasting. The computer-programmed model is the final model chosen to forecast the 10-months ahead in 1996. The forecast projects a steady incline in the generation of electricity, however with smaller seasonal spikes than most of the years before.

Introduction

In this project, monthly time series data regarding the total electricity generation by the US electric industry will be analyzed and used to forecast the direction of the data 10-months into 1996. This project will explore the efficacy of this time series data in accurately predicting 10-months ahead of future values in US electricity generation. The data in this time series represents a fairly linear increase and suggests a forecast of a continual increase in the production of electricity, however with the forecasts made below, it seems probable that further exploration into the future of this data may show that rates of increase in electricity generation may be on the decline. This data is very effective in accomplishing its goal of projecting the 1996 year of values.

Examination of Time Series

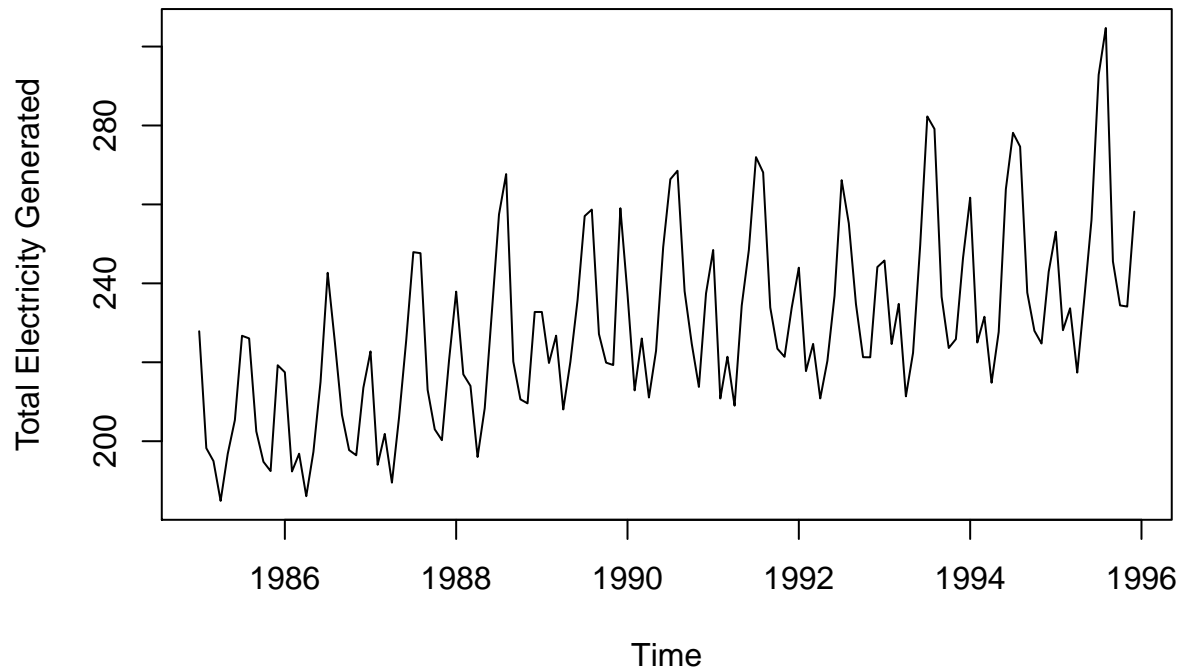
Total Generation of Electricity Per Month in US (85–96), Test



```
## [1] "Summary of Test Set Time Series"
```

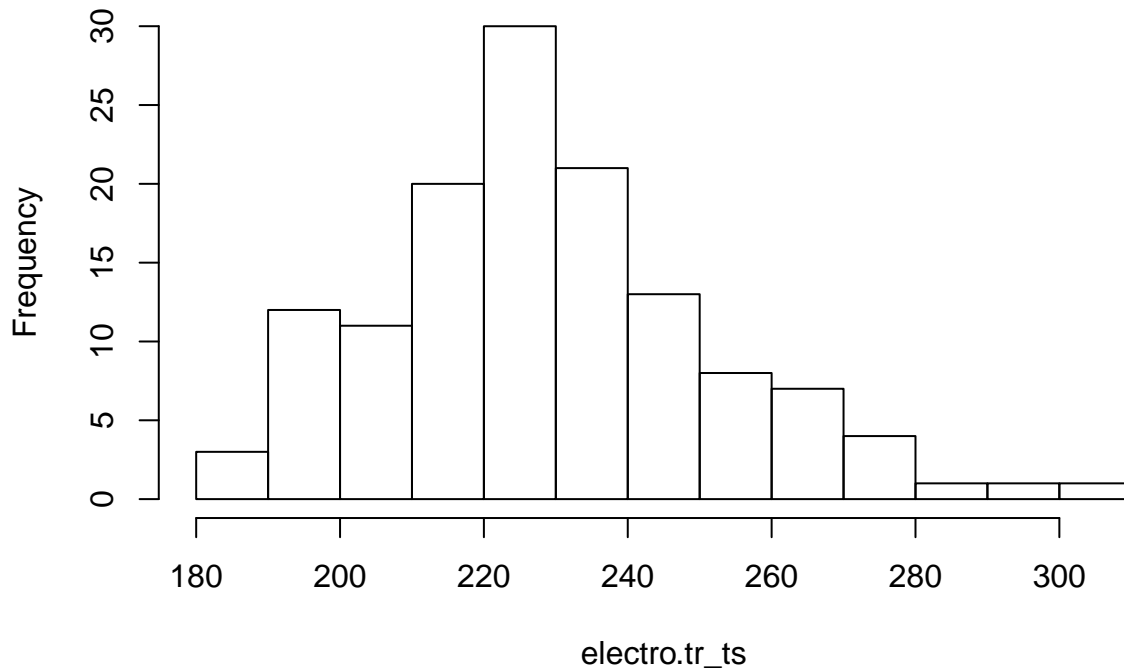
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 184.9   214.2   226.7   231.1   246.3   304.7
```

Total Generation of Electricity Per Month in US (85–96), Train



```
## [1] "Summary of Training Set Time Series"
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 184.9   213.4   225.7   229.1   243.2   304.7
```

Histogram: Time Series of Training Data Set



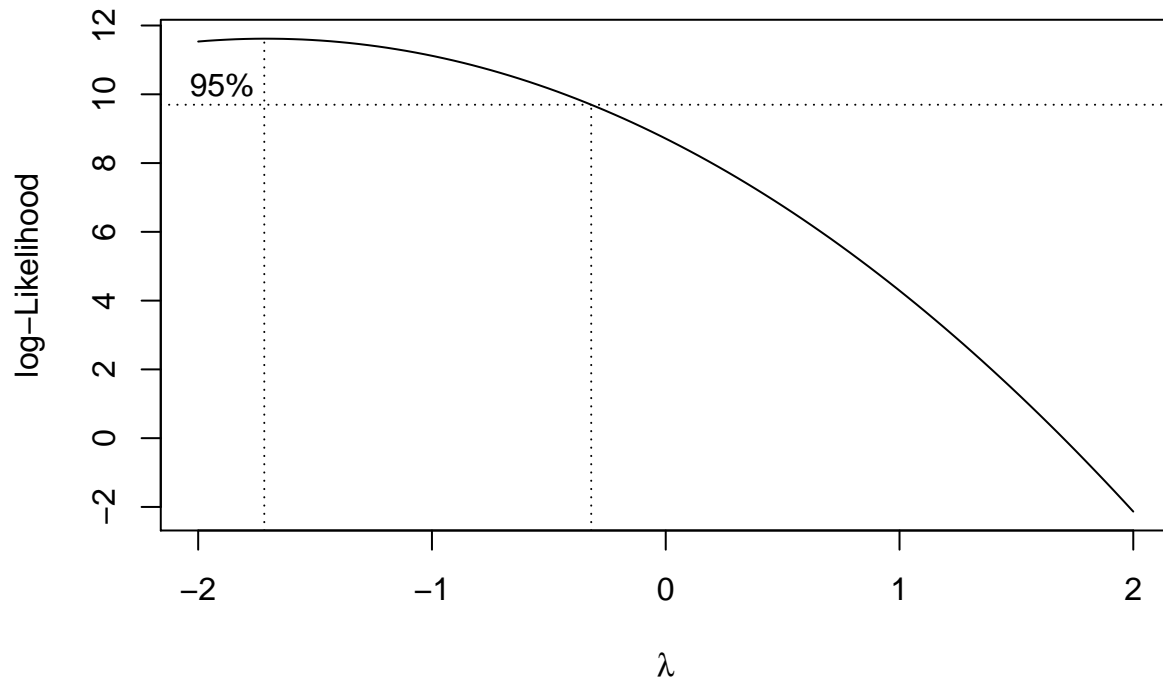
Above, the plotted time series of both the test set (all observations included) and the training set (devoid of the final 10 observations) can be seen. For the purposes of this study, we will use the training dataset ($n = 132$), and compare its resulting forecasted predictions with the those observations in the testing dataset. Observing the above plot of the time series, we can see that there is a slightly-damped upward linear trend with additive seasonality. By interpreting the dataset and this plot of the time series, we can make the assumption that the period for this seasonal trend is 12 (months), or 1 year. The data recorded is monthly data and there are 12 months in a year; it would make sense that electricity generation, which would be likely to correspond with seasonal shifts, would have a year long period. The above plot supports this assumption. There aren't many significant sharp changes or detours in behavior illustrated, however we do see the upward trend plateau for a short period between 1990-1993. The upward trend then resumes during 1993. Additionally, we see slightly larger spikes in 1989, 1993, and 1995 (the largest being 1995).

Transformations

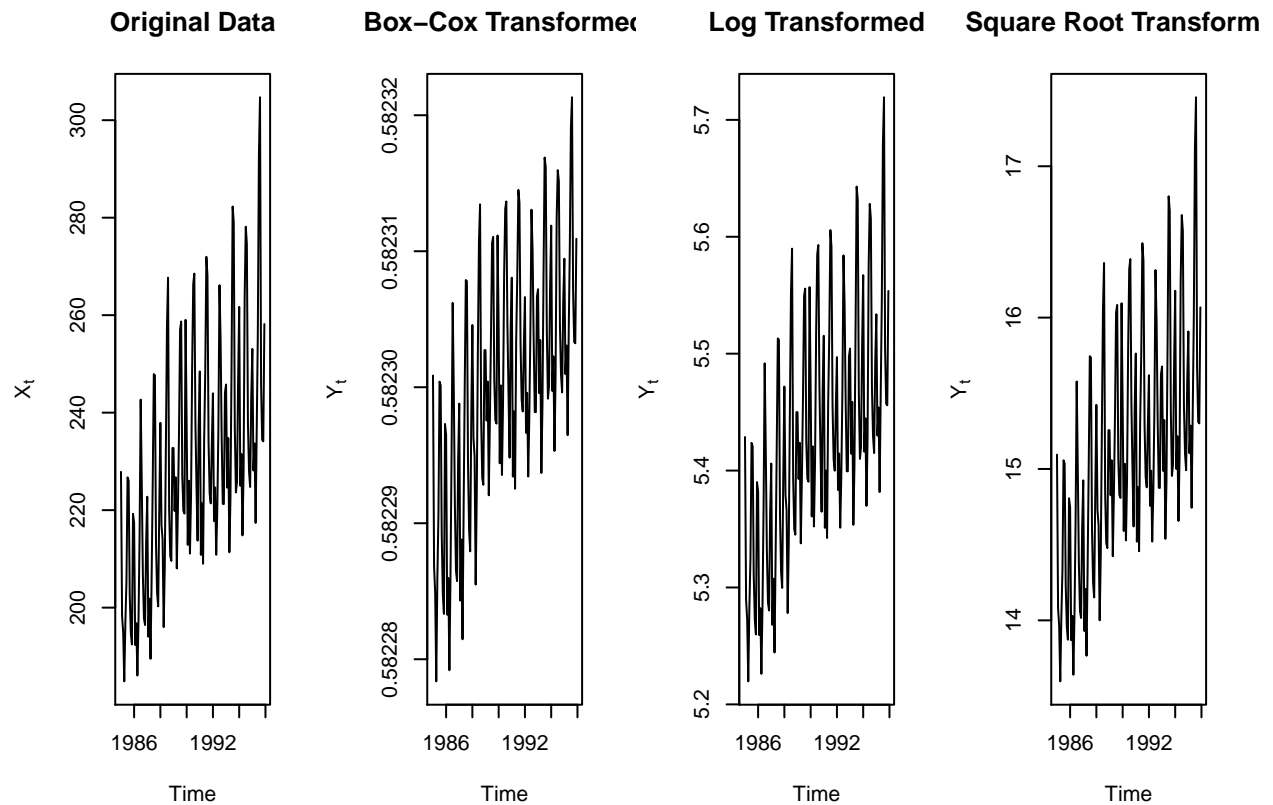
Primary Transformations

In order to stabilize the variance within the data, we will conduct a Box-Cox transformation, a log transformation, and a square root transformation, and decide which is most appropriate.

```
## [1] "Box-Cox Transformation: "
```



```
## [1] "The value of lambda is: -1.71717171717172"
```



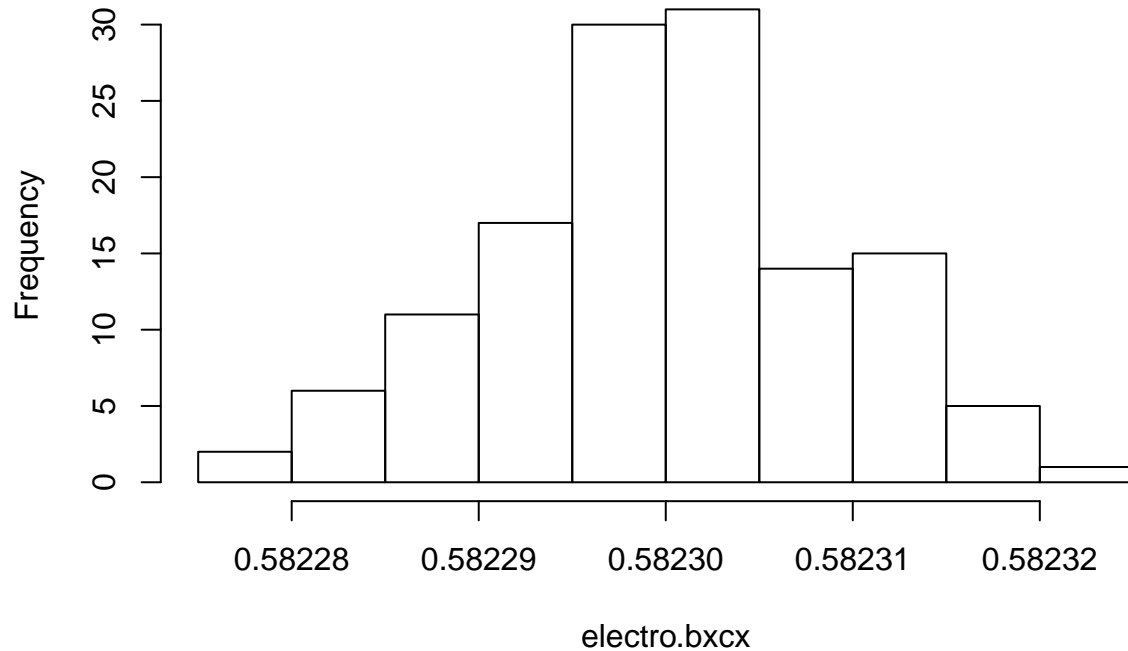
```
## [1] "The variance of the original training data is: 550.748634299098"
```

```
## [1] "The variance of the Box-Cox transformed data is: 8.00514724446019e-11"
```

```
## [1] "The variance of the log transformed data is: 0.0101433349554827"
```

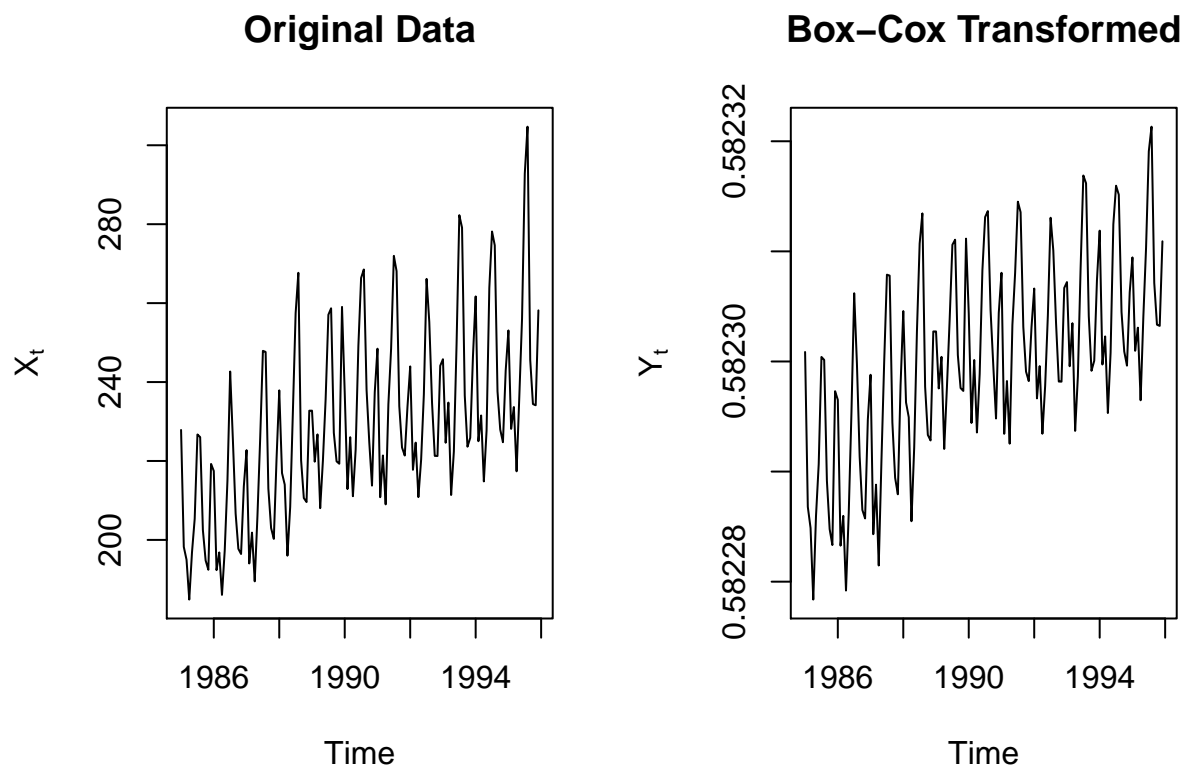
```
## [1] "The variance of the square-root transformed data is: 0.588835055533459"
```

Histogram: Box-Cox Transformed Data

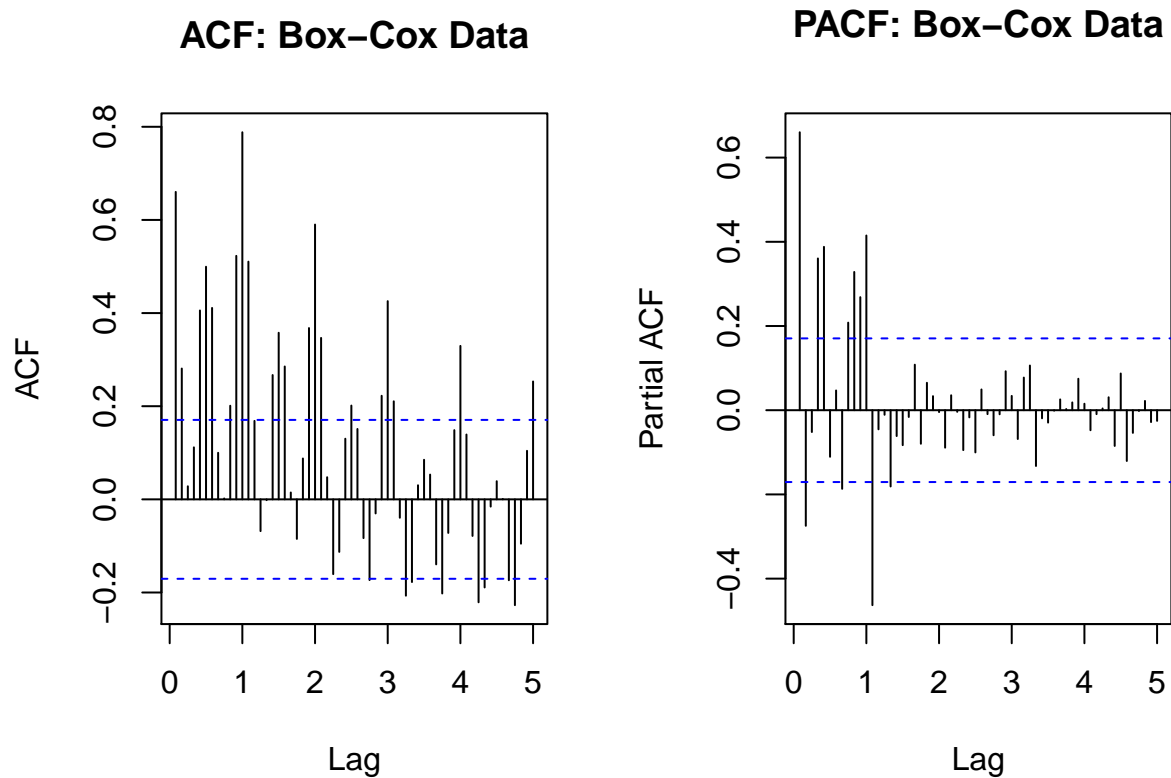


With our results above, we can see that the variance has been dramatically reduced for all 3 transformations. The transformation with the lowest variance, however, is the Box-Cox transformation. The Box-Cox algorithm chose a lambda value of -1.717172, which makes it a reciprocal power transformation. The histogram plot of the Box-Cox transformation of the time series appears to look significantly more like normal, although it is still important to detrend/deseasonalize and stationarize the data. From the above information, the smallest variance and the more normal-looking plot, we can conclude that the Box-Cox transformation is the most suitable method for stabilization of variance.

Chosen Transformation: Box-Cox Transformed Time Series



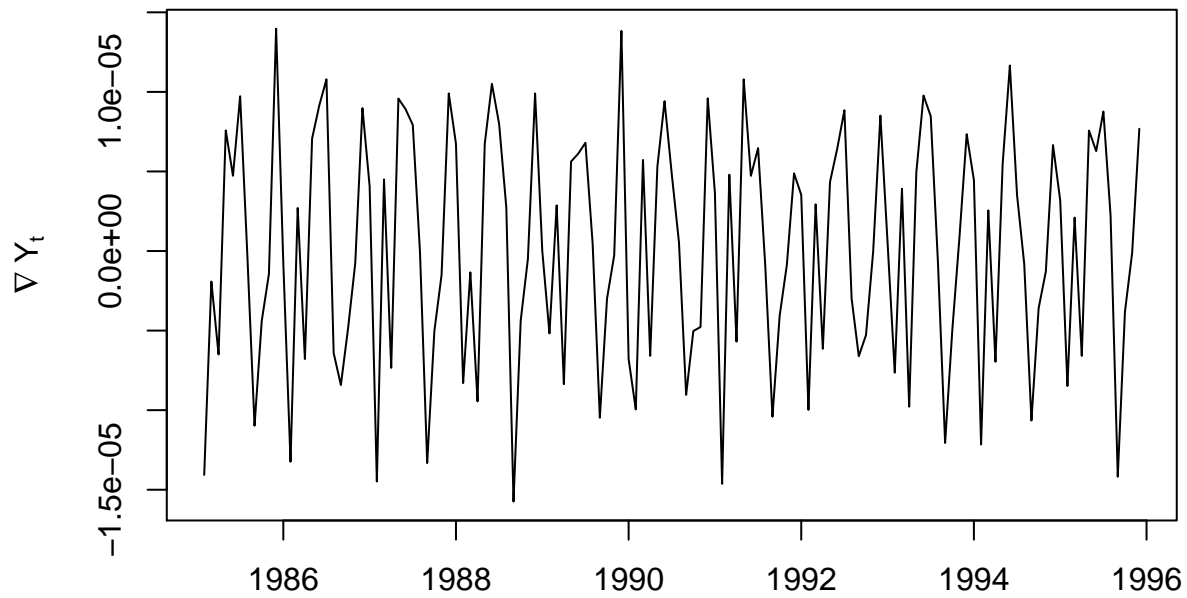
ACF/PACF of Box-Cox Transformed Time Series



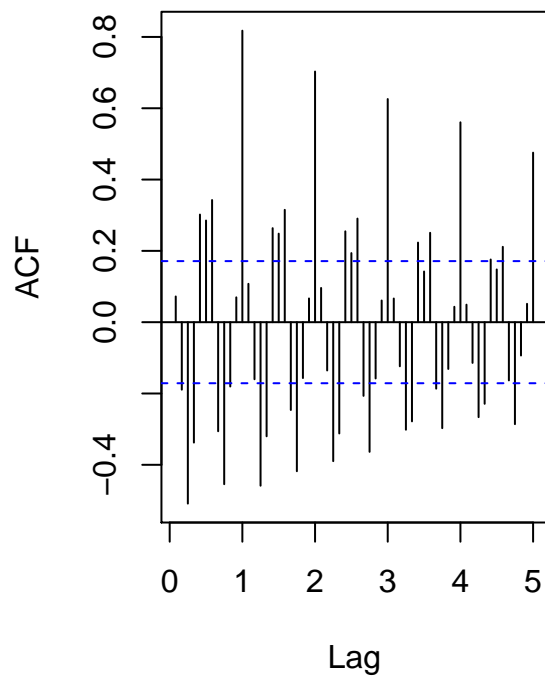
Differencing

Difference to Remove Trend

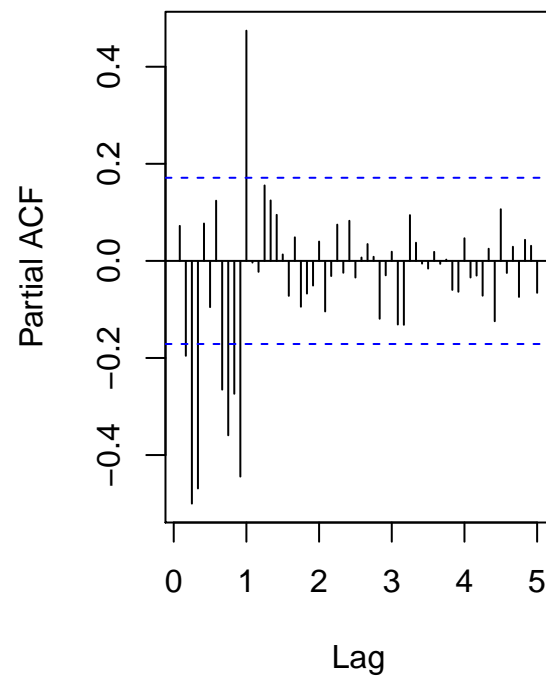
De-trended: Electricity Generation Differenced Lag 1



ACF of TS Differenced Lag 1



PACF of TS Differenced Lag 1

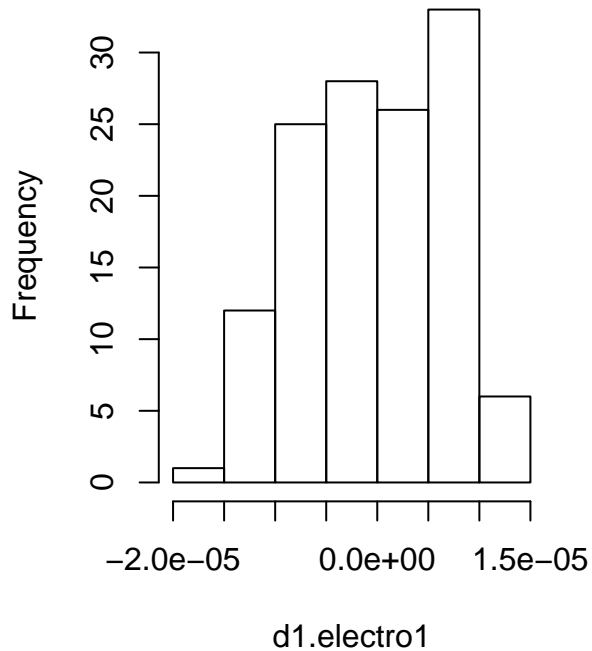


```
## [1] "The variance of the data differenced once at lag 1 is: 5.39172726366207e-11"
```

```
## [1] "The variance of the data differenced twice at lag 1 is: 9.88032626706133e-11"
```

```
## [1] "The variance of the de-trended time series is: 5.39172726366207e-11"
```

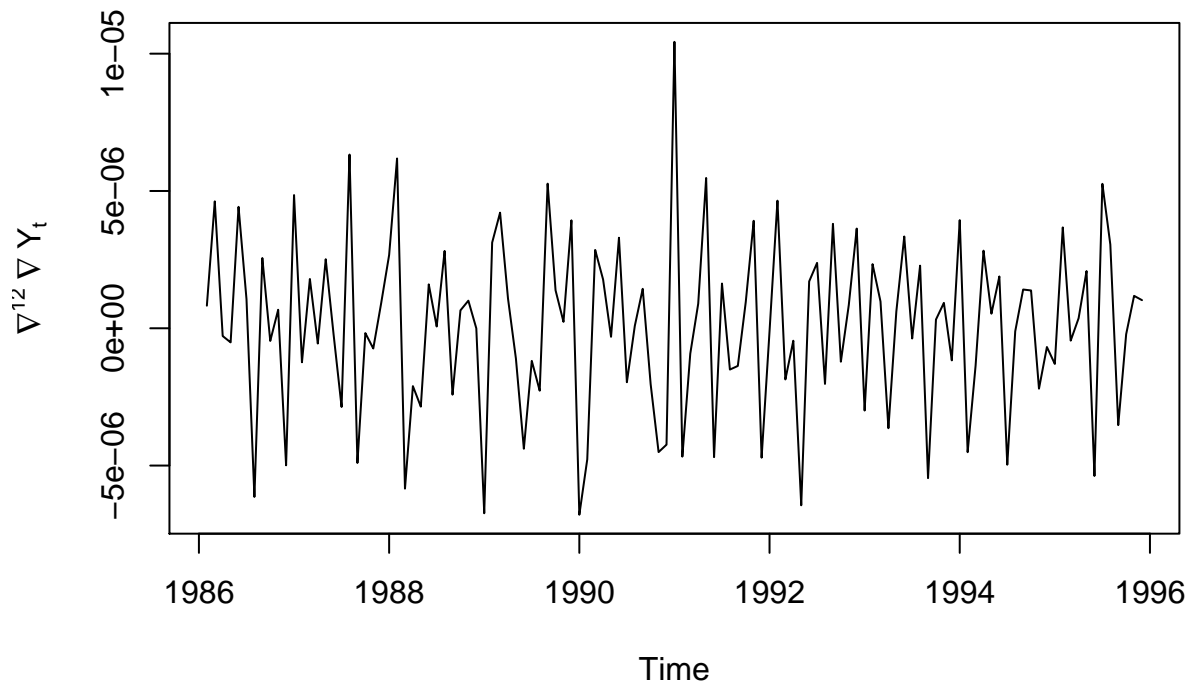
Histogram: Differenced Lag 1



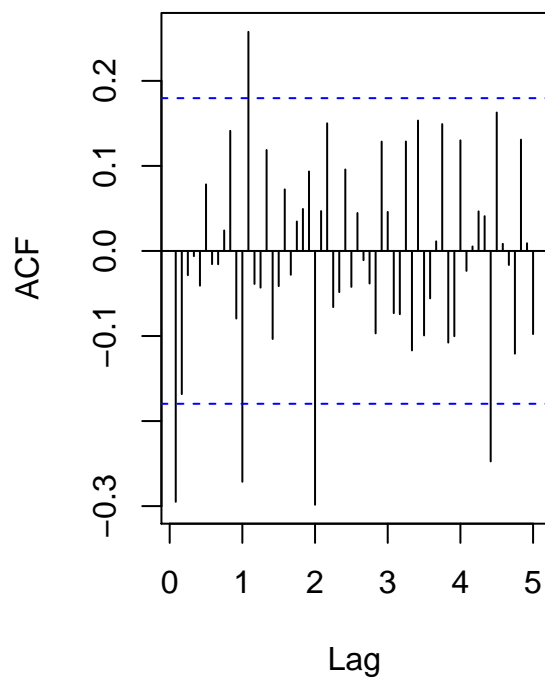
In this section, the time series is differenced at lag 1 in order to remove trend. A second differencing at lag 1 increases the variance, and therefore is considered overdifferencing. Since the variance is smaller when differencing only once at lag 1, we can conclude the time series has a significant linear trend over quadratic trend. With this single differencing at lag 1 of the time series, the linear trend has been significantly removed. The variance is also smaller than that of the prior Box-Cox Transformed data.

Difference to Remove Seasonality

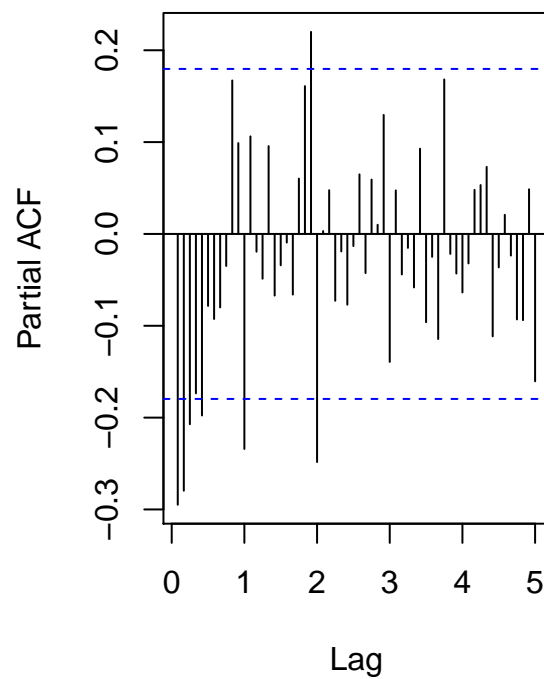
De-seasonalized: Electricity Generation Differenced x2 Lag 12



ACF Differenced x2 Lag 12



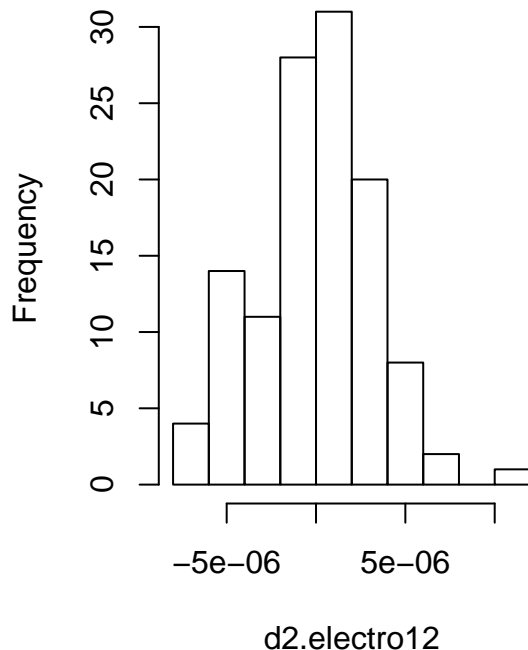
PACF Differenced x2 Lag 12



```
## [1] "The variance of the de-trended de-seasonalized time series is: 1.03172048629293e-11"
## Augmented Dickey-Fuller Test
## alternative: stationary
```

```
##
## Type 1: no drift no trend
##      lag      ADF p.value
## [1,]  0 -14.65    0.01
## [2,]  1 -11.66    0.01
## [3,]  2  -9.77    0.01
## [4,]  3  -8.59    0.01
## [5,]  4  -8.59    0.01
## Type 2: with drift no trend
##      lag      ADF p.value
## [1,]  0 -14.60    0.01
## [2,]  1 -11.62    0.01
## [3,]  2  -9.74    0.01
## [4,]  3  -8.56    0.01
## [5,]  4  -8.55    0.01
## Type 3: with drift and trend
##      lag      ADF p.value
## [1,]  0 -14.54    0.01
## [2,]  1 -11.56    0.01
## [3,]  2  -9.69    0.01
## [4,]  3  -8.52    0.01
## [5,]  4  -8.51    0.01
## ----
## Note: in fact, p.value = 0.01 means p.value <= 0.01
```

Histogram: Differenced x2 Lag 1

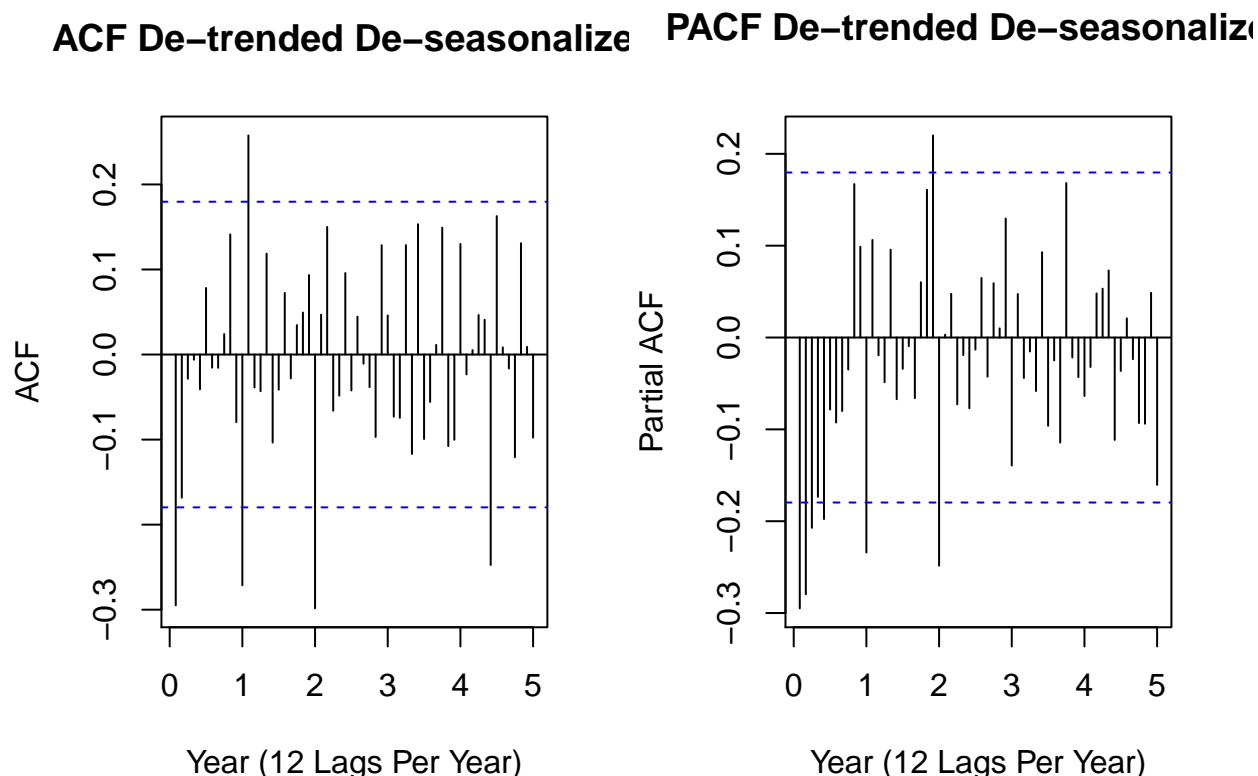


In this section, the time series is differenced for the second time, this time at lag 12, in order to remove the seasonal component. With this second differencing at lag 12, the seasonal component has been significantly removed and the variance has been reduced to $1.031720486e-11$. Since the Augmented Dickey-Fuller Test has passed, we can assume the series is now stationary.

Model Building

Model Interpretation Through Analysis of ACF/PACF

Below, we have the ACF and PACF for the de-trended and de-seasonalized time series from the last section.



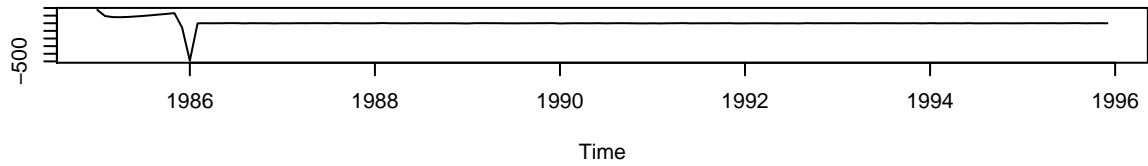
In the above section, we differenced the Electricity Generated time series once at lag 1 to remove the trend and once at lag 12 to remove seasonality. With this respectively, it can be concluded that $d = 1$ (non-seasonal differencing component) and $D = 1$ (seasonal differencing component). Since the Augmented Dickey-Fuller Test passed and the ACF values are mostly within the confidence interval bands, we are able to conclude that both the linear trend and seasonal component have been removed; therefore, the data has been made stationary. Observing the ACF at lags which are multiples of the seasonal period $s = 12$, we notice strong correlations (spikes) at lags 12 and 24 (labeled years 1 and 2). We can conclude from this observation that the time series for our observed model (Model 1) has a seasonal MA component with $Q = 2$. The same pattern occurs within the PACF as well, allowing us to conclude that the data also has a seasonal AR component with $P = 2$. In the ACF, there is also a significant correlation at lag 1, after which the lags cut off; from this, we can assume that there is a non-seasonal MA component with $q = 1$. The PACF on the other hand, shows very significant correlations in both lags 1 and 2, and a somewhat significant correlation at lag 3 before the cut off, suggesting a non-seasonal AR component with $p = 3$. The assumptions made from the ACF and PACF of our transformed time series suggest the Seasonal ARIMA model: $SARIMA(3,1,1) \times (2,1,2)_{12}$.

Parameter Estimation

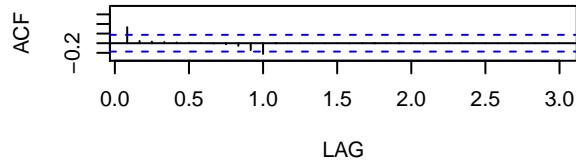
```
## [1] "Model 1 (Observed From ACF/PACF) Summary & Estimated Paramaters:"
```

Model: (3,1,1) (2,1,2) [12]

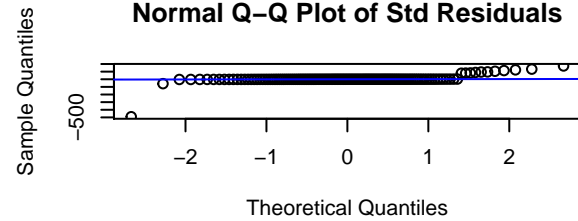
Standardized Residuals



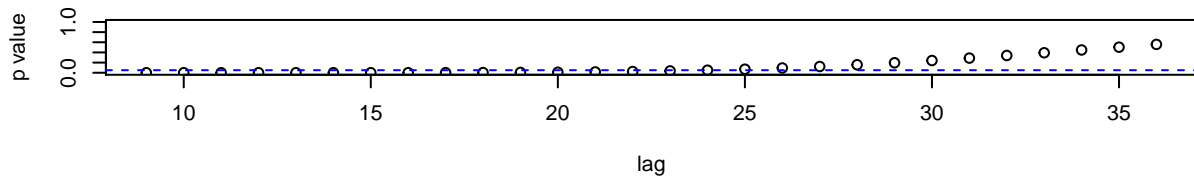
ACF of Residuals



Normal Q-Q Plot of Std Residuals



p values for Ljung-Box statistic



```
## $fit
##
## Call:
## stats::arima(x = xdata, order = c(p, d, q), seasonal = list(order = c(P, D,
##     Q), period = S), xreg = xreg, optim.control = list(trace = trc, REPORT = 1,
##     reltol = tol))
##
## Coefficients:
##          ar1      ar2      ar3      ma1      sar1      sar2      sma1      sma2
##      0.1771 -0.1540 -0.1382 -0.7430  0.0247 -0.3035 -0.6980 -0.0172
## s.e.  0.1229  0.0987  0.1076  0.0904  0.2779  0.1279  0.2961  0.2676
##      xreg
##     -0.0453
## s.e.   0.0200
##
## sigma^2 estimated as 4.302e-12:  log likelihood = 1380.92,  aic = -2741.83
##
## $degrees_of_freedom
## [1] 123
##
## $ttable
##      Estimate      SE t.value p.value
## ar1    0.1771 0.1229  1.4405  0.1523
## ar2   -0.1540 0.0987 -1.5605  0.1212
## ar3   -0.1382 0.1076 -1.2848  0.2013
## ma1   -0.7430 0.0904 -8.2150  0.0000
## sar1    0.0247 0.2779  0.0891  0.9292
## sar2   -0.3035 0.1279 -2.3726  0.0192
```

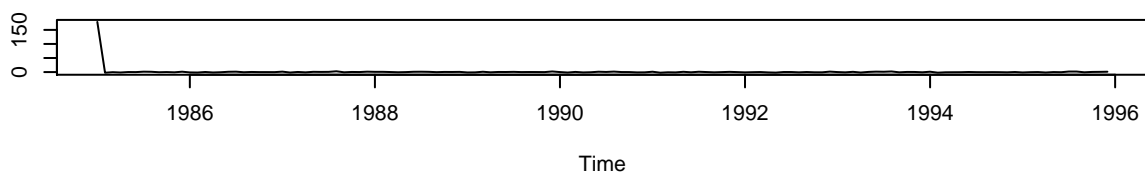
```

## sma1 -0.6980 0.2961 -2.3573 0.0200
## sma2 -0.0172 0.2676 -0.0643 0.9488
## xreg -0.0453 0.0200 -2.2635 0.0254
##
## $AIC
## [1] -25.03554
##
## $AICc
## [1] -25.00662
##
## $BIC
## [1] -25.83899
##
## Call:
## stats::arima(x = xdata, order = c(p, d, q), seasonal = list(order = c(P, D,
##      Q), period = S), xreg = xreg, optim.control = list(trace = trc, REPORT = 1,
##      reltol = tol))
##
## Coefficients:
##          ar1          ar2          ar3          ma1          sar1          sar2          sma1          sma2
##          0.1771 -0.1540 -0.1382 -0.7430 0.0247 -0.3035 -0.6980 -0.0172
## s.e.      0.1229 0.0987 0.1076 0.0904 0.2779 0.1279 0.2961 0.2676
##          xreg
##          -0.0453
## s.e.      0.0200
##
## sigma^2 estimated as 4.302e-12: log likelihood = 1380.92, aic = -2741.83
## [1] "Model 2 (computer-generated model) Summary & Estimated Parameters:"
##
## Series: electro.bxcx
## ARIMA(3,1,1)(0,0,2)[12] with drift
##
## Coefficients:
##          ar1          ar2          ar3          ma1          sma1          sma2          drift
##          0.4530 -0.1881 -0.2571 -0.9545 0.9643 0.4900          0
## s.e.      0.0884 0.0932 0.0851 0.0337 0.0967 0.0716      NaN
##
## sigma^2 estimated as 2.746e-09: log likelihood=1461.35
## AIC=-2906.7 AICc=-2905.52 BIC=-2883.69
## [1] "Re-fit Model 2 Including Xreg:"

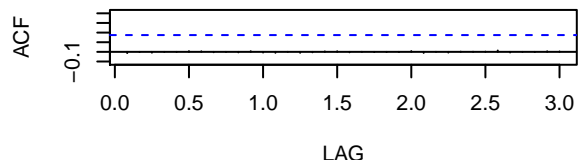
```

Model: (3,1,1) (0,0,2) [12]

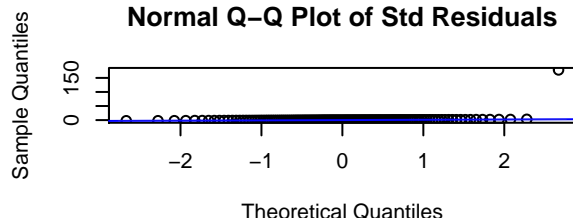
Standardized Residuals



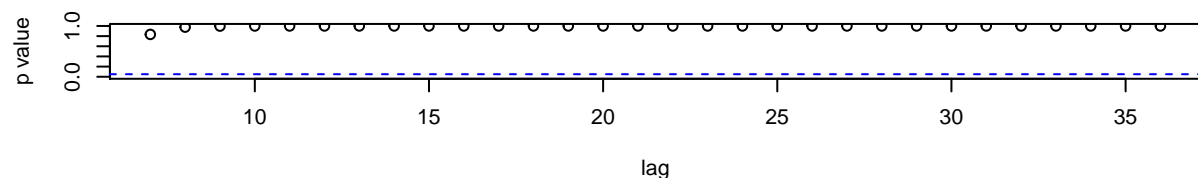
ACF of Residuals



Normal Q-Q Plot of Std Residuals



p values for Ljung-Box statistic



```
## $fit
##
## Call:
## stats::arima(x = xdata, order = c(p, d, q), seasonal = list(order = c(P, D,
##     Q), period = S), xreg = xreg, optim.control = list(trace = trc, REPORT = 1,
##     reltol = tol))
##
## Coefficients:
##          ar1      ar2      ar3      ma1      sma1      sma2      xreg
##          0.4530 -0.1881 -0.2571 -0.9545  0.9643  0.4900      0
## s.e.  0.0884  0.0932  0.0851  0.0337  0.0967  0.0716   NaN
##
## sigma^2 estimated as 1.067e-11:  log likelihood = 1461.35,  aic = -2906.7
##
## $degrees_of_freedom
## [1] 125
##
## $ttable
##      Estimate      SE  t.value p.value
## ar1    0.4530 0.0884   5.1236 0.0000
## ar2   -0.1881 0.0932  -2.0175 0.0458
## ar3   -0.2571 0.0851  -3.0202 0.0031
## ma1   -0.9545 0.0337 -28.2951 0.0000
## sma1    0.9643 0.0967   9.9719 0.0000
## sma2    0.4900 0.0716   6.8444 0.0000
## xreg    0.0000   NaN      NaN    NaN
##
## $AIC
```

```
## [1] -24.15718
##
## $AICc
## [1] -24.13316
##
## $BIC
## [1] -25.00431
##
## Call:
## stats::arima(x = xdata, order = c(p, d, q), seasonal = list(order = c(P, D,
##      Q), period = S), xreg = xreg, optim.control = list(trace = trc, REPORT = 1,
##      reltol = tol))
##
## Coefficients:
##          ar1          ar2          ar3          ma1          sma1          sma2    xreg
##          0.4530    -0.1881    -0.2571    -0.9545    0.9643    0.4900         0
## s.e.    0.0884    0.0932    0.0851    0.0337    0.0967    0.0716      NaN
##
## sigma^2 estimated as 1.067e-11:  log likelihood = 1461.35,  aic = -2906.7
```

Using, the `auto.arima()` function, the following computer-generated seasonal ARIMA model was synthesized: SARIMA(3,1,1)x(0,0,2)₁₂ with drift (0). All of the non-seasonal parameters (p, d, q) are the same as those for Model 1, which was acquired by analysis of the ACF and PACF. Model 2 possesses the same non-seasonal AR component with p = 3, non-seasonal differencing component d = 1, and non-seasonal MA component with q = 1. The seasonal MA component Q = 2 is also identical to that in the model deduced by analysis of the ACF/PACF above, however the other order values (P, D) are different. While the plots above suggested a seasonal AR component with P = 2 and a seasonal differencing component with D = 1, the `auto.arima()` function chose P = 0 and D = 0.

We assess the significance of the above estimated coefficients (parameters) by testing whether they are greater than about $2x(\sigma)$ for that particular coefficient. For Model 1, all coefficients are insignificant, except for those coefficients θ_1 , Φ_2 , and Θ_1 . These are respectively, the coefficients for the `ma1`, `sar2`, and `sma1`. For Model 2 on the other hand, it can be stated that all coefficients are significant. This suggests that the quality of our estimations is much more accurate, and probably safer to depend on for Model 2.

Model Selection

```
## [1] "The AICc for the analytically observed model (using ACF/PACF) is: -2740.35565277998"
## [1] "The AICc for the computer-generated model (using auto-arima) is: -2905.79281992578"
```

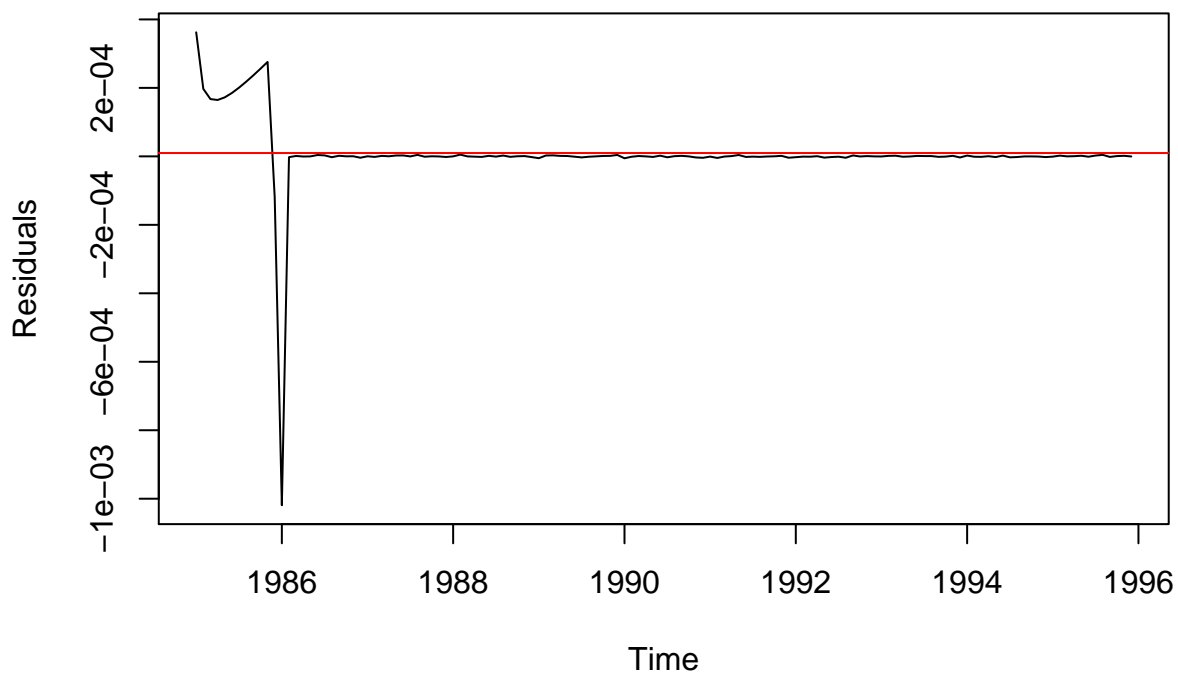
Based on the above AICc values, we can conclude that the best performing model out of the two is the computer-generated model chosen by the `auto.arima()` function by a small difference in AICc. This principle of parsimony also supports this model, as it has fewer parameters.

Diagnostic Checking

We will now perform diagnostic tests in order to check if the discovered models are suitable for forecasting. Diagnostic checks will include analysis of residuals, Portmanteau and normality tests, and checking for causality and invertibility.

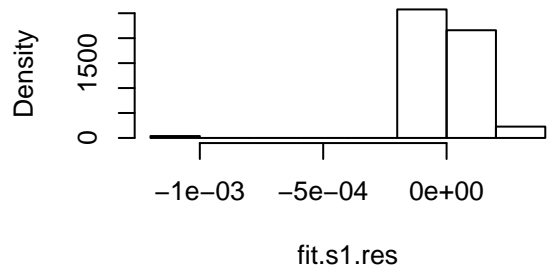
Model 1 (Model Found Through Analysis of ACF/PACF):

Residuals for Observed Model 1

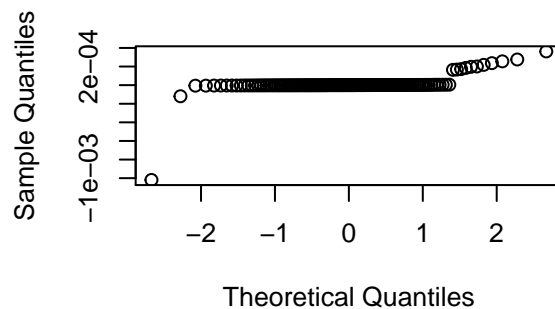


```
## [1] "The sample variance of residuals is: 1.23288563787723e-08"
```

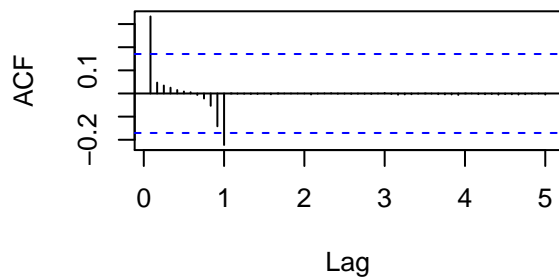
Histogram of Residuals: Model 1



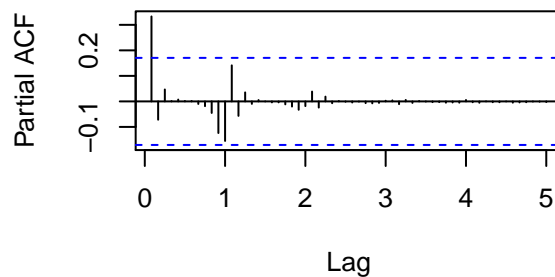
Normal Q-Q Plot: Model 1



ACF of Residuals: Model 1



ACF of Residuals: Model 1



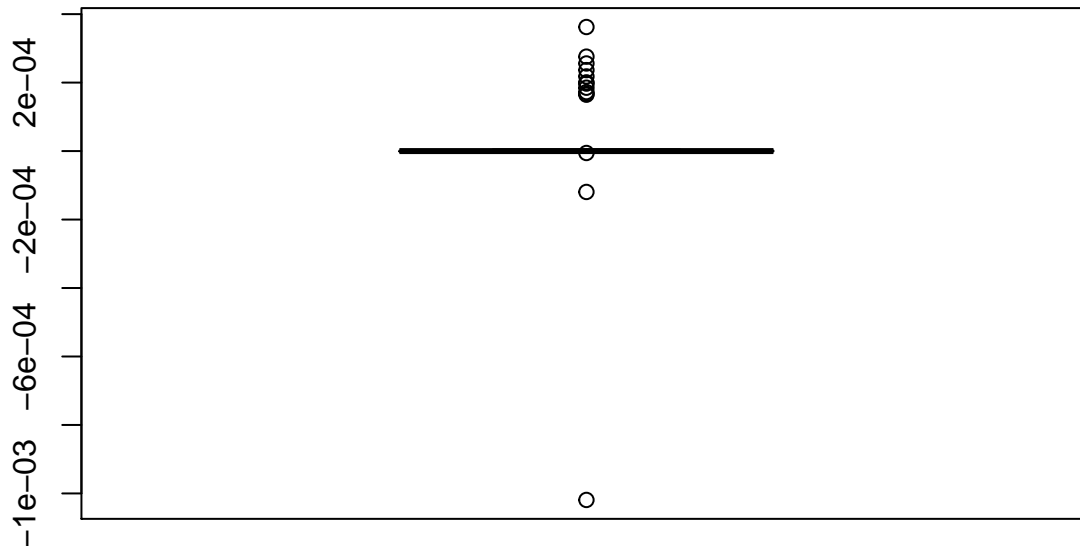
There is at least one outlier in the residuals which is making the plots difficult to assess. Nearly all of the

ACF and PACF values are within the confidence intervals, thus can be counted as zeros. However, there are strange patterns in both the ACF and PACF. The sample mean (red line) in the plot of the residuals is close to 0, supporting the assumption that the residuals resemble white noise. The sample variance is not close to 1.

Examine Outliers

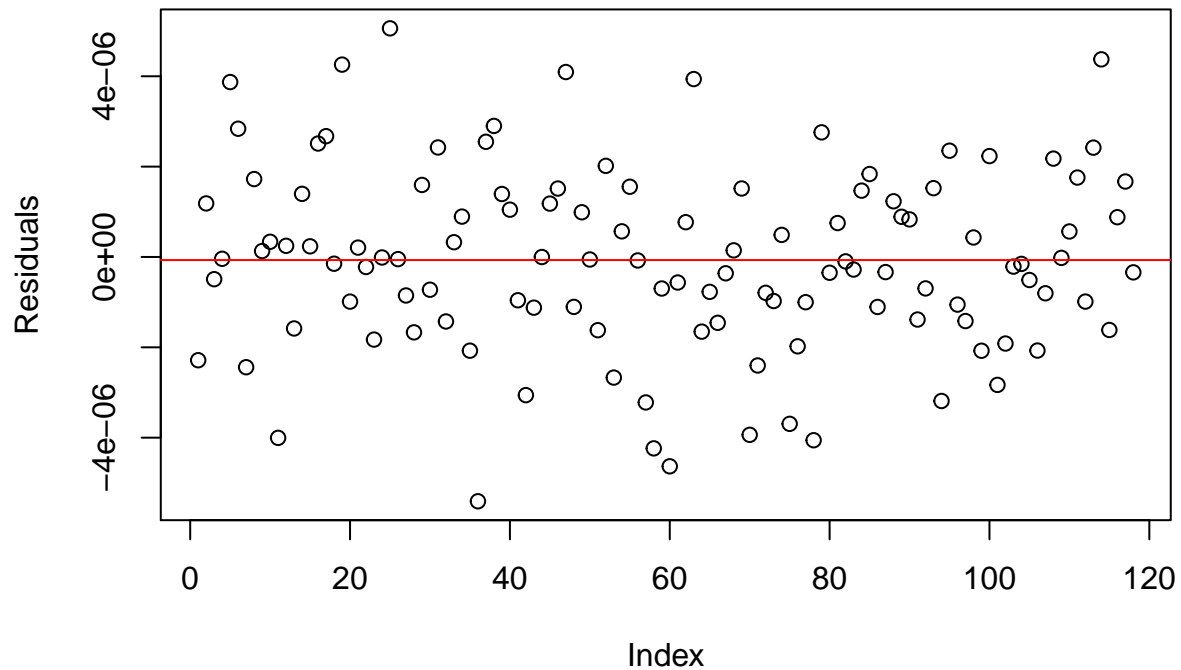
We will temporary remove the outliers from the residuals in order to better analyze their relationship in the plots.

Outlier Examination: Model 1



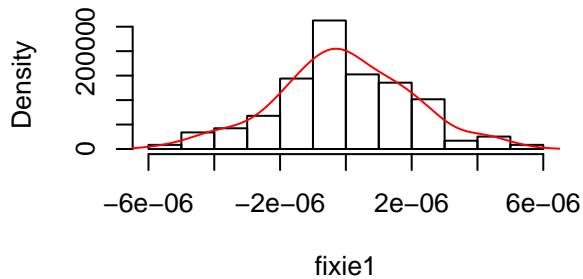
```
## [1] "Boxplot Statistics: "
## $stats
## [1] -5.407660e-06 -1.253446e-06 -3.039629e-08 1.632450e-06 5.060675e-06
##
## $n
## [1] 132
##
## $conf
## [1] -4.272686e-07 3.664760e-07
##
## $out
## [1] 3.623398e-04 1.971141e-04 1.673172e-04 1.646475e-04 1.723168e-04
## [6] 1.850416e-04 2.006322e-04 2.179903e-04 2.365195e-04 2.558804e-04
## [11] 2.758429e-04 -1.193779e-04 -1.018989e-03 -5.634234e-06
```

Adjusted Residuals for Observed Model 1 (Outliers Removed)

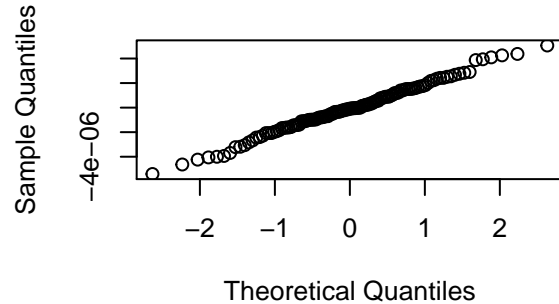


[1] "The sample variance of residuals without outliers is: 4.0996620287897e-12"

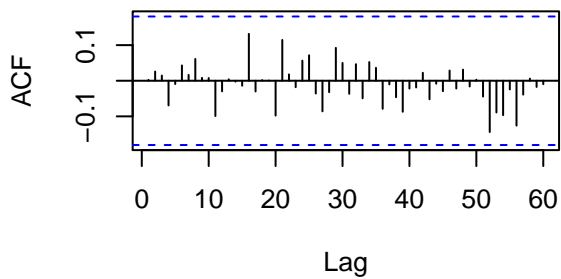
Histogram of Residuals (Outlier Adj): Model



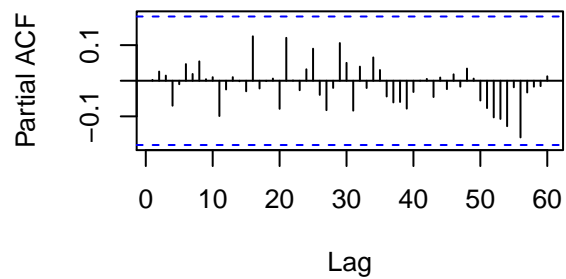
Normal Q-Q Plot (Outlier Adj): Model



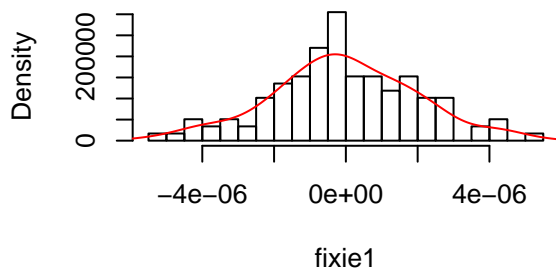
ACF of Residuals (Outlier Adj): Model



ACF of Residuals (Outlier Adj): Model



Histogram of Residuals (Outlier Adj): Mox



With outliers (2) removed it can be seen that there are no strong concrete visible patterns or relationships in the residual plot, however it does seem as if there may be some underlying trend and/or seasonality. The residuals appear for the most part to be random and uncorrelated, however they do seem like they may have a slight negative linear trend. Observing the plot, it can be seen that there may be a slight trend and minor change in variance in some places. The distribution in the histogram appears to be somewhat normal-like, however when we look at a histogram with more breaks, it begins to look less normal. The Normal Q-Q Plot looks normal for the most part with its central values, but has some somewhat significant tail activity which veers from normal. All values in the ACF and PACF remain within the confidence intervals, thus can be counted as zeros; this supports a lack of significant correlation between the residuals. The sample mean (red line) in the residual plot is very close to 0, offering evidence for the residuals resembling white noise. The sample variance of the residuals with outliers removed is still not close to 1, however. Lets take a further look.

Tests (Portmanteau and Normality)

The Box.test function uses a lag value based on the lag autocorrelation coefficients:

$$\text{lag} = h = \sqrt{n} = \sqrt{132} = 11 \text{ (rounded)}$$

The fitdf value is the number of parameters ($p + q$) to be estimated or 0: $p + q = 4$

```
##
## Box-Pierce test
##
## data: fit.s1.res
## X-squared = 18.293, df = 7, p-value = 0.01072
##
## Box-Ljung test
##
## data: fit.s1.res
## X-squared = 18.977, df = 7, p-value = 0.008261
## [1] "Mcleod-Li Test:"
##
## Box-Ljung test
##
## data: (fit.s1.res)^2
## X-squared = 2.4816, df = 11, p-value = 0.996
##
## Shapiro-Wilk normality test
##
## data: fit.s1.res
## W = 0.30621, p-value < 2.2e-16
```

```
## [1] "Residuals Fitted to AR(0), White Noise:"
```

```
##
```

```
## Call:
```

```
## ar(x = fit.s1.res, aic = TRUE, order.max = NULL, method = c("yule-walker"))
```

```
##
```

```
## Coefficients:
```

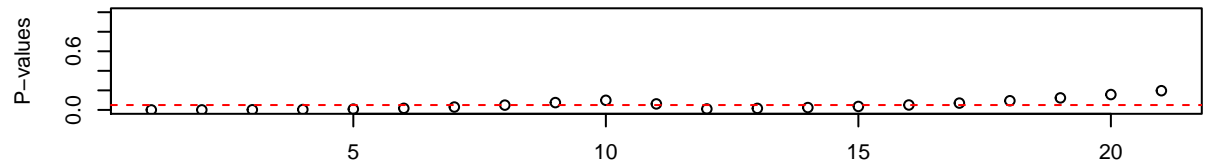
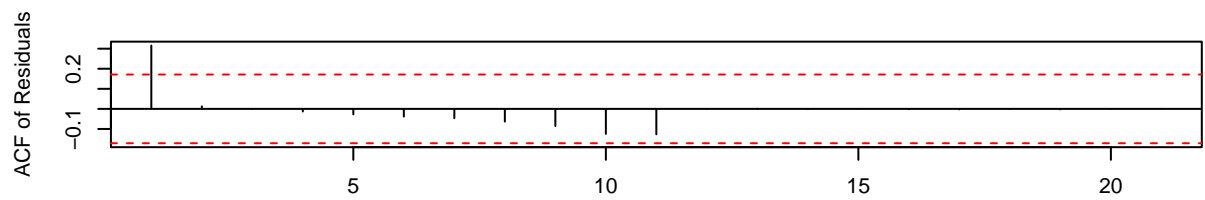
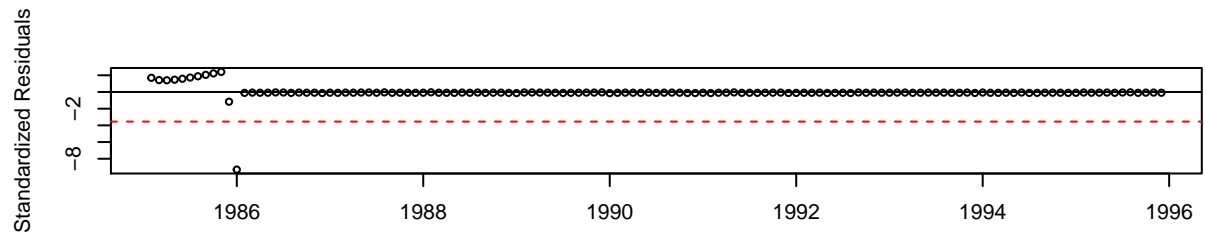
```
##      1
```

```
## 0.3329
```

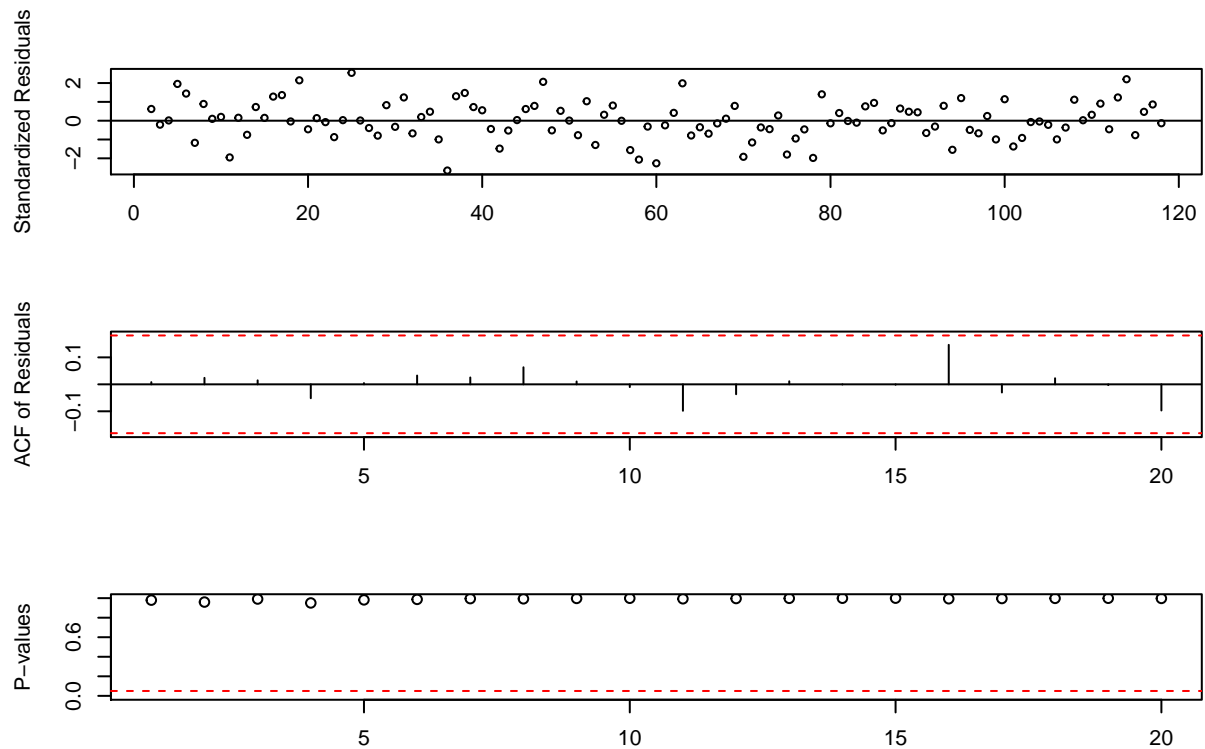
```
##
```

```
## Order selected 1  sigma^2 estimated as  1.105e-08
```

```
## [1] "Residual Plots: "
```



```
## [1] "Residual Plots (outliers removed): "
```



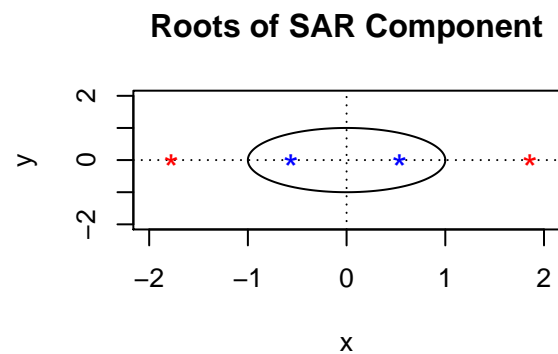
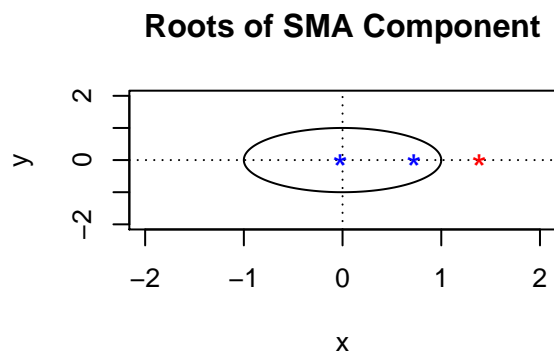
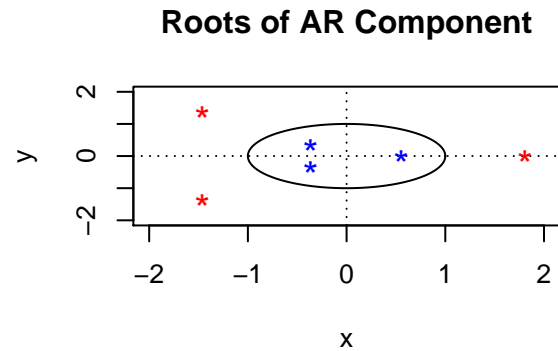
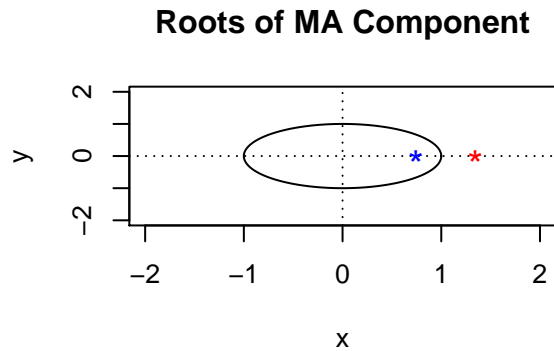
Not all Portmanteau tests passed with p-values larger than 0.05. The Box-Pierce and Box-Ljung tests failed, having p-values of less than 0.05. Due to the failure of the Box-Ljung and Box-Pierce tests, we reject the null hypothesis of independence in the residuals, therefore confirming that there are significant autocorrelations between the residuals. There is minimal support for these results in the plots of the ACF and PACF with outliers included, however most of the values still remain within the confidence intervals. Those values are treated as zeros, suggesting no strong correlations between those residuals. Since the Portmanteau tests have suggested significant autocorrelations in the residuals, we cannot confirm that the residuals resemble white noise, $WN(0, \sigma_z^2)$.

The Shapiro-Wilk normality test failed with a p-value less than 0.05, thus its hypothesis that the residuals are normally distributed is rejected. Although the Shapiro-Wilk test failed, it does not mean that the residuals do not exhibit any near-normal behavior. There is some normal-like behavior, however it is not surprising that this test failed considering the fairly drastic tail behavior. The Shapiro-Wilk test is very sensitive to slight departures in normality; the rejection of normality in this test could be due solely to tail behavior. This test is also particularly sensitive to small sample sizes, and could be rejected for that reason. The McLeod-Li test for heteroskedasticity passed, meaning we fail to reject the hypothesis that the series of residuals is IID, or in other words the hypothesis that the variability of the residuals is equal (homogenous) throughout the series. This test could potentially pass for the reason that pattern of the changing variance is seasonal and therefore cancels itself out. In the more zoomed out plot of the residuals without outliers, we see a DNA-like sinusoidal pattern which may be indicative of seasonality. There doesn't appear to be a significant trend, however there are cyclical changes in variance and it does seem as if they are not governed by randomness. There are also two significant outliers near 1986. This model has failed primary diagnostic testing and is not approved for forecasting.

Causality and Invertibility

```
## [1] "MA Roots:"
## [1] 1.345895+0i
## [1] "AR Roots:"
```

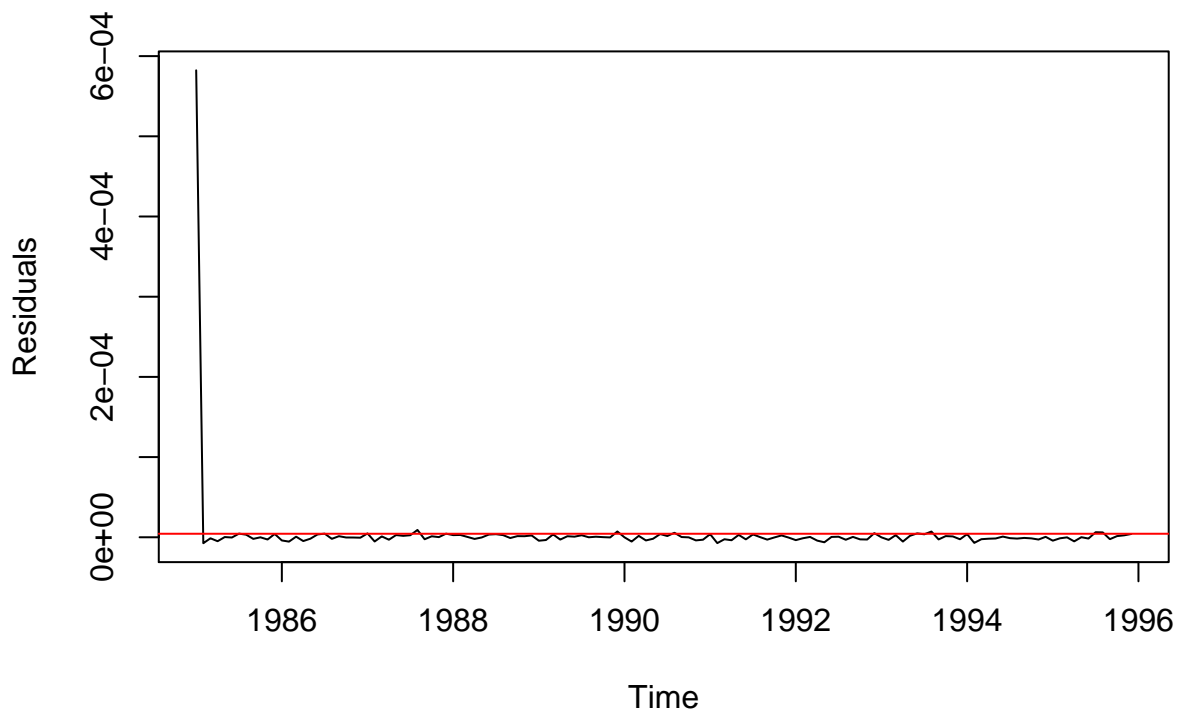
```
## [1] 1.808009-0.000000i -1.461168+1.366426i -1.461168-1.366426i
## [1] "SMA Roots:"
## [1] 1.385371+0i -41.966766-0i
## [1] "SAR Roots:"
## [1] 1.856332-0i -1.774948+0i
```



The MA component of Model 1 possesses one real root which falls outside of the unit circle (> 1), suggesting that this component is invertible; it is also causal since all $MA(q)$ is always causal. The AR component of the model possesses one real root and two imaginary roots that all fall outside the bounds of the unit circle. This suggests that the AR component is causal. The AR component is invertible since all $AR(p)$ is always invertible. The SMA component has two real roots which lie outside of the unit circle, and can be called invertible and causal. One of these roots cannot be seen (only its reflection is visible) because it is too large for the scale of the plot. The SAR component has two real roots which both lie outside of the unit circle. Altogether, we can assume invertibility and causality for the entire model.

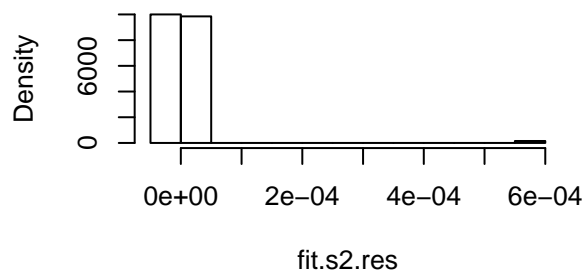
Preferred Model by AICc: Model 2 (Model Chosen By auto.arima() Function):

Residuals for auto.arima() Model 2

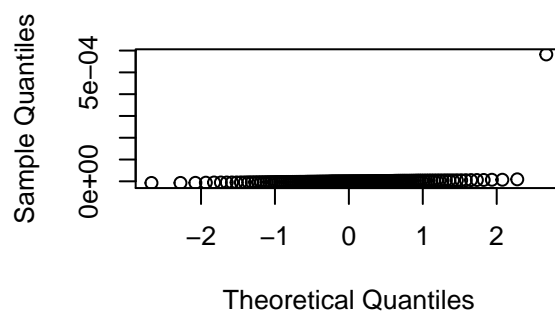


```
## [1] "The sample variance of residuals is: 2.57955800304829e-09"
```

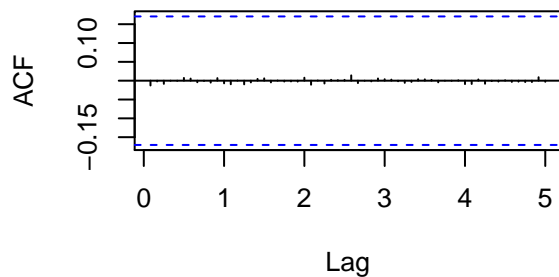
Histogram of Residuals: Model 2



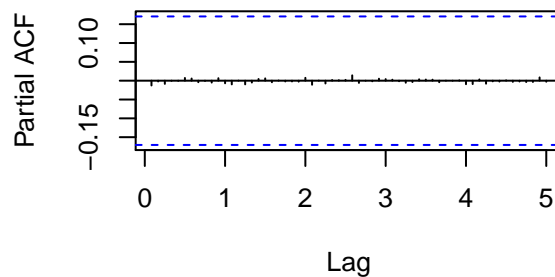
Normal Q-Q Plot: Model 2



ACF of Residuals: Model 2



ACF of Residuals: Model 2



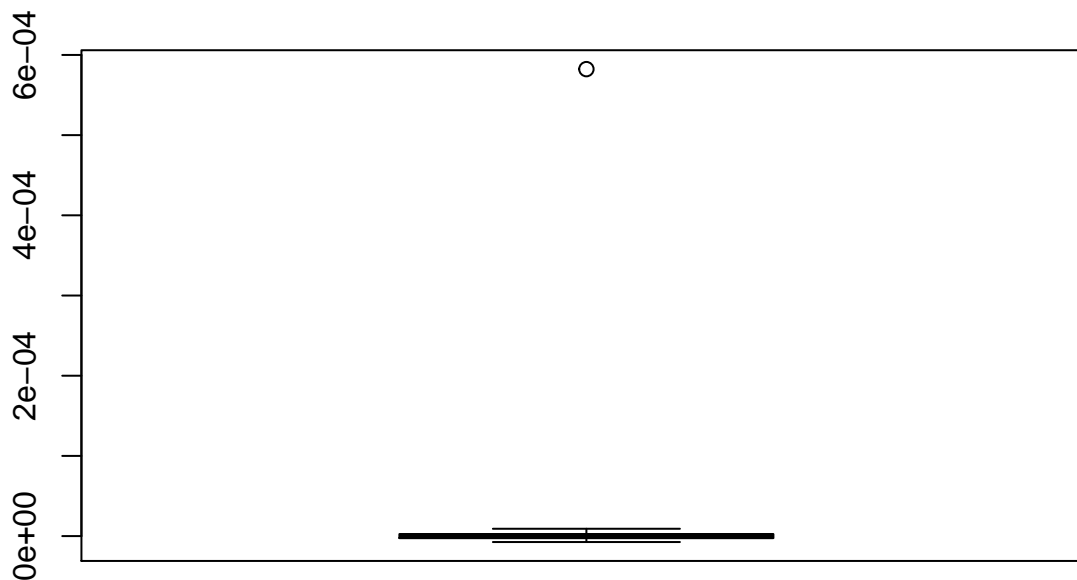
There is at least one outlier in the residuals which is making the plots difficult to assess. All of the ACF and

PACF values are within the confidence intervals, thus can be counted as zeros. The sample mean (red line) in the plot of the residuals is close to 0, supporting the assumption that the residuals resemble white noise. The sample variance is not close to 1.

Examine Outliers

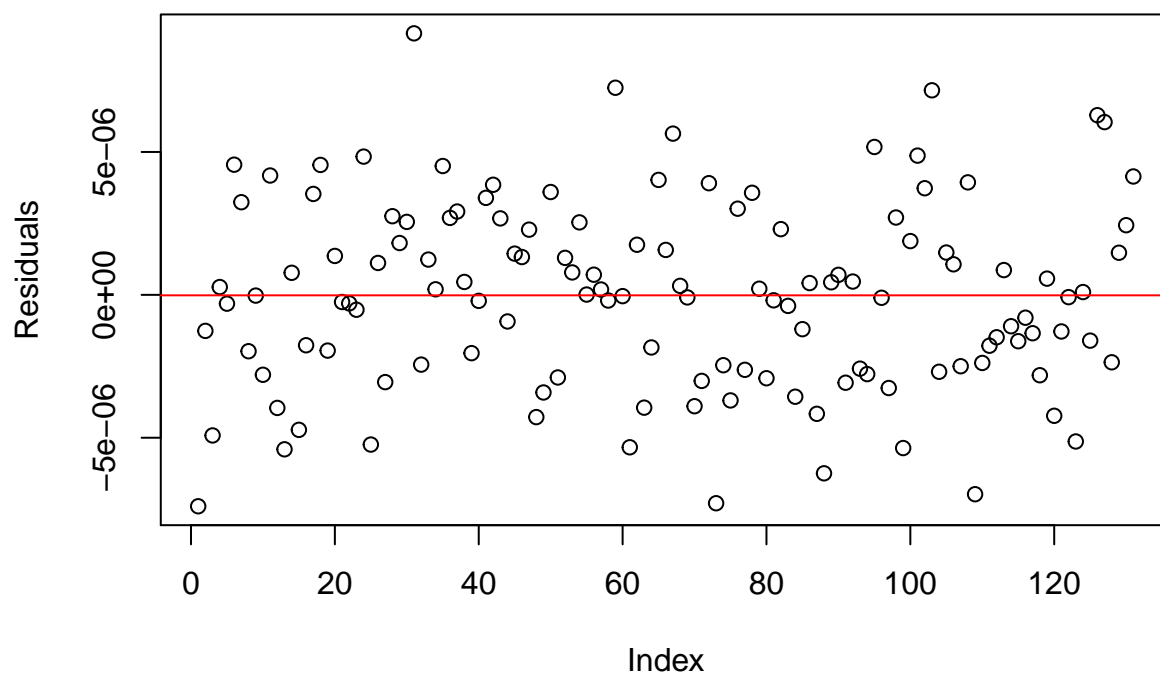
We will temporary remove the outliers from the residuals in order to better analyze their relationship in the plots.

Outlier Examination: Model 2



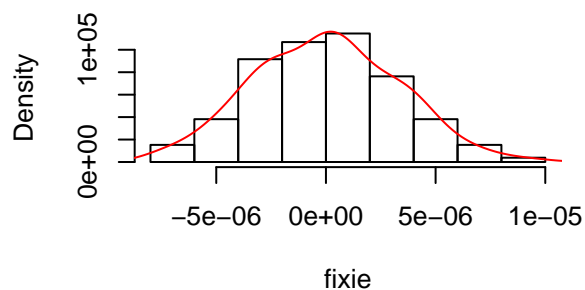
```
## [1] "Boxplot Statistics: "  
## $stats  
## [1] -7.397540e-06 -2.480079e-06 -7.995096e-09 2.486300e-06 9.149759e-06  
##  
## $n  
## [1] 132  
##  
## $conf  
## [1] -6.909783e-07 6.749881e-07  
##  
## $out  
## [1] 0.0005822997
```


Adjusted Residuals for auto.arima() Model 2 (Outliers Removed)

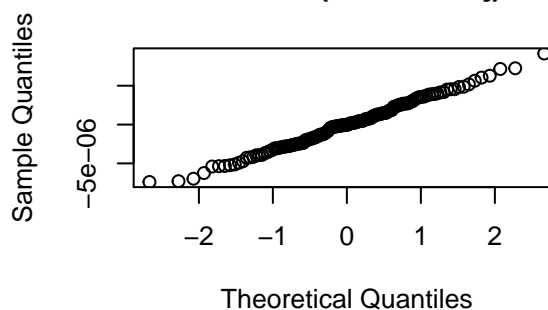


```
## [1] "The sample variance of residuals without outliers is: 1.07554768102652e-11"
```

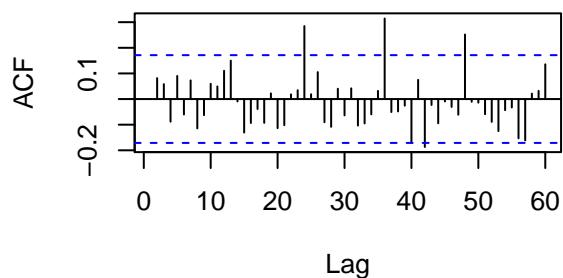
Histogram of Residuals (Outlier Adj): Model



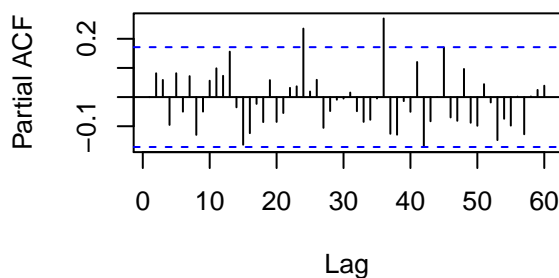
Normal Q-Q Plot (Outlier Adj): Model



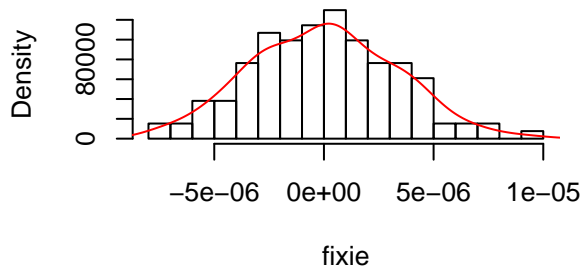
ACF of Residuals (Outlier Adj): Model



ACF of Residuals (Outlier Adj): Model



Histogram of Residuals (Outlier Adj): Model 2



With outliers (1) removed it can be seen that there are no visible patterns or relationships in the residual plot. The residuals appear to be random and uncorrelated, resembling white noise. Observing the plot, it can be seen that there is no trend, no seasonality, and no change in variance. The distribution in the histogram appears to be fairly close to normal, much better than that for Model 1. The Normal Q-Q Plot looks very strongly normal, but has some slight almost negligible tail activity. The Normal Q-Q Plot is more strongly normal than that for Model 1. Mostly all values in the ACF and PACF remain within the confidence intervals, thus can be counted as zeros; this supports the claim that there is no significant correlation between the residuals. The sample mean (red line) in the residual plot is very close to 0, offering evidence for the residuals resembling white noise. The sample variance of the residuals with outliers removed is still not close to 1, however.

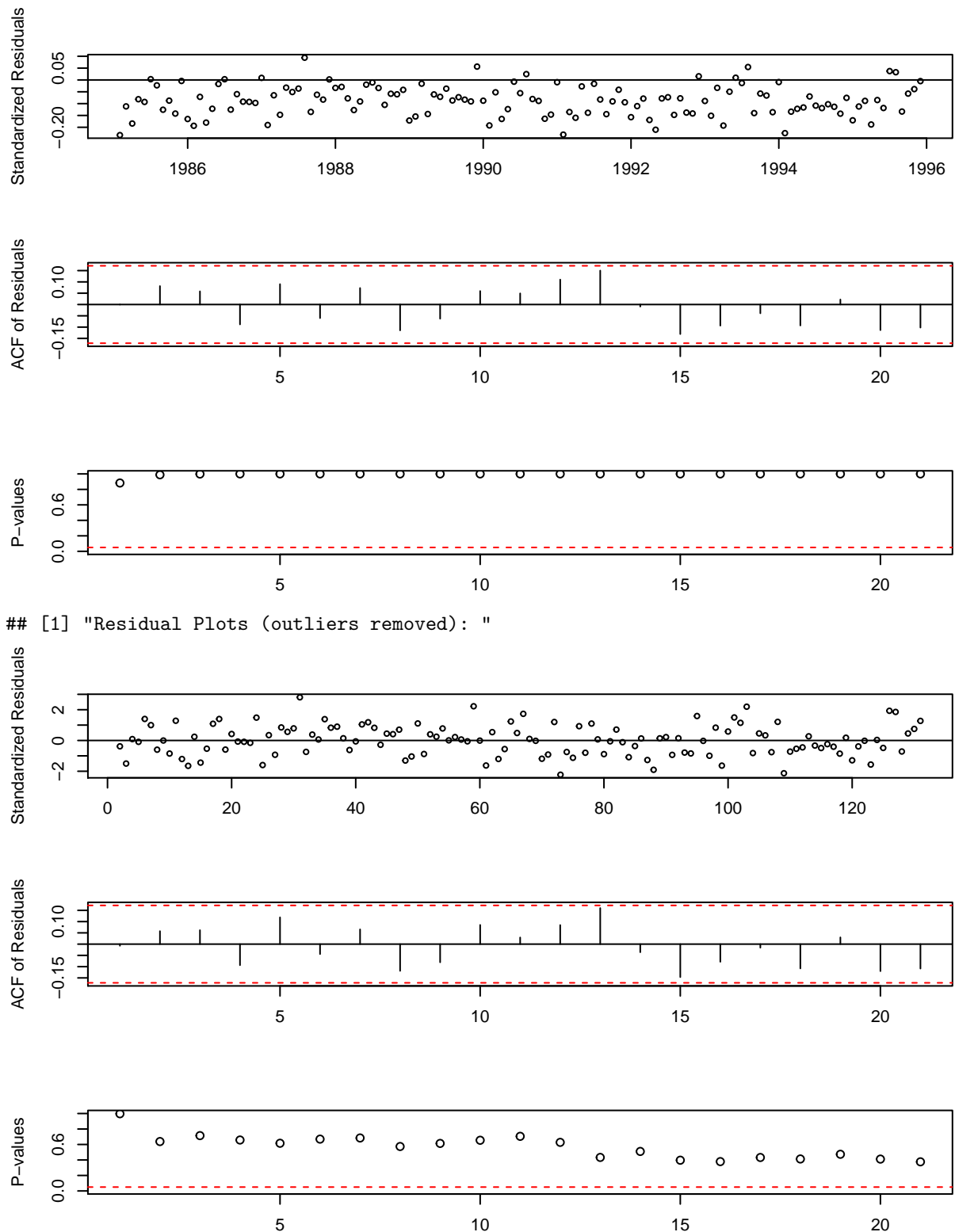
Tests

The Box.test function uses a lag value based on the lag autocorrelation coefficients:

$$\text{lag} = h = \sqrt{n} = \sqrt{132} = 11 \text{ (rounded)}$$

The fitdf value is the number of parameters ($p + q$) to be estimated or 0: $p + q = 4$

```
##
## Box-Pierce test
##
## data: fit.s2.res
## X-squared = 0.054289, df = 7, p-value = 1
##
## Box-Ljung test
##
## data: fit.s2.res
## X-squared = 0.057085, df = 7, p-value = 1
## [1] "Mcleod-Li Test:"
##
## Box-Ljung test
##
## data: (fit.s2.res)^2
## X-squared = 0.00024832, df = 11, p-value = 1
##
## Shapiro-Wilk normality test
##
## data: fit.s2.res
## W = 0.095018, p-value < 2.2e-16
## [1] "Residual Plots: "
```

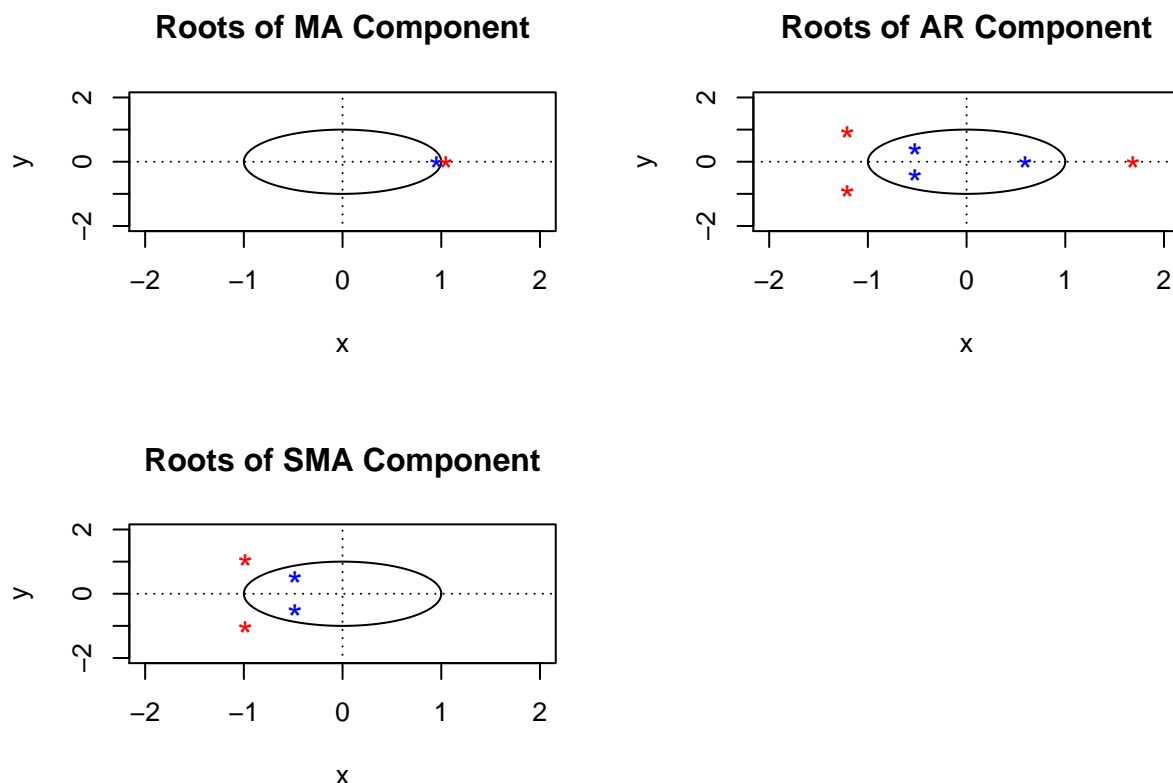


All Portmanteau tests passed with p-values larger than 0.05. Due to the passing of the Box-Ljung and Box-Pierce tests, we fail to reject the null hypothesis of independence in the residuals, therefore confirming that there is not significant autocorrelations in the residuals. This is supported by the plots of the ACF and PACF, in which nearly all values remain within the confidence intervals (treated as zeros), suggesting no strong correlations between residuals. Since there are no significant autocorrelations in the residuals, it is

confirmed that the residuals resemble white noise, $WN(0, \sigma_z^2)$. The residuals are fitted to $AR(0)$, i.e. WN. The Shapiro-Wilk normality test failed with a p-value less than 0.05, thus its hypothesis that the residuals are normally distributed is rejected. Although the Shapiro-Wilk test failed, it does not mean that the residuals do not exhibit normal behavior. The Shapiro-Wilk test is very sensitive to slight departures in normality; the rejection of normality in this test could be due solely to tail behavior. This test is also particularly sensitive to small sample sizes, and could be rejected for that reason. The McLeod-Li test for heteroskedasticity passed, meaning we fail to reject the hypothesis that the series of residuals is IID, or in other words the hypothesis that the variability of the residuals is equal (homogenous) throughout the series. In the plot of the residuals, there does not appear to be trend, seasonality or change in variance except for 1 outlier.

Causality and Invertibility

```
## [1] "MA Roots:"
## [1] 1.047669+0i
## [1] "AR Roots:"
## [1] 1.684686-0.000000i -1.208154+0.921479i -1.208154-0.921479i
## [1] "SMA Roots:"
## [1] -0.98398+1.035664i -0.98398-1.035664i
```



The MA component of the the model possesses one real root which falls outside of the unit circle (> 1), suggesting that this component is invertible; it is also causal since all $MA(q)$ is always causal. The AR component of the model possesses one real root and two imaginary roots that all fall outside the bounds of the unit circle. This suggests that the AR component is causal. The AR component is invertible since all $AR(p)$ is always invertible. The SMA component has two imaginary roots which lie outside of the unit circle, and can be called invertible and causal. Altogether, we can assume invertibility and causality for the entire model.

Final Model: Fitted Model

Both Model 1 and Model 2 do not pass all of the diagnostic tests, however Model 2 passes more tests. Model 1 failed all of the Portmanteau tests except for the McLeod-Li test, while Model 2 passed all of the Portmanteau tests. Both models failed the Shapiro-Wilk test, which is not a significant problem since they still both exhibit strong normal behavior in the visual tests. Model 1's failure of the Portmanteau tests suggests correlation between the residuals; this correlation may be due to a seasonal relationship. Model 2 has uncorrelated, approximately normal residuals, which lack trend, seasonality, or change in variance and resemble white noise. Both models are both causal and invertible.

With analysis, I conclude Model 2 to be the best model since it performs significantly better in every diagnostic test. Model 1's failure of the Portmanteau tests is a strong indication that it is not cut out for accurate forecasting. Model 2 also has a lower AICc and fewer parameters (parsimony). Lastly, Model 1 has only 3/8 coefficients which are significant, while Model 2 has all 6 coefficients significant. After parameter estimation and diagnostic checking, the final fitted model for the monthly time series of total electricity generated in the US from 1985 to 1996 is concluded to be:

$$SARIMA(3, 1, 1)(0, 0, 2)_{12} :$$

$$(1 - \phi_1 B - \phi_2 B^2 - \phi_3 B^3)(1 - B)X_t = (1 + \theta_1 B)(1 + \Theta_1 B^{12})(1 + \Theta_2 B^{12})Z_t$$

$$(1 - 0.453B + 0.1881B^2 + 0.2571B^3)(1 - B)X_t = (1 - 0.9545B)(1 + 0.9643B^{12})(1 + 0.49B^{12})Z_t$$

$$(1 - 0.453B + 0.1881B^2 + 0.2571B^3)\nabla X_t = (1 - 0.9545B)(1 + 0.9643B^{12} + 0.49B^{24})Z_t$$

where,

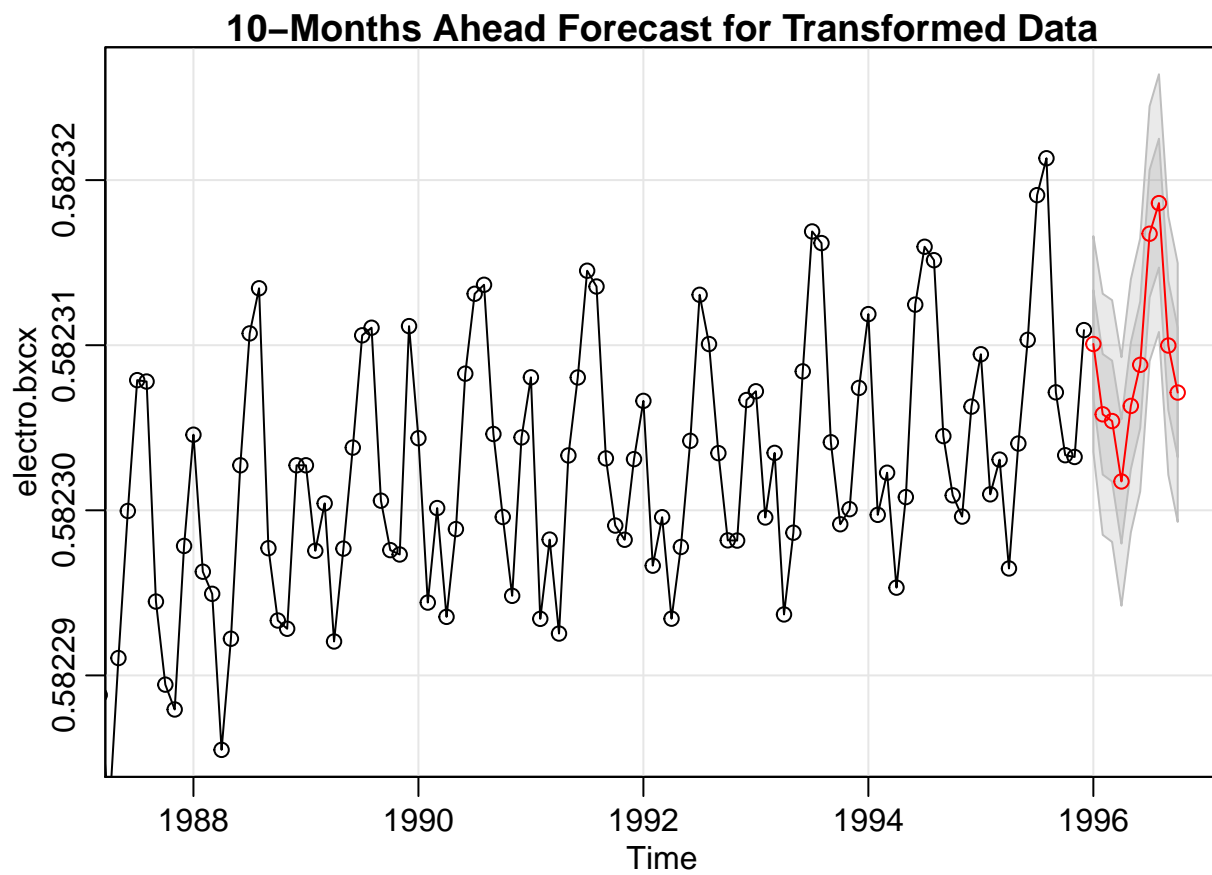
$$Z_t = WN(0, \sigma_z^2)$$

$$\sigma_z^2 = 0.000000003$$

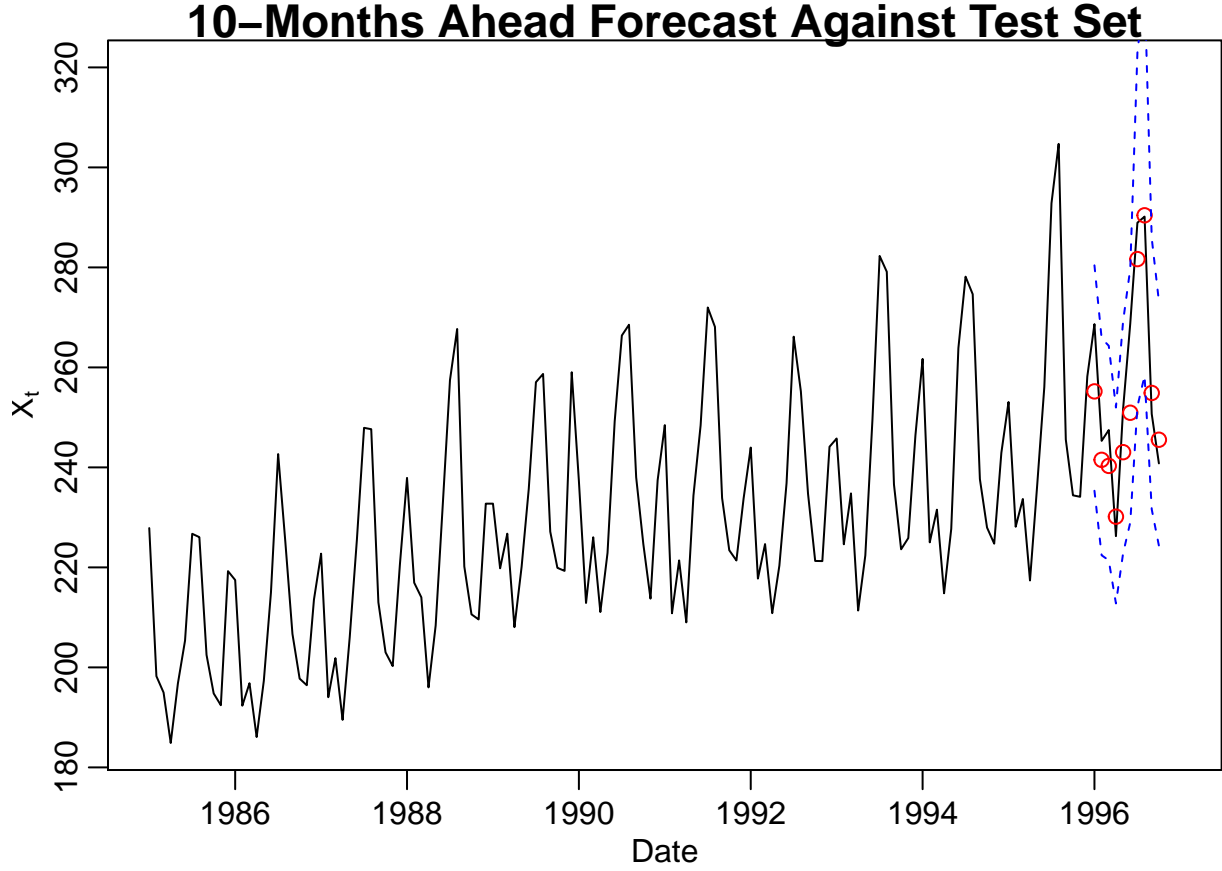
```
## [1] "AR(0) fit, White Noise Variance:"
##
## Call:
## ar(x = fit.s2.res, aic = TRUE, order.max = NULL, method = c("yule-walker"))
##
##
## Order selected 0  sigma^2 estimated as  2.58e-09
## [1] "Final Model Fit, White Noise Variance:"
## Series: electro.bxcx
## ARIMA(3,1,1)(0,0,2)[12] with drift
##
## Coefficients:
## Warning in sqrt(diag(x$var.coef)): NaNs produced
##          ar1      ar2      ar3      ma1      sma1      sma2      drift
##          0.4530 -0.1881 -0.2571 -0.9545  0.9643  0.4900         0
## s.e.  0.0884  0.0932  0.0851  0.0337  0.0967  0.0716      NaN
##
## sigma^2 estimated as 2.746e-09:  log likelihood=1461.35
## AIC=-2906.7  AICc=-2905.52  BIC=-2883.69
```

Forecasting

```
## [1] "10-Step Ahead Forecast On Transformed Data:"
```



```
## [1] "10-Step Ahead Forecast On Un-Transformed Data:"
```



Conclusions

With analysis of the ACF and PACF of the transformed time series, we found one model, $SARIMA(3,1,1) \times (2,1,2)_{12}$. Another model, $SARIMA(3,1,1) \times (0,0,2)_{12}$, was found using the `auto.arima()` computer-generation function. It was found that model chosen by `auto.arima()` was the most effective model in its performance in AICc and diagnostic checking, also possessing the minimum parameters. The final fitted model for the $SARIMA(3,1,1) \times (0,0,2)_{12}$ was written:

$$(1 - 0.0078B - 0.153315B^2 - 0.0219073B^3)\nabla X_t = (1 - 0.9545B)(1 + 1.4543B^{12} + 0.472507B^{24})Z_t$$

This model was found to accomplish the purpose of the study and provide a very accurate forecast for the 10-months of 1996 left out of the training data. Because of the accuracy of the prediction, we can assume the model is a very good one. We can also make further conclusions to assume that this forecast may suggest the rates of increase in US industry electricity generation are beginning to decline, although the generation of electricity is still increasing. With rapid growth in and growing access to technology, it would make sense that electricity generation would continue to rise to accommodate demand. However, it is interesting that the data was able to predict the lower values in seasonal trend in the year of 1996. Perhaps, rates of increase are declining because of new conservation techniques that also emerge with new technology. Another probable reason for a decline in rate of increase could be the bubble the technology world has come to in modern times. It is still going somewhere, but it seems to be taking its time or unsure of the next step.

Acknowledgments

Thank you to Nhan Hyunh for the advice on certain complicated decisions throughout the process!

References

- R Studio Software Version 1.0.136
- Utilities, Source: Makridakis, Wheelwright and Hyndman (1998), in file: data/elecnew, Description: The total generation of electricity by the U.S. electric industry (monthly data for the period Jan 1985 - Oct. 1996) For recent data, click here, <https://datamarket.com/data/set/22wj/the-total-generation-of-electricity-by-the-us-electric-industry-monthly-data-for-the-period-jan-1985-oct-1996-for-recent-!ds=22wj&display=line>
- Stack Overflow: <https://stackoverflow.com/questions/4787332/how-to-remove-outliers-from-a-dataset>
<http://www.unige.ch/ses/sococ/cl/r/tasks/outliers.e.html>
- Stack Exchange (Mcleod-Li Test): <https://stats.stackexchange.com/questions/174934/difference-between-ljung-box-and-r>
- Cran.r-Project (heteroskedasticity): <https://cran.r-project.org/web/packages/olsrr/vignettes/heteroskedasticity.html>
- Why Xreg Is Used: <http://www.stat.pitt.edu/stoffer/tsa2/Rissues.htm>

Appendix

```
# Libraries

library(data.table) # rename column
library(MASS)       # boxcox transformations
library(forecast)   # for auto.arima()
library(qpcR)
library(astsa)      # for sarima()
library(dse)        # plot roots
library(plotrix)    # draw.circle
library(TSA)
library(aTSA)

# Read in data
# Test Set
electric <- read.csv("generation-electric.csv")

# Use data.table library to rename column/variable
setnames(electric, "The.total.generation.of.electricity.by.the.U.S..electric.industry..monthly.data.for")

# Test Set: Create Time Series Object
electro_ts <- ts(electric[,2], start = c(1985, 1), frequency = 12)
# Plot Time Series for Test Set
electro_tsplot <- ts.plot(electro_ts, ylab = "Total Electricity Generated",
                          main = "Total Generation of Electricity Per Month in US (85-96), Test")
print.default(paste0("Summary of Test Set Time Series"))
summary(electro_ts)
```



```

# Column 2 Only
electricity <- electric[,2]

# Training Set (- 10 observations)
electric.tr <- electricity[1:132]

# Training Set: Create Time Series Object
electro.tr_ts <- ts(electric.tr, start = c(1985, 1), frequency = 12)
# Plot Time Series for Training Set
electro.tr_tsplot <- ts.plot(electro.tr_ts, ylab = "Total Electricity Generated",
                             main = "Total Generation of Electricity Per Month in US (85-96), Train")
print.default(paste0("Summary of Training Set Time Series"))
summary(electro.tr_ts)

hist(electro.tr_ts, main = "Histogram: Time Series of Training Data Set")

# Box-Cox Transformation
time <- 1:length(electro.tr_ts)
fit <- lm(electro.tr_ts ~ time)
boxcoxTrans <- boxcox(electro.tr_ts ~ time, plotit = TRUE)

lambda <- boxcoxTrans$x[which(boxcoxTrans$y == max(boxcoxTrans$y))]
print.default(paste0("The value of lambda is: ",
lambda))
electro.bbox <- (1/lambda)*(electro.tr_ts^lambda-1)

# Log Transformation
lgTrans <- log(electro.tr_ts)

# Square Root Transformation
squareTrans <- sqrt(electro.tr_ts)

# Plot Transformations vs. Original Data
op <- par(mfrow = c(1,4))
ts.plot(electro.tr_ts, main = "Original Data", ylab = expression(X[t]))
ts.plot(electro.bbox, main = "Box-Cox Transformed", ylab = expression(Y[t]))
ts.plot(lgTrans, main = "Log Transformed", ylab = expression(Y[t]))
ts.plot(squareTrans, main = "Square Root Transformed", ylab = expression(Y[t]))
par(op)

# Variances
# Original Data
print.default(paste0("The variance of the original training data is: ",
var(electro.tr_ts)))
# Box-Cox Transformed Data
print.default(paste0("The variance of the Box-Cox transformed data is: ",
var(electro.bbox)))
# Log Transformed Data
print.default(paste0("The variance of the log transformed data is: ",
var(lgTrans)))
# Square Root Transformed Data
print.default(paste0("The variance of the square-root transformed data is: ",

```

```

var(squareTrans)))

# Histogram of Box-Cox Algorithm (Lowest Variance Transformation)
hist(electro.bxcx, main = "Histogram: Box-Cox Transformed Data")

# Chosen Transformation
# Box-Cox Transformed Data vs. Original Data
op <- par(mfrow = c(1,2))
ts.plot(electro.tr_ts, main = "Original Data", ylab = expression(X[t]))
ts.plot(electro.bxcx, main = "Box-Cox Transformed", ylab = expression(Y[t]))
par(op)

# ACF/PACF for Box-Cox Transformation
par(mfrow = c(1,2))
acf(electro.bxcx, lag.max = 60, main = "ACF of Electricity Generation TS")
pacf(electro.bxcx, lag.max = 60, main = "PACF of Electricity Generation TS")

# Difference @ Lag 1 -> Remove Trend
d1.electro1 <- diff(electro.bxcx, lag = 1, differences = 1)

# Test for linear/quadratic trend
d2.electro1 <- diff(d1.electro1, lag = 1, differences = 1)

# If var(d1) < var(d2) => linear
# If var(d2) < var(d1) => quadratic
print.default(paste0("The variance of the data differenced once at lag 1 is: ",
var(d1.electro1)))
print.default(paste0("The variance of the data differenced twice at lag 1 is: ",
var(d2.electro1)))

d1_ts_plot <- ts.plot(d1.electro1, ylab = expression(nabla-Y[t]),
                      main = "De-trended: Electricity Generation Differenced Lag 1")

par(mfrow = c(1,2))
acf(d1.electro1, lag.max = 60, main = "ACF of TS Differenced Lag 1")
pacf(d1.electro1, lag.max = 60, main = "PACF of TS Differenced Lag 1")

print.default(paste0("The variance of the de-trended time series is: ",
var(d1.electro1)))

hist(d1.electro1, main = "Histogram: Time Series Differenced Lag 1")

# Difference @ Lag 12 -> Remove Seasonality
d2.electro12 <- diff(d1.electro1, lag = 12, differences = 1)

```

```

d2_ts_plot <- ts.plot(d2.electro12, ylab = expression(nabla^{12}~\nabla Y[t]),
                     main = "De-seasonalized: Electricity Generation Differenced x2 Lag 12")

par(mfrow = c(1,2))
acf(d2.electro12, lag.max = 60, main = "ACF of TS Differenced x2 Lag 12")
pacf(d2.electro12, lag.max = 60, main = "PACF of TS Differenced x2 Lag 12")

print.default(paste0("The variance of the de-trended de-seasonalized time series is: ",
var(d2.electro12)))

# Augmented Dickey-Fuller Test for Stationarity
adf.test(d2.electro12, nlag = NULL, output = TRUE)

hist(d2.electro12, main = "Histogram: Time Series Differenced x2 Lag 12")

acf(d2.electro12, lag.max = 60, xlab = "Year (12 Lags Per Year)", main = "ACF of De-trended De-seasonalized Time Series")
pacf(d2.electro12, lag.max = 60, xlab = "Year (12 Lags Per Year)", main = "PACF of De-trended De-seasonalized Time Series")

# SARIMA modeling: (with orders chosen above)
set.seed(10)

# SARIMA parameter estimation
print.default(paste0("Model 1 (Observed From ACF/PACF) Summary & Estimated Parameters:"))
fit.sarima1 <- sarima(electro.bxcx, p=3, d=1, q=1, P=2, D=1, Q=2, S=12, xreg=1:length(electro.bxcx), de=1)
fit.sarima1
fit.sarima1$fit

# Acquire computer-generated model with auto.arima():
print.default(paste0("Model 2 (computer-generated model) Summary & Estimated Parameters:"))
auto.fit2 <- auto.arima(electro.bxcx)
auto.fit2

# Re-fit Model 2 Including Xreg
print.default(paste0("Re-fit Model 2 Including Xreg:"))
fit.sarima2 <- sarima(electro.bxcx, p=3, d=1, q=1, P=0, D=0, Q=2, S=12, xreg=1:length(electro.bxcx), de=1)
fit.sarima2
fit.sarima2$fit

# Analytically-selected model (from ACF/PACF)
print.default(paste0("The AICc for the analytically observed model (using ACF/PACF) is: ",
AICc(fit.sarima1$fit)))

# Computer-selected model (from auto.arima() function)
print.default(paste0("The AICc for the computer-generated model (using auto-arima) is: ",
AICc(fit.sarima2$fit)))

```

```

fit.s1.res <- fit.sarima1$fit$residuals

plot(fit.s1.res, ylab = "Residuals", main = "Residuals for Observed Model 1")
abline(h=mean(fit.s1.res),col="red")

# Sample Variance
print.default(paste0("The sample variance of residuals is: ",
var(fit.s1.res)))

op = par(mfrow = c(2,2))
hist(fit.s1.res, main = "Histogram of Residuals: Model 1", probability = TRUE)
qqnorm(fit.s1.res, main = "Normal Q-Q Plot: Model 1")
acf(fit.s1.res, lag.max = 60, main = "ACF of Residuals: Model 1")
pacf(fit.s1.res, lag.max = 60, main = "ACF of Residuals: Model 1")

# Residual Outlier Examination With Boxplot
boxplot(fit.s1.res, main = "Outlier Examination: Model 1")
print.default(paste0("Boxplot Statistics: "))
boxplot.stats(fit.s1.res)

# Remove Outlier(s) (assign adjusted residual dataset to variable)
fixie1 <- fit.s1.res[!fit.s1.res %in% boxplot.stats(fit.s1.res)$out]

# Residual Plot (corrected for outliers)
plot(fixie1, ylab = "Residuals", main = "Adjusted Residuals for Observed Model 1 (Outliers Removed)")
abline(h=mean(fixie1),col="red")

# Sample Variance
print.default(paste0("The sample variance of residuals without outliers is: ",
var(fixie1)))

# Other Plots for Residuals (corrected for outliers)
op = par(mfrow = c(2,2))
hist(fixie1, main = "Histogram of Residuals (Outlier Adj): Model 1", probability = TRUE)
lines(density(fixie1),col="red")
qqnorm(fixie1, main = "Normal Q-Q Plot (Outlier Adj): Model 1")
acf(fixie1, lag.max = 60, main = "ACF of Residuals (Outlier Adj): Model 1")
pacf(fixie1, lag.max = 60, main = "ACF of Residuals (Outlier Adj): Model 1")

hist(fixie, main = "Histogram of Residuals (Outlier Adj): Model 1", breaks = 20, probability = TRUE)
lines(density(fixie),col="red")

# Portmanteau Tests
# Box Pierce Test
Box.test(fit.s1.res, lag = 11, type = c("Box-Pierce"), fitdf = 4)

# Ljung-Box Test
Box.test(fit.s1.res, lag = 11, type = c("Ljung-Box"), fitdf = 4)

```

```

# Mcleod-Li Test
print.default(paste0("Mcleod-Li Test:"))
Box.test((fit.s1.res)^2, lag = 11, type = c("Ljung-Box"), fitdf = 0)

# Shapiro-Wilk Normality Test
shapiro.test(fit.s1.res)

# Check fitted residuals to AR(0), i.e. White Noise
print.default(paste0("Residuals Fitted to AR(0), White Noise:"))
ar(fit.s1.res, aic = TRUE, order.max = NULL, method = c("yule-walker"))

print.default(paste0("Residual Plots: "))
tsdiag(arima(fit.s1.res, order=c(0,0,0)))
print.default(paste0("Residual Plots (outliers removed): "))
tsdiag(arima(fixie1, order=c(0,0,0)))

# Pre-Loaded Function (given by Professor Feldman)
source("plot.roots.R")

ARs.coefs <- c(1, 0.1771, -0.1540, -0.1382)
MAs.coefs <- c(1, -0.7430)
SARs.coefs <- c(1, 0.0247, -0.3035)
SMAs.coefs <- c(1, -0.6980, -0.0172)

# MA Roots
print.default(paste0("MA Roots:"))
polyroot(MAs.coefs)
# AR Roots
print.default(paste0("AR Roots:"))
polyroot(ARs.coefs)
# SMA Roots
print.default(paste0("SMA Roots:"))
polyroot(SMAs.coefs)
# SAR Roots
print.default(paste0("SAR Roots:"))
polyroot(SARs.coefs)

plot.roots(NULL, polyroot(MAs.coefs), main="Roots of MA Component")
plot.roots(NULL, polyroot(ARs.coefs), main="Roots of AR Component")
plot.roots(NULL, polyroot(SMAs.coefs), main="Roots of SMA Component")
plot.roots(NULL, polyroot(SARs.coefs), main="Roots of SAR Component")

fit.s2.res <- fit.sarima2$fit$residuals

plot(fit.s2.res, ylab = "Residuals", main = "Residuals for auto.arima() Model 2")
abline(h=mean(fit.s2.res), col="red")

```

```

# Sample Variance
print.default(paste0("The sample variance of residuals is: ",
var(fit.s2.res)))

op = par(mfrow = c(2,2))
hist(fit.s2.res, main = "Histogram of Residuals: Model 2", probability = TRUE)
qqnorm(fit.s2.res, main = "Normal Q-Q Plot: Model 2")
acf(fit.s2.res, lag.max = 60, main = "ACF of Residuals: Model 2")
pacf(fit.s2.res, lag.max = 60, main = "ACF of Residuals: Model 2")

# Residual Outlier Examination With Boxplot
boxplot(fit.s2.res, main = "Outlier Examination: Model 2")
print.default(paste0("Boxplot Statistics: "))
boxplot.stats(fit.s2.res)

# Remove Outlier(s) (assign adjusted residual dataset to variable)
fixie <- fit.s2.res[!fit.s2.res %in% boxplot.stats(fit.s2.res)$out]

# Residual Plot (corrected for outliers)
plot(fixie, ylab = "Residuals", main = "Adjusted Residuals for auto.arima() Model 2 (Outliers Removed)",
abline(h=mean(fixie),col="red"))

# Sample Variance
print.default(paste0("The sample variance of residuals without outliers is: ",
var(fixie)))

# Other Plots for Residuals (corrected for outliers)
op = par(mfrow = c(2,2))
hist(fixie, main = "Histogram of Residuals (Outlier Adj): Model 2", probability = TRUE)
lines(density(fixie),col="red")
qqnorm(fixie, main = "Normal Q-Q Plot (Outlier Adj): Model 2")
acf(fixie, lag.max = 60, main = "ACF of Residuals (Outlier Adj): Model 2")
pacf(fixie, lag.max = 60, main = "ACF of Residuals (Outlier Adj): Model 2")

hist(fixie, main = "Histogram of Residuals (Outlier Adj): Model 2", breaks = 20, probability = TRUE)
lines(density(fixie),col="red")

# Box Pierce Test
Box.test(fit.s2.res, lag = 11, type = c("Box-Pierce"), fitdf = 4)

# Ljung-Box Test
Box.test(fit.s2.res, lag = 11, type = c("Ljung-Box"), fitdf = 4)

# Mcleod-Li Test
print.default(paste0("Mcleod-Li Test:"))
Box.test((fit.s2.res)^2, lag = 11, type = c("Ljung-Box"), fitdf = 0)

# Shapiro-Wilk Normality Test
shapiro.test(fit.s2.res)

```

```

# Check fitted residuals to AR(0), i.e. White Noise
ar0mod <- ar(fit.s2.res, aic = TRUE, order.max = NULL, method = c("yule-walker"))

print.default(paste0("Residual Plots: "))
tsdiag(arima(fit.s2.res, order=c(0,0,0)))
print.default(paste0("Residual Plots (outliers removed): "))
tsdiag(arima(fixie, order=c(0,0,0)))

# Pre-Loaded Function (given by Professor Feldman)
source("plot.roots.R")

AR.coefs <- c(1, 0.4530, -0.1881, -0.2571)
MA.coefs <- c(1, -0.9545)
SMA.coefs <- c(1, 0.9643, 0.4900)

# MA Roots
print.default(paste0("MA Roots:"))
polyroot(MA.coefs)
# AR Roots
print.default(paste0("AR Roots:"))
polyroot(AR.coefs)
# SMA Roots
print.default(paste0("SMA Roots:"))
polyroot(SMA.coefs)

plot.roots(NULL, polyroot(MA.coefs), main="Roots of MA Component")

plot.roots(NULL, polyroot(AR.coefs), main="Roots of AR Component")

plot.roots(NULL, polyroot(SMA.coefs), main="Roots of SMA Component")

print.default(paste0("AR(0) fit, White Noise Variance:"))
ar0mod

print.default(paste0("Final Model Fit, White Noise Variance:"))
auto.fit2

# Forecast for Model 2 (Best Model)
print.default(paste0("10-Step Ahead Forecast On Transformed Data:"))
sarima.pred <- sarima.for(electro.bxcx, 10, p=3, d=1, q=1, P=0, D=0, Q=2, S=12)
title(main = "10-Months Ahead Forecast for Transformed Data")

pred.val <- sarima.pred$pred
pred.se <- sarima.pred$se

u.lim <- pred.val+1.96*pred.se
l.lim <- pred.val-1.96*pred.se

```

```

pred.unbox <- InvBoxCox(sarima.pred$pred, lambda = lambda)

U.unbox <- InvBoxCox(u.lim, lambda = lambda)
L.unbox <- InvBoxCox(l.lim, lambda = lambda)

print.default(paste0("10-Step Ahead Forecast On Un-Transformed Data:"))
par(mfrow = c(1,1))
plot(electro_ts, xlim=c(1985,1997), ylim=c(min(electro_ts),320), ylab = expression(X[t]),
      xlab = "Date", main = "10-Months Ahead Forecast Against Test Set", cex.main=1.4)
space=1/12
k=0:11*space
indextoadd=(1996)+k[1:10]
points(indextoadd, pred.unbox, pch=1, col="red")
lines(U.unbox, lty="dashed", col="blue")
lines(L.unbox, lty="dashed", col="blue")

```