

STA 137
Final Project

Analysis of Daily Average Receipts per Theater
for the Movie *Chicago* over Time

Ian Xu (915463457)
ianxu@ucdavis.edu

Truc Le (914920690)
trlle@ucdavis.edu

Professor Prabir Burman
STA 137: Applied Time Series Analysis
University of California, Davis
March 16, 2021

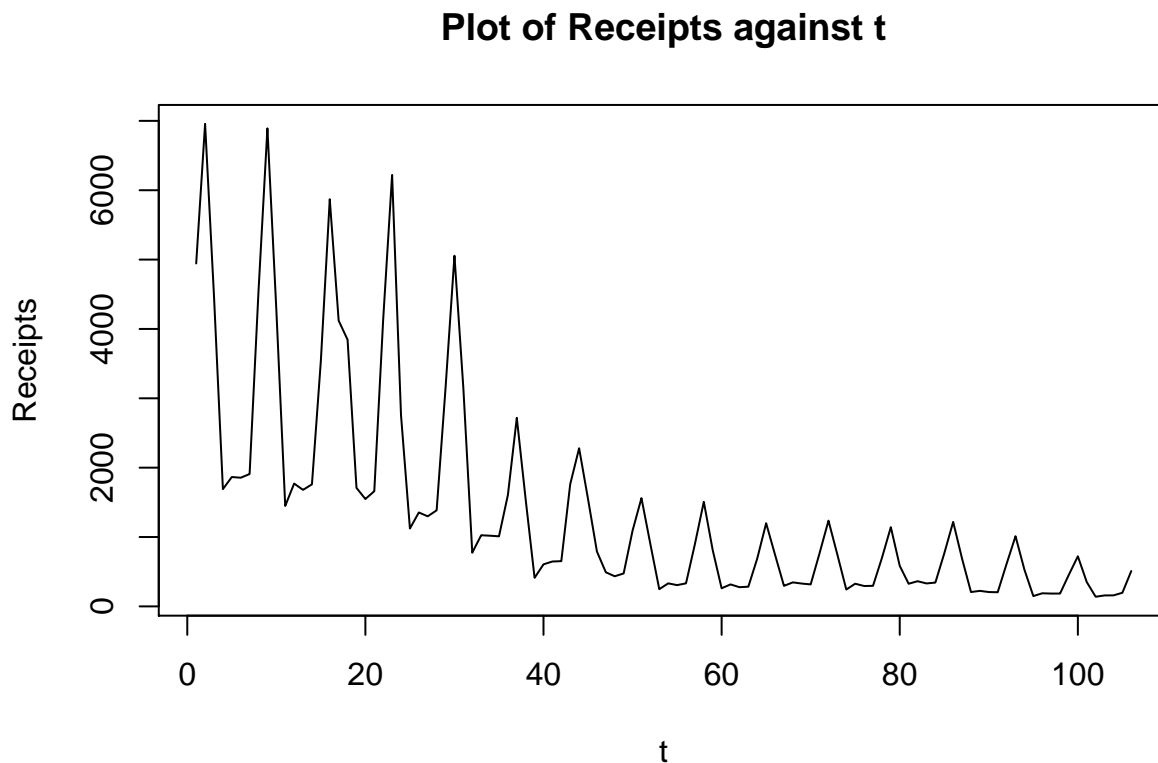
I. Introduction: Statement of Problem

The goal of this report is to fit a time series model to the Chicago dataset in order to predict the daily average receipts per theatres for the movie Chicago for the given four and a half months time period. The dataset contains the daily average receipts per theatres for the movie Chicago from January 3, 2003 up until April 18, 2003, and it is considered a time series since the data was taken over a period of time within equaled spaced increments. Time series modeling is important for this dataset as it can forecast the data points outside of the given time intervals by using model it already created from the current data points it possess. However, there are constraints to this time series prediction model, it is illogical to forecast before the movie was released and after it is no longer in theatre as it would give inaccurate predictions for the daily average receipts per theatres. In addition, like any time series model, it should be noted that any forecasting farther away from the given model would have a higher probability of inaccuracy in predicting the daily average receipts per theatres for the movie Chicago.

II. Materials and Methods: Description of the Data and Methods Used in the Analysis

For this project, we were provided a file (`Chicago.txt`), which contains daily average receipts (**Receipts**) per theater for the movie Chicago. The data covers the time period January 3, 2003 to April 23, 2003 (time $t = 1, \dots, 106$).

A simple plot of **Receipts** against **t**:



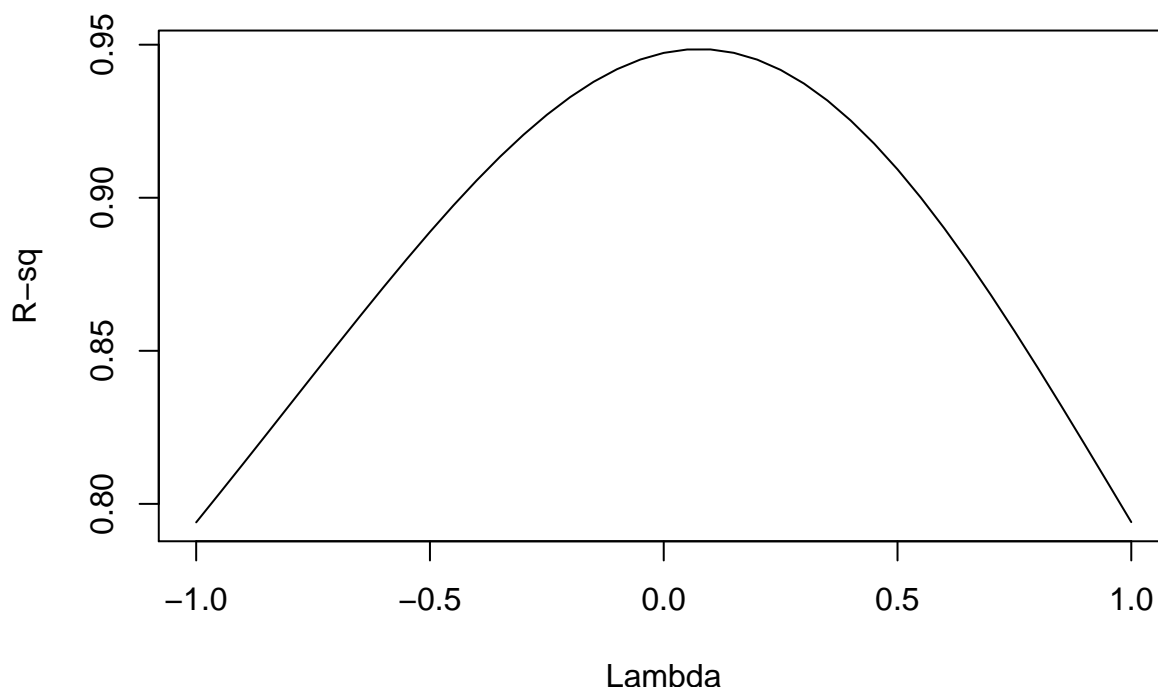
From the **Plot of Receipts against t**, it is observed that **Receipts** seems to have trend, seasonal, and rough components, with a cycle of roughly 7 days. A noticeable problem observed from the above plot is that the spread of **Receipts** changes over time (i.e. the variance is not constant). To address this problem, a Box-Cox transformation will be apply to the dataset, which will make the spread of **Receipts** more constant overtime.

To fit a time-series model to the dataset, we will use the basic model of $Y_t = m_t + s_t + X_t$, $t = 1, \dots, n = 106$,

where Y_t denotes **Receipts**, m_t denotes the trend, s_t denotes the seasonal component, and X_t denotes the rough part. To estimate m_t and s_t , the function `trndseas()`, which was provided in the file `trndseas.R`, will be used. After obtaining the \hat{m}_t and \hat{s}_t (the estimated trend and seasonal effects respectively), the rough part is then estimated by subtracting the estimated trend and seasonal effects from the observed values (i.e. $\hat{X}_t = Y_t - \hat{m}_t - \hat{s}_t$). After obtaining the estimated rough, the `auto.arima()` function will be used to find the best ARMA model. This ARMA model will then be used to estimate the rough in the time-series prediction model.

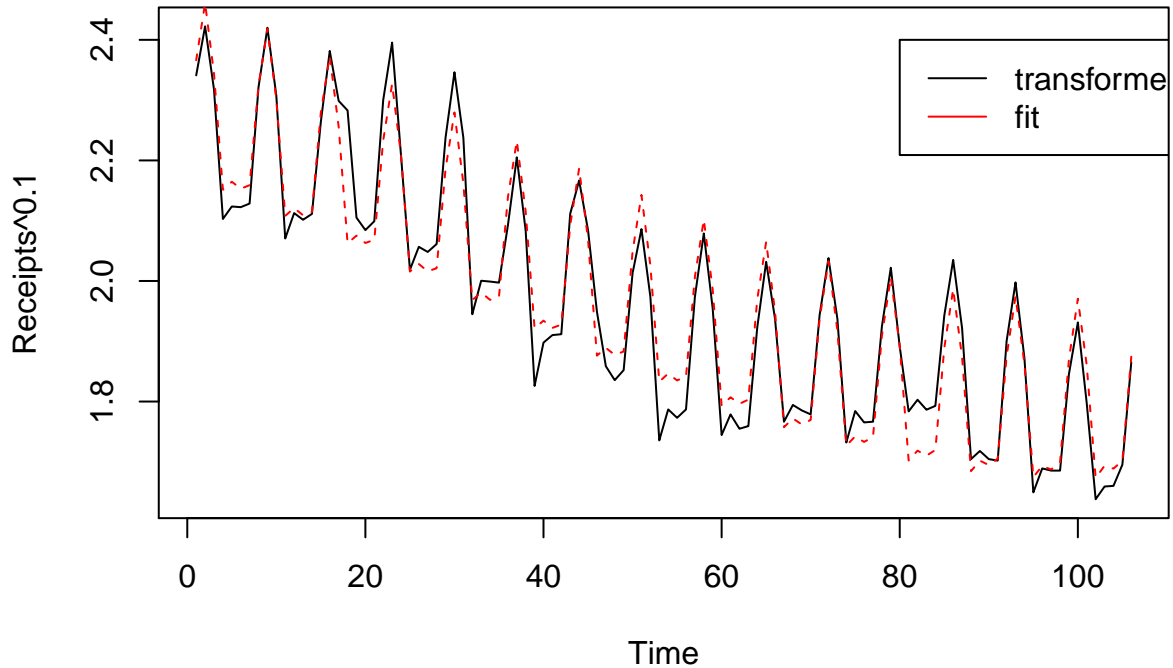
III. Results: Explanation of the Results in our Analysis

Average Ticket sales: R-square



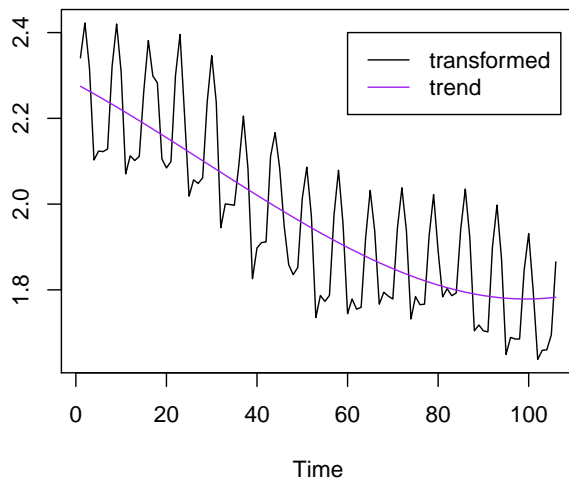
As previously discussed, due to the inconsistent variance observed from the **Plot of Receipts against t**, it is necessary to perform a Box-Cox transformation on the original data before a time series model could be fitted. In order to choose an appropriate Box-Cox transformation, the optimal lambda must be picked first. To find the optimal lambda value, a **Lambda vs. R-squared plot** is constructed to approximate the interval of where the optimal lambda is in. As seen in the **Lambda vs. R-squared plot**, it is observed that the optimal lambda is within the interval of $[-1,1]$, therefore, every values in increment of 0.05 from -1 to 1 was taken to find the optimal lambda values using the `trndseas()` function. The result from the `trndseas()` function shows that the optimal lambda value is 0.1. The data is transformed using the X_t^λ transformation model since the calculated optimal lambda in this case is not equal to 0.

Plot: Transformed and Fitted

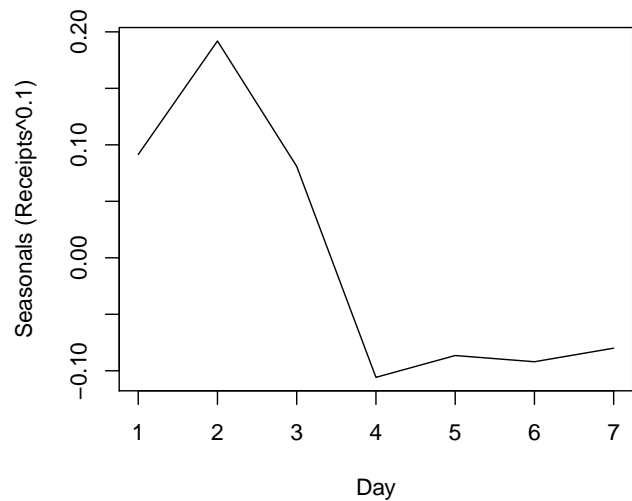


As observed from the black line (transformed) in **Plot: Transformed and Fitted** the transformed data gives a more consistent variance when compared to the original data as seen in the **Plot of Receipts against t**. Therefore, this newly transformed data will be use to fit the time series model moving forward.

Estimated Trend

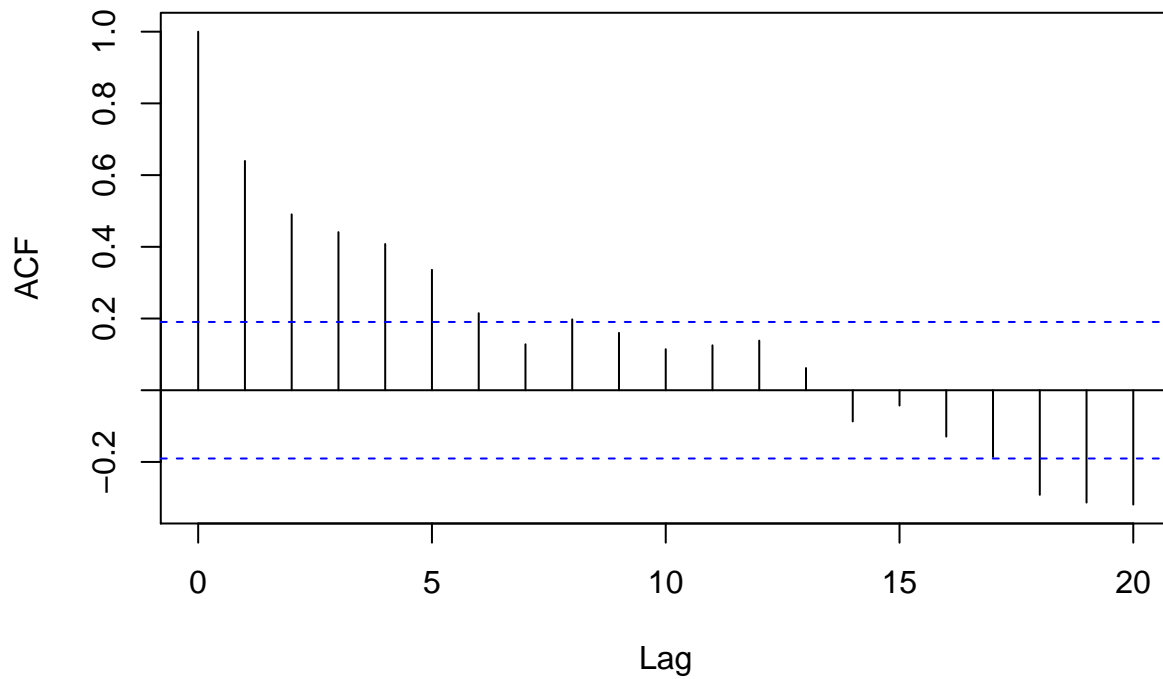


Estimated Seasonal Component



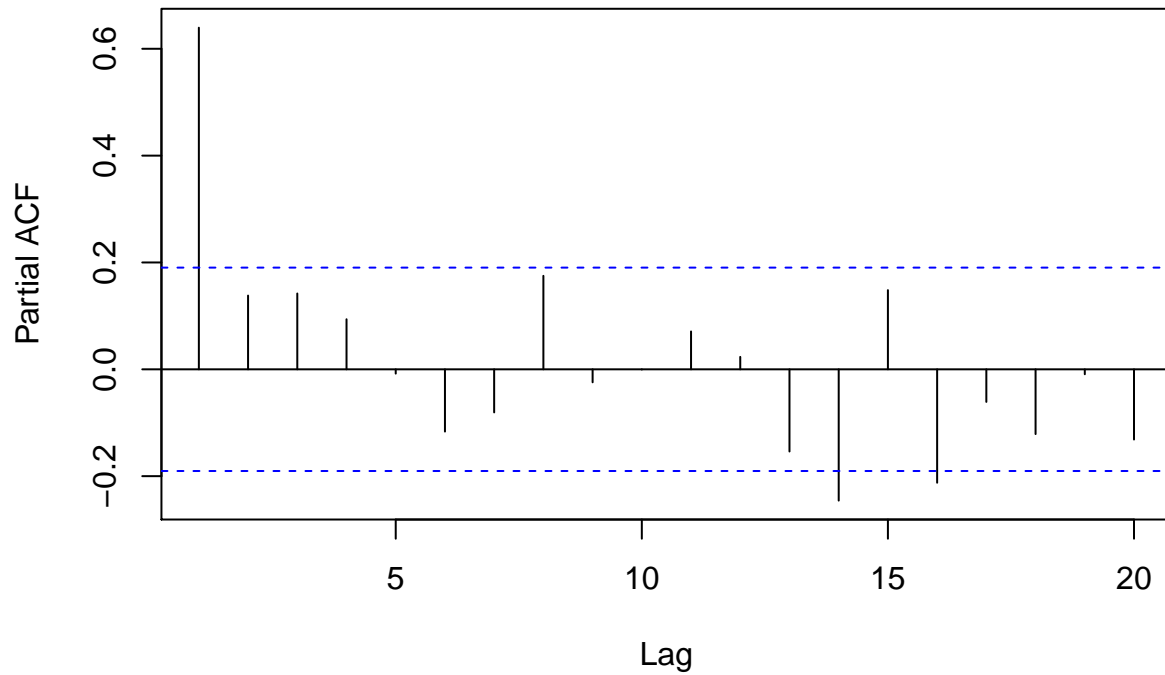
The next step after transforming the data using the optimal lambda value of 0.1, the trend and the seasonality components of the transformed dataset, m_t and s_t , were estimated using the `trndseas()` function. As observed in the **Estimated Trend** plot above, the estimated polynomial of degree 3 trend line fits the transformed data very well, as it lays superimposed over the middle of the transformed data while following the decreasing trend that the transformed line go through overtime. In addition, from the **Estimated Seasonals Component** plot, which only show one cycle, it can be observed that seasonality occurred in a cycle of 7 days before the cycle restarts again.

Series y_0.1 – fit



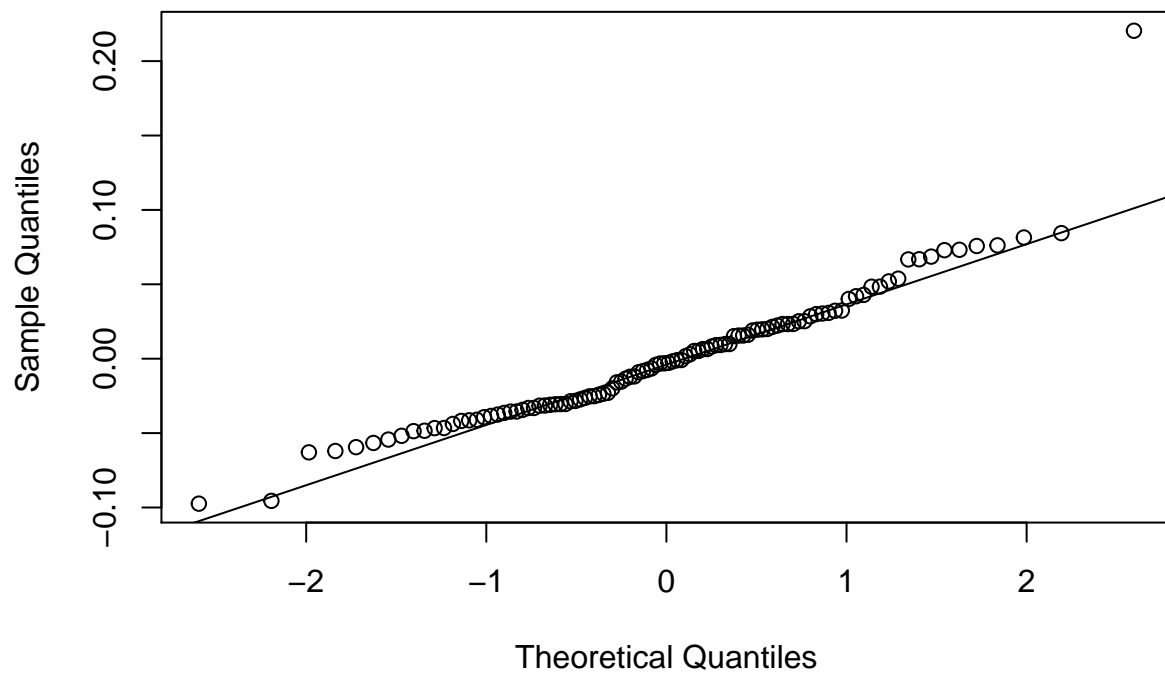
From the ACF plot of X_t , it is observed that the ACF values almost seem to be monotonously decreasing. There are some lags here and there where this observation does not hold true, but in general, the ACF values from lags 1 through 20 seem to be decreasing. Initially, lags 1 through 6 are significant. The lags then become insignificant, but they return to being significant at lag 18, and continue to be significant until lag 20. Based on the ACF plot, we should not model X_t using an MA model.

Series y_0.1 – fit



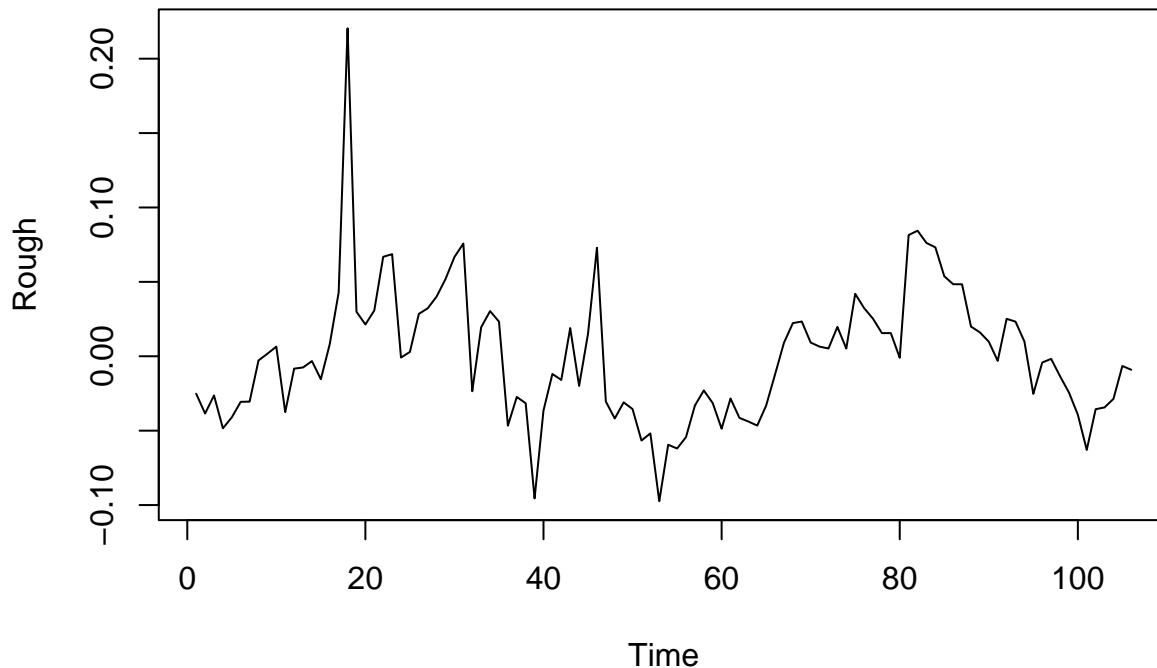
From the PACF plot of X_t , it is observed that there are not many significant lags. Lags 1, 14, and 16 are significant. If it is believed that lags 14 and 16 are not truly significant, but only appear significant due to randomness, then X_t could be modeled using an AR(1) model. However, if lags 14 and 16 are truly significant, then X_t cannot be modeled using an AR(1) model.

Normal Q–Q Plot



The normal Q-Q plot of \hat{X}_t shows that the roughs seem fairly normal, with the exception of the single outlier on the top right hand corner. While not all the points lie perfectly on the theoretical line, they also do not deviate too much from the line. However, because of the outlier and slight deviation from the theoretical line, we cannot say with confidence that \hat{X}_t is normally distributed.

Rough part



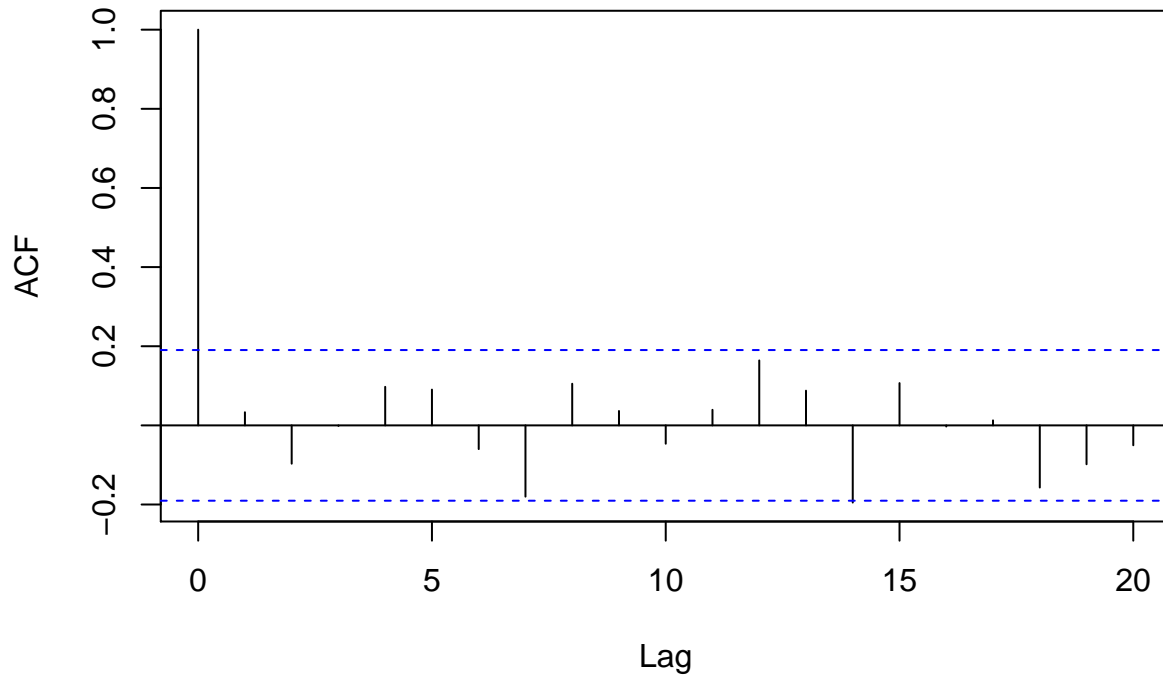
From the plot of \hat{X}_t against time, it is observed that the rough seems to have constant mean. The variance of the rough, on the other hand, does not appear to be constant. \hat{X}_{18} appears to be an “outlier” when compared with the other estimated roughs because it is much further away than the other estimated roughs. This is concerning because it suggests that the roughs do not have constant variance, and would thus violate the requirements for a weak stationary series. Other than this one outlier, however, the variance seems to be constant.

Based auto.arima and using the AICC criterion, the most appropriate model for \hat{X}_t is ARMA(1,1).

```
## Series: y_0.1 - fit
## ARIMA(1,0,1) with zero mean
##
## Coefficients:
##      ar1      ma1
##      0.8249 -0.3428
## s.e.  0.0862  0.1541
##
## sigma^2 estimated as 0.001129: log likelihood=209.96
## AIC=-413.92  AICc=-413.69  BIC=-405.93
```

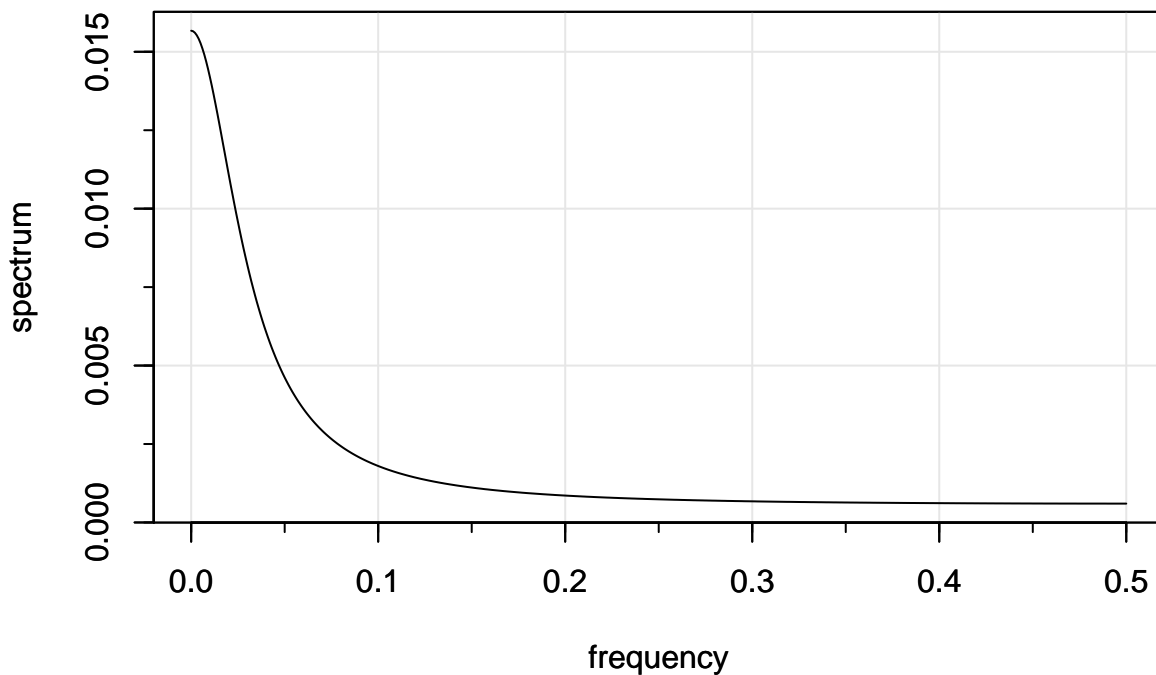
The ACF plot of the residuals:

Series res

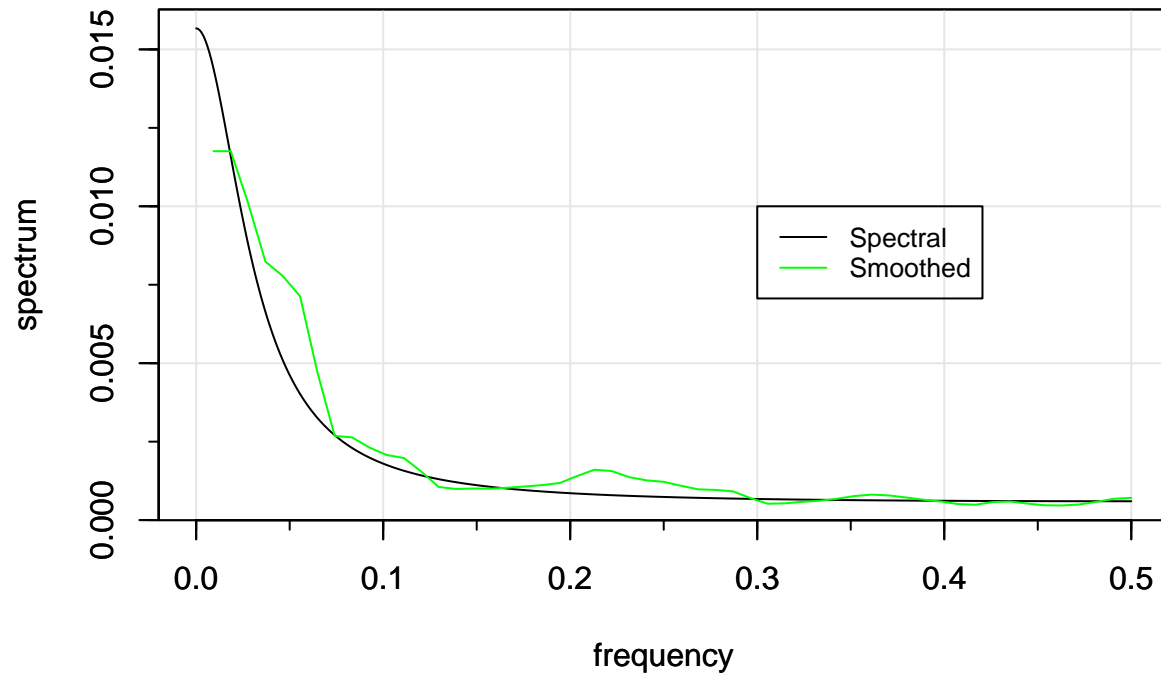


The ACF plot of the residuals of the ARMA(1,1) model show that the residuals seem to be independent. The only lag that is significant (other than lag 0) is lag 14. However, lag 14 is only barely significant and can be attributed to randomness. Therefore, it seems that the residuals of the ARMA(1,1) model is white noise and independent of one another. Thus, it seems that the ARMA(1,1) is an appropriate model for \hat{X}_t .

Spectral Density Graph for the ARMA(1,1) model

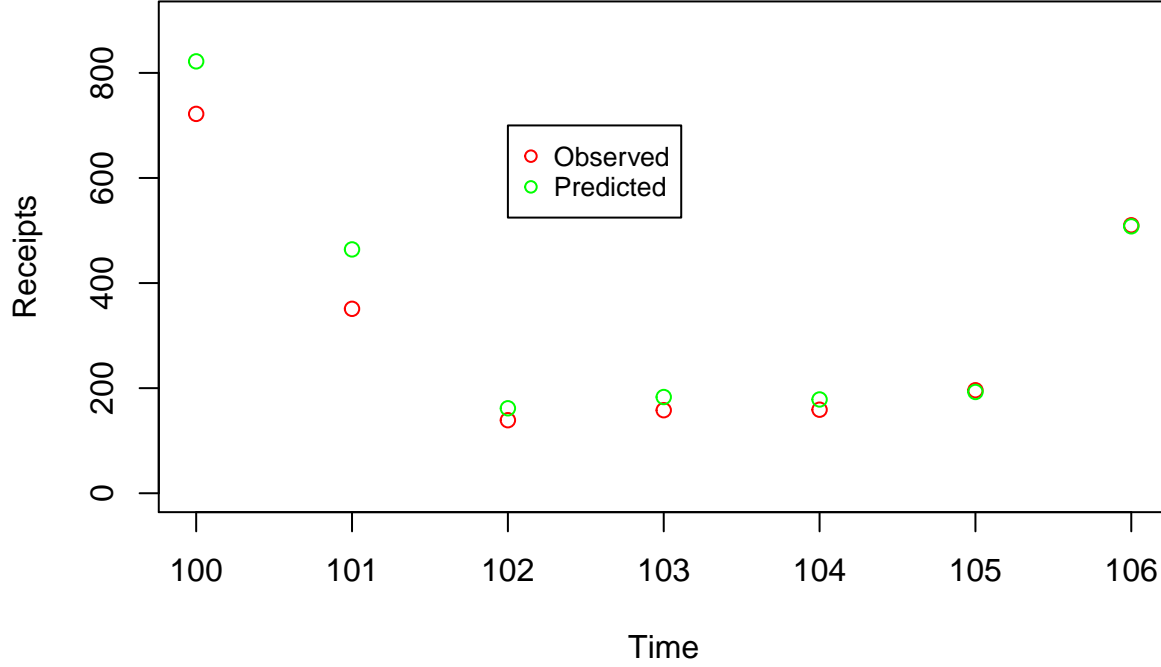


Spectral Density and Smooth Periodogram



The Spectral Density and Smooth Periodogram plot shows the smooth periodogram and the spectral density function of the fitted ARMA(1,1) model that was chosen to be the best model through the AIC criterion. As observed from the graph above, the smooth periodogram estimates the spectral density function fairly accurately, with the two biggest deviations located approximately at the frequency intervals of $[0.05, 0.075]$ and $[0.2, 0.25]$. Aside from the two deviations, the smooth periodogram does not deviate much from the spectral density functions line as shown above.

Observed and Predicted Values



In the **Observed and Predicted Values** plot above, the ARMA(1,1) model is fitted to the estimated rough with the exception of the last 7 days. The last 7 days were forecasted by using the estimated rough, and adding the 7 observations from the estimated rough to the estimated trend and estimated seasonality to get the prediction values. The 7 forecasted values were then plotted against the observed values on the same graph to observe how close the predictions were to observed values. As seen from the **Observed and Predicted Values** plot, the forecasted values using the ARMA(1,1) model gave a fairly accurate prediction for the last 7 days when compared to the observed values.

IV. Conclusion and Discussion

The objective of this report is to analyze daily average receipts per theater for the movie Chicago (Receipts) using time series methods. First, a Box-Cox transformation was applied to Receipts so that the variance of Receipts would be more constant overtime. Then, the trend m_t and the seasonal components s_t were estimated using the provided function `trndseas()`. After obtaining \hat{m}_t and \hat{s}_t , the estimated rough \hat{X}_t was obtained by subtracting \hat{m}_t and \hat{s}_t from transformed values of Receipts.

Next, using the AICC criterion and `auto.arima()`, an ARMA(1,1) model was selected to model X_t . Subsequent plots of the residuals of the ARMA(1,1) model show that there are no significant lags. This indicates that the ARMA(1,1) model is a good model for X_t .

Then, the spectral density function of the ARMA(1,1) model and a smoothed periodogram were plotted on the same graph. Using `specselect()`, a function provided in the discussion handout, it is observed that the most appropriate smoothed periodogram is when `span = 11` ($k = 5$). At this span, the smoothed periodogram achieves a balance between smoothing the periodogram and preserving detail.

Finally, using all the data except for the last 7 days, the ARMA(1,1) model was refitted. This model was used to forecast the last 7 observations. While the predictions were not perfect, they came considerably close to the true observed values.

Code Appendix

```
knitr::opts_chunk$set(echo = TRUE)
library(MASS)
source("https://raw.githubusercontent.com/Truc-T-Le/Time-Series-Project/main/trndseas.R")
chicago <- read.csv("https://raw.githubusercontent.com/Truc-T-Le/Time-Series-Project/main/chicago.csv")
library(forecast)
library(astsa)
library(Hmisc)
y <- chicago$reciepts
t <- 1:106
plot(t, y, type = "l", ylab = "Receipts", main = "Plot of Receipts against t")
#3/4/5
y<- chicago$reciepts
tm<- 1:106
seas <- 7
lam=seq(-1,1,by=0.05)
ff=trndseas(chicago$reciepts,seas,lam,3)
rsq=ff$rsq
ff=trndseas(chicago$reciepts,seas,0.1,3)
trend=ff$trend
season=ff$season
fit=ff$fit
Day=1:7
plot(lam,rsq,type="l",xlab="Lambda",ylab="R-sq",main="Average Ticket sales: R-square")

#fitted
{plot.ts(y^0.1,ylab="Receipts^0.1",main='Plot: Transformed and Fitted')
points(tm,fit,type='l',lty=2, col="red")
legend(80,2.4, c("transformed","fit"), lty=c(1,1,2), col=c("black", "red"))}

par(mfrow=c(1,2))
#trend line
{plot.ts(y^0.1,ylab="",main='Estimated Trend')
points(tm,trend,type='l', col="purple")
legend(60,2.4, c("transformed","trend"), lty=c(1,1,2), col=c("black", "purple"))}

plot(Day,season,type='l',ylab='Seasonals (Receipts^0.1)',main='Estimated Seasonal Component')
#6

y_0.1 <-(y)^0.1
acf(y_0.1-fit)
pacf(y_0.1-fit)
qqnorm(y_0.1-fit);qqline(y_0.1-fit)
plot(y_0.1-fit, type = "l", xlab = "Time", ylab = "Rough", main = "Rough part")
#7
auto.arima(y_0.1-fit, stepwise = F, approximation = F, ic = "aicc", max.order = 10)

mod_ARMA<- arima(y^0.1-fit, order =c(1,0,1))
c<-mod_ARMA$coef
d<-mod_ARMA$var.coef
e<-mod_ARMA$sigma2
```

```

res <- mod_ARMA$residuals

acf(res)

#8

specselect=function(y,kmax)
{
  ii=spec.pgram(y,log="no",plot=FALSE)
  ii=ii$spec
  cc=norm(as.matrix(ii),type="F")^2
  ctr=rep(1,kmax) ###criterion function
  for(k in 1:kmax)
  {
    ss=2*k+1; kk=1/(2*k)
    ff=spec.pgram(y,spans=ss,log="no",plot=FALSE)
    fspec=ff$spec
    ctr[k]=norm(as.matrix(ii-fspec),type="F")^2+kk*cc
  }
  kopt=which.min(ctr)
  result=list(ctr=ctr,kopt=kopt)
  return(result)
}

b<-specselect(y_0.1-fit,18)

coef.ar <- mod_ARMA$coef[1]
coef.ma <- mod_ARMA$coef[2]
sigma2 <- mod_ARMA$sigma2
mod_spec <-arma.spec(ar=coef.ar, ma=coef.ma, var.noise=sigma2, log='no', main = "Spectral Density Graph

#8.a
{arma.spec(ar=coef.ar, ma=coef.ma,var.noise = sigma2, log = "no",main="Spectral Density and Smooth Peri
smooth<-spec.pgram(y_0.1-fit, log = "no", spans = 11, main="", xlab="", ylab="", plot=F )
lines(smooth$freq, smooth$spec, col = "green")
legend(0.3, 0.010, legend = c("Spectral", "Smoothed"), col = c("black", "green"), lty = 1, cex = 0.8)}

#8.b

rough_no7 = y_0.1-fit
rough_no7 = rough_no7[-c(100,101,102,103,104,105,106)]
mod_ARMA11 = arima(rough_no7, order = c(1,0,1))

#forecast trend

row_old = 1:99
row_pred = c(100,101,102,103,104,105,106)
trend_pred = approxExtrap(row_old, ff$trend, xout= row_pred)$y

season_pred = rep(ff$season, length.out = length(chicago$reciepts))[-(1:99)]

forecast = predict(mod_ARMA11, n.ahead = 7)

```

```

x_fc = forecast$pred

truepredicted = x_fc+trend_pred+season_pred
untransform_predict = truepredicted^10
{plot(c(100,101,102,103,104,105,106), chicago$reciepts[100:106], col = "red", main = "Observed and Pred.
points(c(100,101,102,103,104,105,106),untransform_predict, col = "green")
legend(102, 700, legend = c("Observed", "Predicted"), col = c("red", "green"), pch = 1, cex = 0.8)}
# this is the code appendix

```