

PSTAT 174 Final Project

Time Series Analysis



Milk Production January 1962 - December 1975

Authors:

Whitney Pemberton
Melanie Flandes
Tomio Sugawara
Runchen Zhang
Jensen Newlin

Supervised By:

Professor Sudeep Bapat

Contents

Abstract	3
Introduction	3
Exploratory Analysis	3
Initial Data	3
Data Decomposition	4
Data Transformation	4
Changing Variance	4
Detrending	6
Deseasonalizing	7
Model Selection and Estimation	8
Causality and Invertibility	9
Diagnostics	10
Normality of Residuals	10
Serial Correlation	11
Heteroskedasticity	11
Forecasting	12
Conclusion	13
Future Study	13
References	14
Appendix	14

Abstract

Due to the necessity and mass consumption of milk in the U.S., the production of milk has always been an important topic to consider. The goal of this project is to apply time series analysis on milk production data, to predict future monthly production. Exploratory analysis through multiple plots is used to identify non stationary components. Next, methods such as Box-Cox and differencing are applied to transform the data into a stationary process. Then, model selection was conducted using AICc values. After, diagnostic checks were conducted to make sure the models satisfied assumptions. Finally, this SARIMA model was used to predict future monthly milk production.

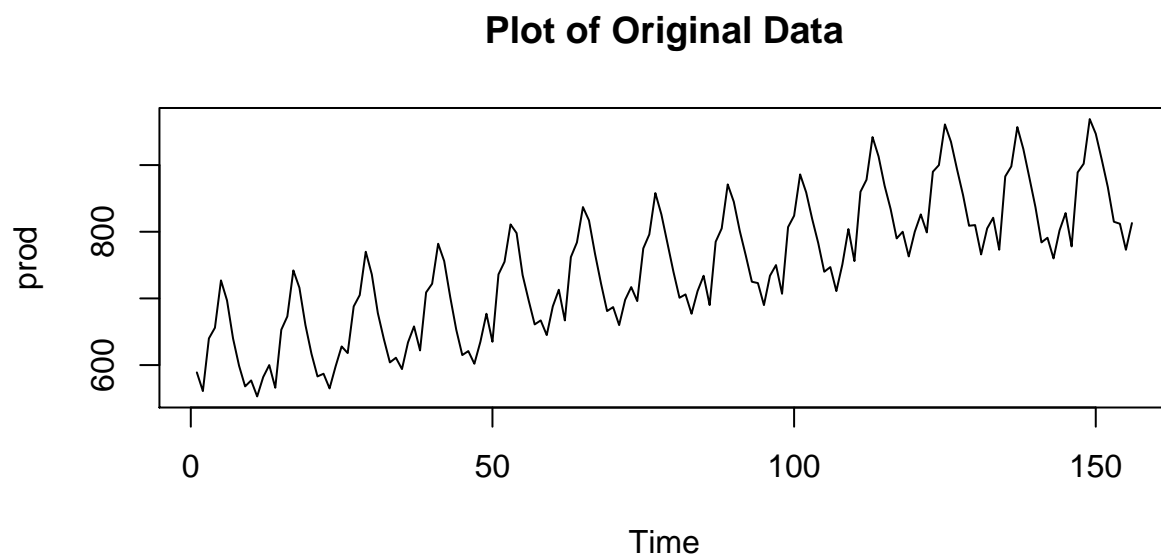
Introduction

Over the years, the United States population has increased drastically. Due to this increase, the production rates of all food products have also gone up. One product that is in constant production, due to its high demand, is milk. This increased demand makes it important to analyze trends to better predict future needs. The data provides monthly milk production (pounds per cow) in the United States starting from January 1962 to December 1975. We are interested in conducting time series analysis on this milk production dataset because we want to create forecasted values for future milk production. These forecasted values help producers understand how much milk needs to be produced in order to meet the future demand of consumers. To conduct this analysis we use R and R studio since it has many built in statistical packages which include functions needed to conduct time series analysis. Box-Cox transformation, detrending and deseasonalizing help us create a stationary model. After, we fit multiple SARIMA models using AICc to determine our final model, it is used to forecast future production. Our final model $SARIMA(1, 1, 1)x(1, 1, 1)_{12}$ is fitted on a subset of our milk production data (156 months).

Exploratory Analysis

Initial Data

The data set we are using has two variables: date, on a monthly basis and pounds per cow of milk produced. There are 168 observations in the data set but we leave out the last 12 months of data so we can asses performance later on. Here is the plot of the original non transformed data:

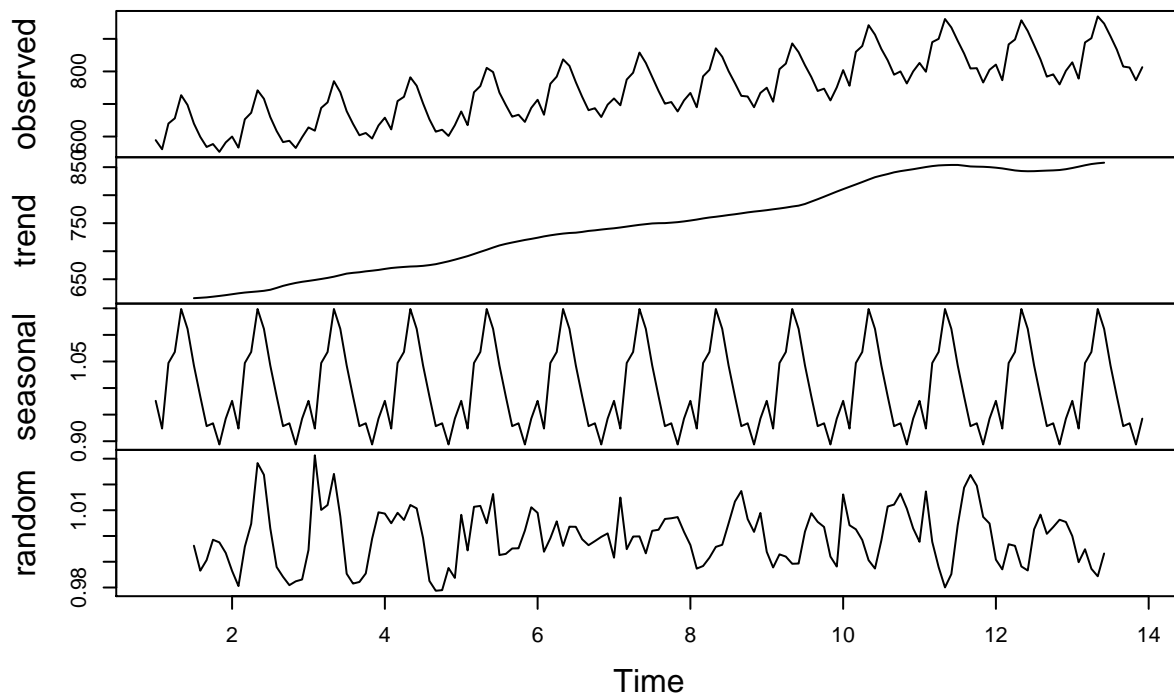


The data shows the presence of variations that occur at specific regular intervals which suggests there is seasonality. You can also see a positive linear trend and a changing variance.

Data Decomposition

The trend, seasonality and changing variance is further confirmed by our decomposition plot below. Here we can see the clear presence of a positive linear trend. Next, the seasonality repeats in intervals of 12 months. To correct these first two problems we will difference the data. Finally, there is the possibility of changing variance. A Box-Cox transformation may correct this final problem to make the data stationary.

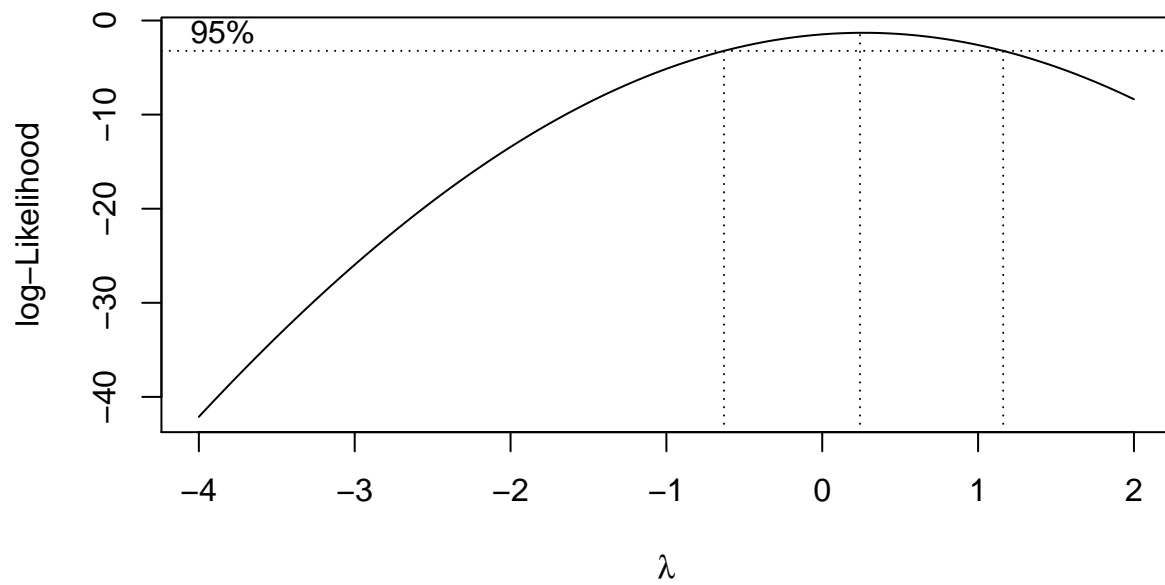
Decomposition of multiplicative time series



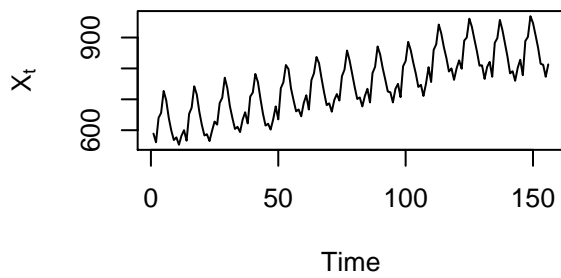
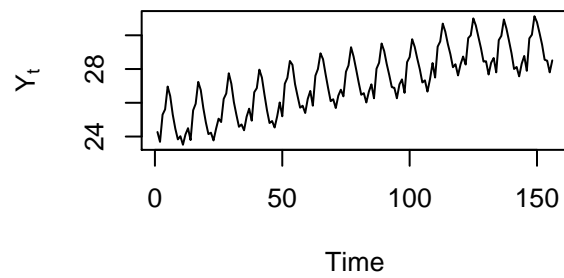
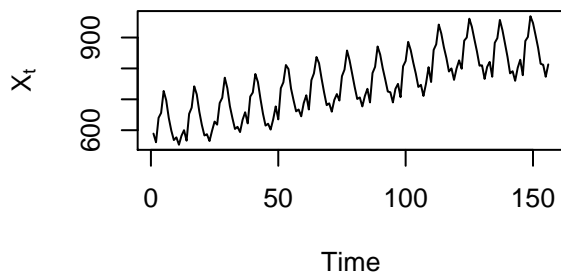
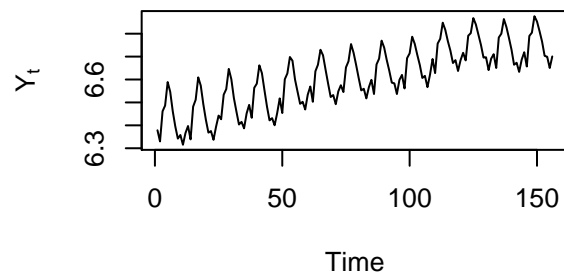
Data Transformation

Changing Variance

First, We are going to transform the data to try and make the variance portion stationary. Even though the variance looks to be constant we apply the Box-Cox transformation. After taking the Box-Cox transformation, the lambda value calculated is 0.24 which is close to 0.5 and 0. Therefore, since 0.5 and 0 are within the confidence interval both log and square root transformations are valid choices. Below is the plot of the original data vs. Sqrt transformed data and as you can see there is not much of a difference between them. Under that is the plot of the original data vs. Log transformed data which has the same result.



$$\lambda = 0.242424242424242$$

Plot Original Data**Plot Sqrt transformed Data****Plot Original Data****Plot Log transformed Data**

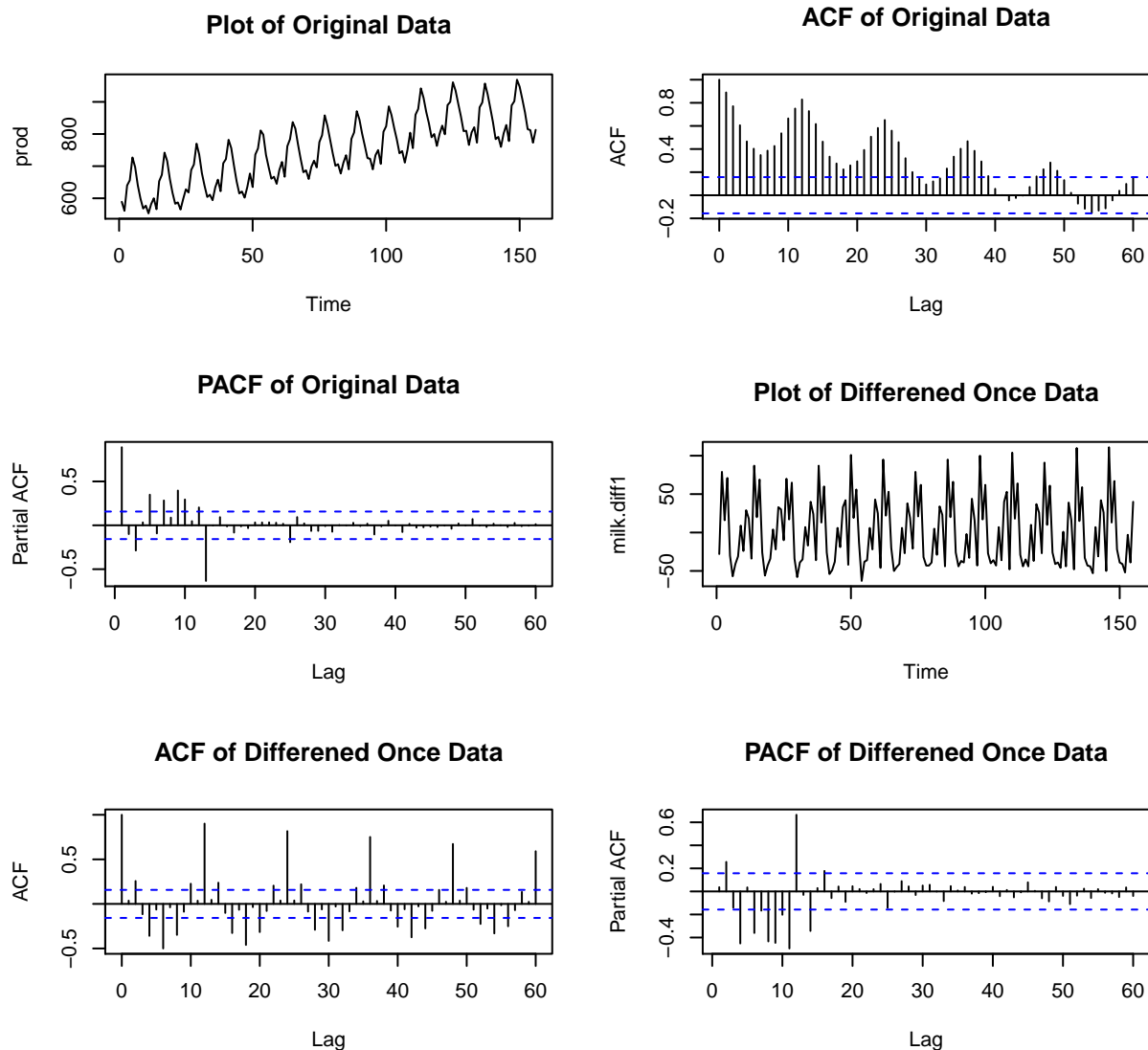
Since both the Sqrt transformation and Log transformation did not show much of a difference from the original data, we will not apply either of these transformations. This again shows that we have constant variance.

Detrending

First, we plotted the ACF and PACF for our original data. Both plots suggest non-stationarity, implying that we need to apply differencing techniques to detrend and deseasonalize our raw data. The time series plot of our original data shows an obvious upward trend. Therefore, we decide to difference our data once, then compare the original data and the differenced data. The variance goes down from 10055.58 to 2042.924 which shows this is a valid transformation to apply to our data. After, we analyze the differenced time series plot we can see that the upward trend vanishes, but there still exists a strong periodic trend, namely seasonality. We then plot the detrended data ACF and PACF, from which we find a spike pattern existing at lag 12, 24, 36, ... in the ACF plot and numbers of spikes before lag 12 in the PACF plot. These two findings suggest that the seasonal component has a period of 12.

$$\text{Variance}(\text{Original}) = 10055.5842431762$$

$$\text{Variance}(\text{DifferencedOnce}) = 2042.92392124005$$

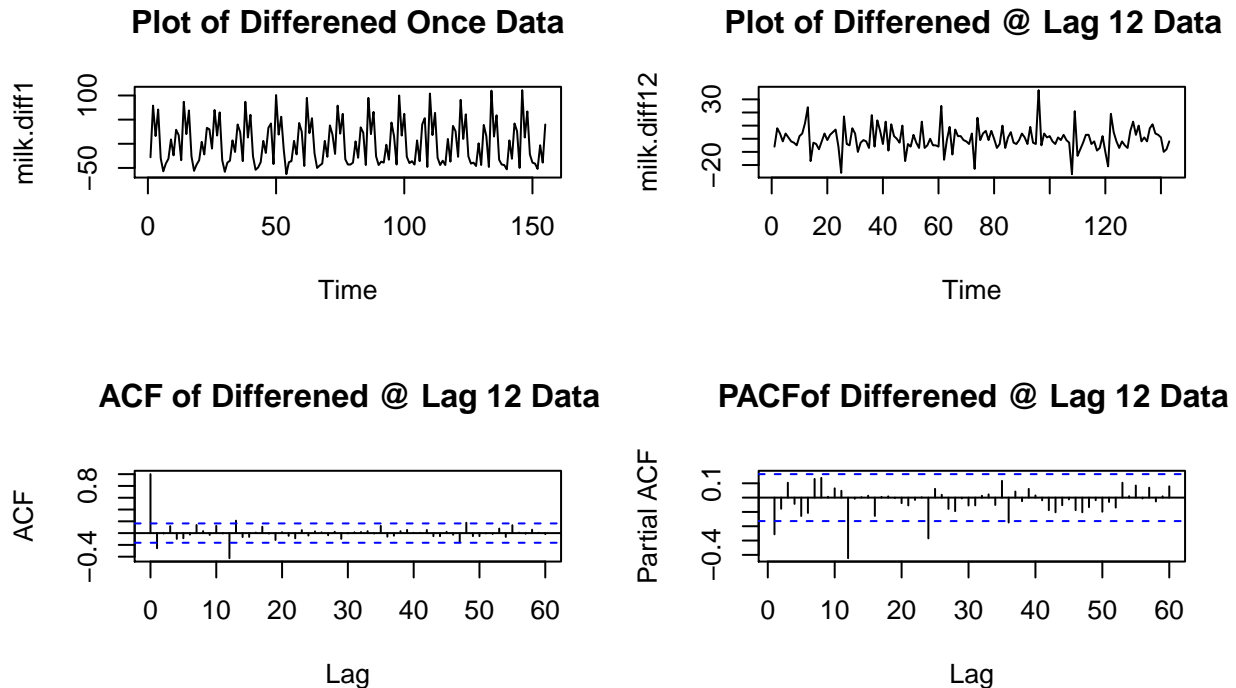


Deseasonalizing

We then proceed to remove seasonality from our data. To do this we difference our data in addition at $lag = 12$. The variance goes down significantly, from 2042.924 to 83.04225, indicating that the differencing was a reasonable option. By plotting the differenced data, we now see that the seasonality is removed and the spread of observations is much narrower than before differencing.

$$Variance(DifferencedOnce) = 2042.92392124005$$

$$Variance(Differenced@Lag12) = 83.0422535211268$$



There are large spikes at lag=12,24,36. Differencing at any of these lags in addition to our already differenced data results in an increase in variance. This would overfit our data so we will not apply any further differencing.

Now, through obtaining ACF and PACF plots for the detrended and deseasonalized data, we may come to the conclusion that the ACF suggests that the data is now stationary and that the cut-off is at lag 12 and at lag 1 so that we obtain a $Q = q = 1$. Also, the PACF suggests that there exist cut-off at lag 1 so that $p = 1$ and cut-off at lag 12, 24, and 36 so that $P = 1, 2$ or 3.

Model Selection and Estimation

Our decision of parameters will be: $p = 1$, $P = (1,2,3)$, $d = 1$, $D = 1$, $q = 1$, $Q = 1$. Therefore, the selection of parameters brings up our three final model candidates, stated as below:

$$SARIMA(1, 1, 1)x(1, 1, 1)_{12}$$

$$SARIMA(1, 1, 1)x(2, 1, 1)_{12}$$

$$SARIMA(1, 1, 1)x(3, 1, 1)_{12}$$

We transform the above expressions into R-code and name them as *fit1*, *fit2*, and *fit3* respectively. When calling them separately, they report their own AIC of 991.93, 992.83 and 994.87, with *fit1* having the smallest AIC. Instinctively, we would decide that *fit1* = $SARIMA(1, 1, 1)x(1, 1, 1)_{12}$ will be our final model choice, but before we go any further on determining the final model, we would like to check the causality and invertibility of all the possible models. Below are the ML estimations for each model.

fit1: $SARIMA(1, 1, 1)x(1, 1, 1)_{12}$

```
##
## Call:
## arima(x = prod, order = c(1, 1, 1), seasonal = list(order = c(1, 1, 1), period = 12),
##      method = "ML")
##
## Coefficients:
##          ar1          ma1          sar1          sma1
##      -0.0937  -0.1723  -0.0468  -0.5845
## s.e.    0.2732   0.2668   0.1214   0.0978
##
## sigma^2 estimated as 53.97:  log likelihood = -490.97,  aic = 991.93
```

fit2: $SARIMA(1, 1, 1)x(2, 1, 1)_{12}$

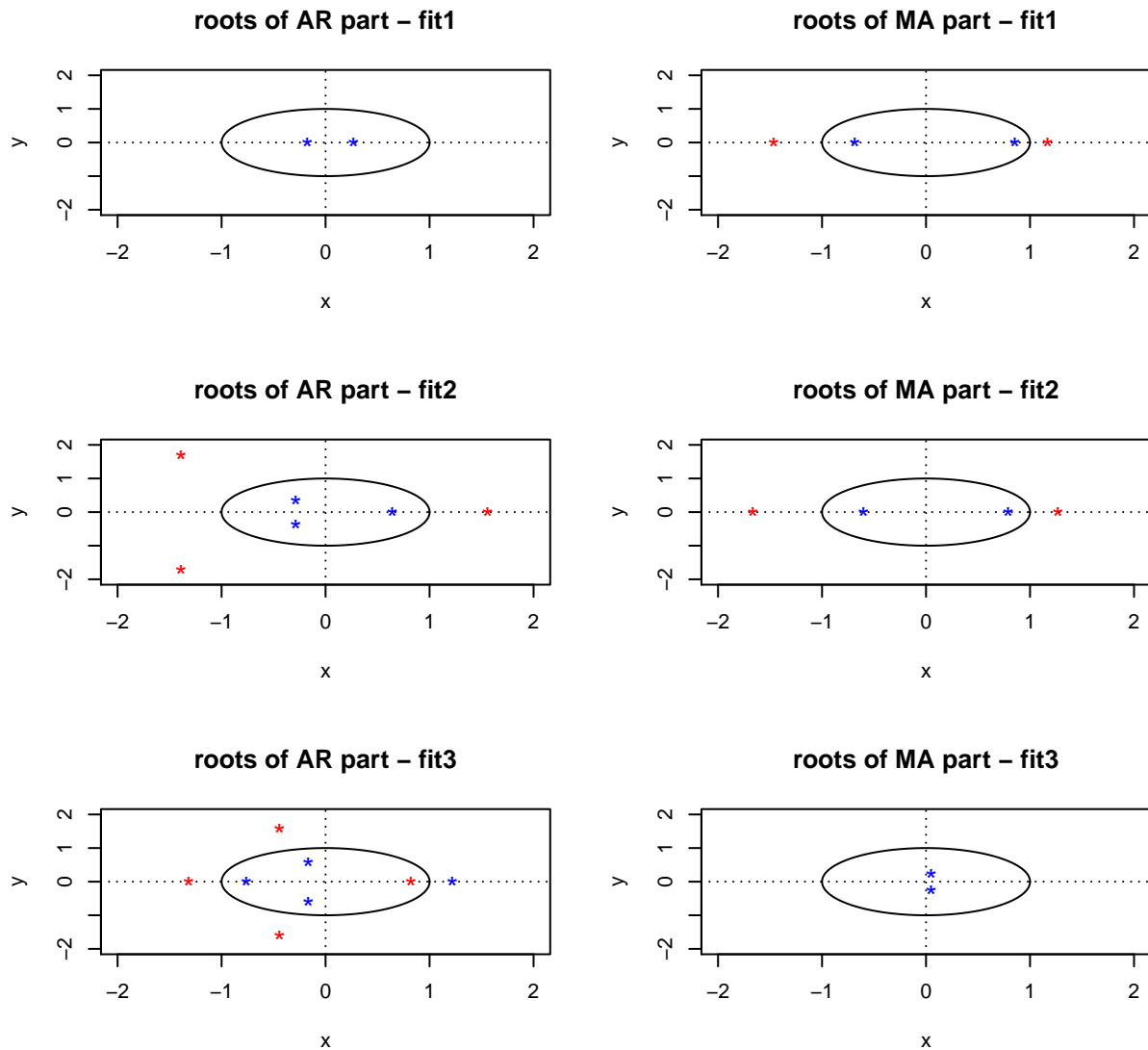
```
##
## Call:
## arima(x = prod, order = c(1, 1, 1), seasonal = list(order = c(2, 1, 1), period = 12),
##      method = "ML")
##
## Coefficients:
##          ar1          ma1          sar1          sar2          sma1
##      -0.0721  -0.1880  -0.1612  -0.1317  -0.4741
## s.e.    0.2346   0.2351   0.1612   0.1196   0.1531
##
## sigma^2 estimated as 53.39:  log likelihood = -490.42,  aic = 992.83
```

fit3: $SARIMA(1, 1, 1)x(3, 1, 1)_{12}$

```
##
## Call:
## arima(x = prod, order = c(1, 1, 1), seasonal = list(order = c(3, 1, 1), period = 12),
##      method = "ML")
##
## Coefficients:
##          ar1          ma1          sar1          sar2          sar3          sma1
##      -0.1241  -0.0805  -0.33   -0.1079  -0.1422  -0.3547
## s.e.         NaN         NaN         NaN    0.0019    0.0043         NaN
##
## sigma^2 estimated as 53.34:  log likelihood = -490.44,  aic = 994.87
```


Causality and Invertibility

Below we can see the plots of the polynomial roots both for the AR(p) and MA(q) portion of each model. The red dots stand for the root and the blue dots are 1/root. If the point is on an axis that means it is a real root, if not it is imaginary. Also, if there is no red dot on the plot that means the root is greater than our plot bounds, which would make those models Causal and Invertible. Here we can see that only one of the red points lie inside the unit circle (AR(p) for fit 3) which implies that all other models are Causal and Invertible. Therefore, we can disregard fit3 and look at fit1 and fit2 as candidates for our final model.



Now we proceed to choose the final model. We adopt the Second-order Akaike Information Criterion (AIC_c) to further help us with choosing a best model. From running the AIC_c function we found that our AIC_c for Model 1 (fit1) to be 992.36 and Model 2 AIC_c (fit2) was 993.447. We found that AIC and AIC_c were lowest for fit1 thus we have decided to use it as our final model.

$$fit1 : SARIMA(1, 1, 1)x(1, 1, 1)_{12}$$

$$(1 - .0937B)(1 - .0468B^{12}) \nabla_{12} \nabla X_t = (1 - .1723B)(1 - .5845B^{12})Z_t$$

$$where Z_t \sim WN(0, 49.814)$$

Diagnostics

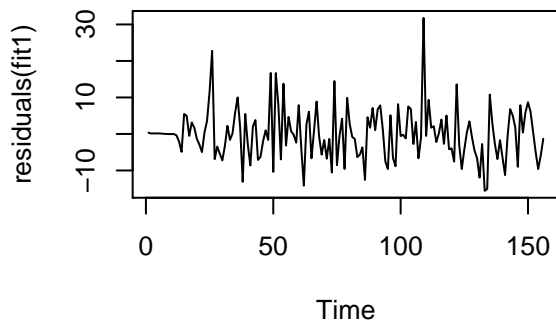
Normality of Residuals

To ensure that our final model satisfies all necessary assumptions, we conduct the diagnostic checks by taking the characteristics of the final model residuals into account. Then we plot the histogram of the residuals, from which we are able to find out if there are any problems regarding the distribution of the residuals. The histogram indicates a slight skew towards the right, but for the most part it looks like the residuals are normally distributed.

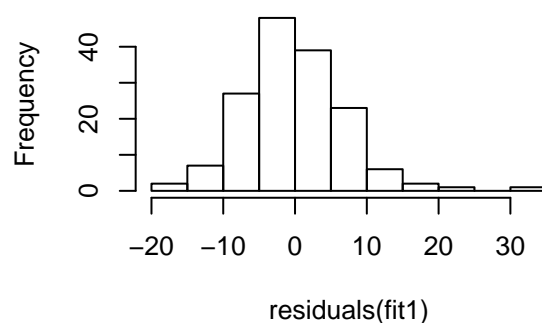
To find out if the residuals are indeed normally distributed, we then apply a more accurate test to see if the normality assumption is met, i.e. Normal Q-Q Plot. Similar to the histogram, the plot seems a bit skewed around the right tail, with a few observations deviating from the main track of the observations. However, after discreet consideration, we deem these observations as simple outliers, considering that the major portion of the points is normally distributed.

What is more, the Shapiro Wilk test reports a $p - value = 0.0002331$, which is less than .05, indicating a rejection of assumption of normality but eventually regarded acceptable considering our histogram and Normal Q-Q Plot are a bit skewed.

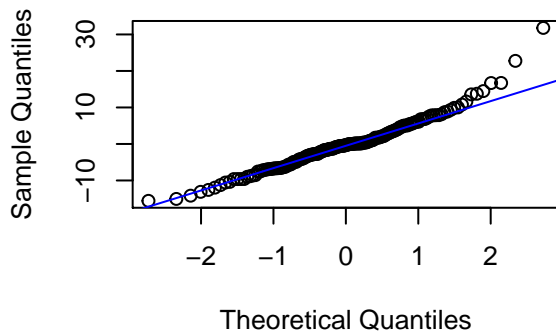
Plot of Final Model Residuals



Histogram of Final Model Residuals



QQ-Plot of Final Model Residuals



```
##
## Shapiro-Wilk normality test
##
## data: residuals(fit1)
## W = 0.96122, p-value = 0.0002332
```

Serial Correlation

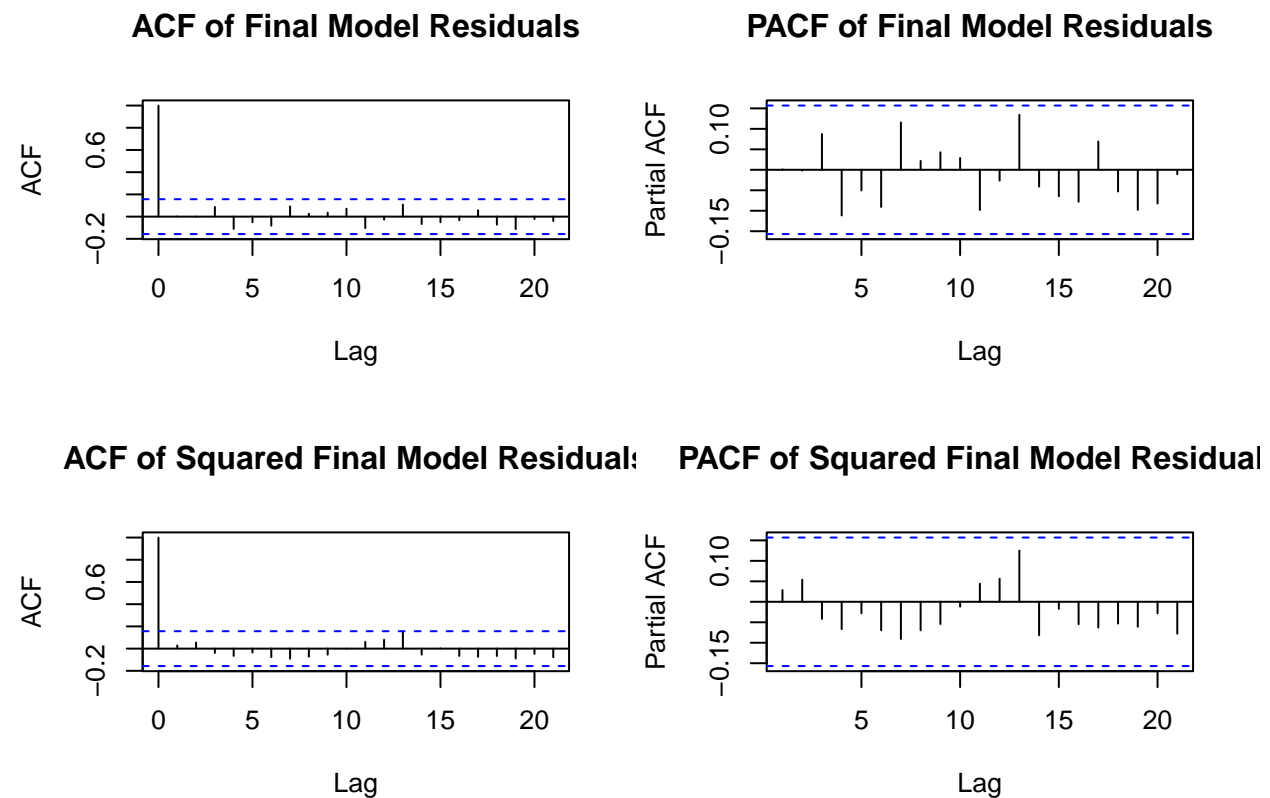
Next, we perform the Box-Pierce and the Ljung-Box tests on our model. Before testing, it is crucial that we set up and declare relevant parameters. We decide the parameters as follow: $lag = h = \sqrt{n}$, $n = 156$, $\sqrt{156} \approx 12$, $p = 1$, $q = 1$, $fit(df) = 1 + 1 = 2$.

We first conduct the Box-Pierce test with $lag = 12$, $fit(df) = 2$ and the reported test statistics, namely the $p - value = 0.558$. Then we test the Ljung-Box test with $lag = 12$, $fit(df) = 2$, with a resulting $p - value = 0.5091$. Finally, we tested the Ljung-Box test with a $fit(df) = 0$, and this time the $p - value = 0.682$. To conclude, we have p-values for all the box tests greater than .05, indicating that the autocorrelations are within the lags 1 through 11, therefore no serial correlation. Also, we may conclude that the residuals are independent.

```
##
## Box-Pierce test
##
## data: residuals(fit1)
## X-squared = 8.7289, df = 10, p-value = 0.558
```

Heteroskedasticity

Lastly, we check Heteroskedasticity by plotting the ACF and PACF of the error terms and of the squared error terms. Below we can see that in both sets the ACF and PACF all have spikes within the confidence interval after lag=1. This means that we do have constant variance of the error terms. Our final model has now passed all diagnostic checks.



Forecasting

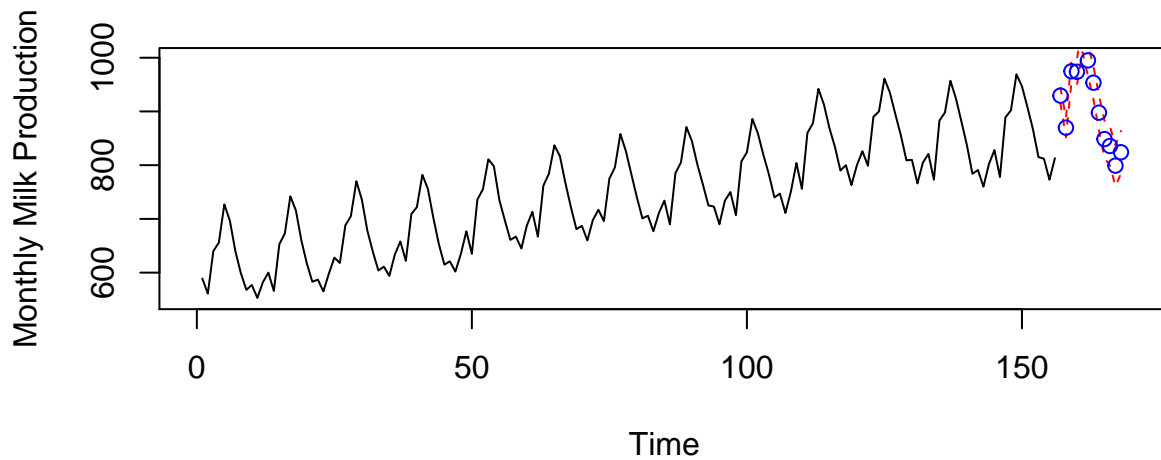
After checking for normality of the residuals and performing diagnostic tests on our final model, we deemed that we could now perform forecasting on data using our *Model* : $SARIMA(1, 1, 1)x(1, 1, 1)_{12}$

$$(1 - .0937B)(1 - .0468B^{12}) \nabla_{12} \nabla X_t = (1 - .1723B)(1 - .5845B^{12})Z_t$$

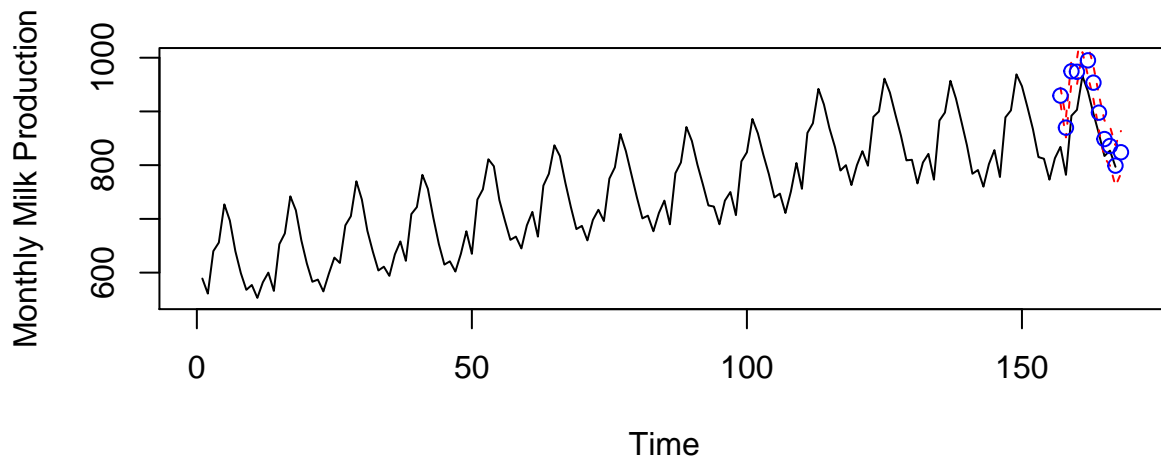
$$\text{where } Z_t \sim WN(0, 49.814)$$

To perform forecasting we used the predict function in R and set the function to predict for the next 12 months of data: January 1975 to December 1975. We then calculated a 95 percent confidence interval which included predictions for lower and upper bounds of the next 12 months of data. We then plotted our predicted values. The red lines indicate the bounds of our confidence interval and the blue dots indicate the actual predicted values. To assess how well our model actually predicted the data we plotted our predictions versus the actual 12 months of data for January 1975 to December 1975. We can see from our plot that our model tended to overestimate a bit compared to actual data.

Forecasted Values for Milk Production



Forecasted Values for Milk Production vs Actual Values



Conclusion

The goal of our project was to fit a model that would forecast milk production. To accomplish this we used the Box Jenkins approach to model building. We started by plotting the original milk production. The data had not been transformed or altered in any way. We noticed issues with our original data such as seasonality, trend, and possibly high variance. To remedy the variance issue we first looked into transforming the data. To figure out what transformation is best to use on our data we used the Box-Cox method. Using Box-Cox we found our lambda value to be .24. Based on our lambda value we decided a sqrt or log transformation would be best since they were both within the confidence interval. We plotted our transformed data and found that there was not a significant difference between our transformed data vs our non transformed data. Therefore, we decided that it would not be worthwhile to continue with our transformed data. We then proceeded to try and difference on our original data in hopes of removing the trend component. We differenced twice in our project, once to remove the trend and again at $lag = 12$ to remove the seasonality. By differencing we also noticed that the variance had went down significantly, so we decided that at this point our data was stationary and could be used for model fitting. We came up with three possible SARIMA models for our data. To compare models and assess which model was the best we used diagnosis such as AIC and AICc. After settling on one specific model we proceeded to perform further diagnosis checking on our model.

Model : SARIMA(1, 1, 1)x(1, 1, 1)₁₂

$$(1 - .0937B)(1 - .0468B^{12}) \nabla_{12} \nabla X_t = (1 - .1723B)(1 - .5845B^{12})Z_t$$

$$where Z_t \sim WN(0, 49.814)$$

This included plotting their roots to determine invertibility and causality as well as residuals to assess normality. We also used traditional diagnostic checking tests such as Box-Pierce, Box-Ljung, and Shapiro Wilk normality test. Our final model passed all diagnostic tests except for Shapiro Wilk normality test. From this we deemed that our model was appropriate and could be used for forecasting. We then proceeded to forecast the next 12 months of data from January 1975 to December 1975. Comparing our forecasts to the actual data from those 12 months showed that our model had for the most part accurately predicted milk production. The weaknesses in our model were that it tended to overestimate milk product across the 12 months. In the end our goal of forecasting for January 1975 to December 1975 was achieved. However, there is definitely room to improve upon our model. We also should look into testing our model and forecasting for more recent observations. Perhaps data within the past 10 years: 2010-2019.

Future Study

As stated earlier, our models forecasts tended to overestimate the production of milk when compared to the actual data. Over estimation is a weakness of our model, but also our model does not account for significant world events such as possible disease outbreak among cows. To remedy this, we need to implement model fitting methods that incorporate more variables to account for these possible influences. Also our model was fitted using data from 1962 to 1975; if want to use our model to predict forecasts for say 2020, 2021, . . . we would probably need to expand the dataset we are using to include observations from 1976-2019. It would also be interesting to see how our model fares for predicting milk production in other countries, not just the U.S. This would require collecting more data from third party studies and adding it our current dataset. Overall, global population size is continuing to rise rapidly and food products, such as milk, will need to be produced in mass quantities. Therefore, the forecasting of milk production will remain relevant and important to meet the needs of society.

References

1. Milk Production 1962-1975, <https://datamarket.com/data/list/?q=provider:tsdl> DataMarket

Appendix

R Code:

```
library(MASS)
library(MuMIn)
library(sarima)
library(astsa)
library(stats)
source("plot.roots.R")
milk <- read.csv("milk.csv")
colnames(milk)[2] <- "production"
full.prod<- milk$production[-168]
prod<- milk$production[-c(157:170)] #Subset data so that we leave out last 12 months of data
#That way we can access performance of our model later on
ts.plot(prod, main = "Plot of Original Data")
#Data is not stationary(clear upward trend) and seasonality component
#Variance looks to be constant
#Apply box cox even though variance looks to be constant

#decompose()
ts <- ts(prod, frequency = 12)
decompose <- decompose(ts, "multiplicative")

plot(as.ts(decompose$seasonal))
plot(as.ts(decompose$trend))
plot(as.ts(decompose$random))

plot(decompose)

#Box Cox
require(MASS)
bctrans <- boxcox(prod~as.numeric(1:length(prod)),
                  plotit=TRUE,
                  lambda=seq(-4,2,0.1))
lambda <- bctrans$x[which(bctrans$y==max(bctrans$y))]
#print(paste("Lambda is equal to:",lambda))
#lamda = .42 which is close to .5 so should try sqrt transformation
#zero is also in the interval so lets try log transformation as well

prod.sqrt <- sqrt(prod)
prod.log <- log(prod)

# Plot original data vs Box-Cox transformed data
par(mfrow = c(2, 2))
ts.plot(prod,main = "Plot Original Data",ylab = expression(X[t]))
ts.plot(prod.sqrt,main = "Plot Sqrt tranformed Data", ylab = expression(Y[t]))
#Does not seem to be much of a difference between sqrt transformed and orginal
#Try plotting orginal vs log transformed
ts.plot(prod,main = "Plot Original Data",ylab = expression(X[t]))
```

```

ts.plot(prod.log,main = "Plot Log tranformed Data", ylab = expression(Y[t]))
#Not much of a difference in this case either
#Decide to not pursue transformation further
#Just work with original non transformed data

par(mfrow = c(3,2))

#Plot original data
ts.plot(prod, main = "Plot of Original Data")
#Lets plot acf and pacf of milk production data
acf(prod, lag.max=60, main="ACF of Original Data")
#From acf can tell that data is defintely non stationary
pacf(prod ,lag.max=60, main="PACF of Original Data")
#pacf also suggests data is non stationary

#Differencing once at lag 1 to remove trend
milk.diff1 <- diff(prod, lag = 1)
#print(paste("The variance with no differencing is: ", var(prod)))
#print(paste("The variance after differencing is: ", var(milk.diff1)))#Variance goes down which is good
ts.plot(milk.diff1, main= "Plot of Differened Once Data") #No trend is evident now, seasonality is stil
acf(milk.diff1 ,lag.max=60, main="ACF of Differened Once Data")
#Seasonal component is apparent with ACF, notice pattern of spikes at 12, 24,...
pacf(milk.diff1 ,lag.max=60, main="PACF of Differened Once Data")
#Lots of spikes in PACF before 12

par(mfrow=c(2,2))
#Difference again at lag 12 to remove seasonality
milk.diff12 <- diff(milk.diff1, lag=12)
#print(paste("The variance after differencing once is: ", var(milk.diff1)))#Variance goes down which is

plot.ts(milk.diff1, main = "Plot of Differened Once Data")
plot.ts(milk.diff12, main = "Plot of Differened @ Lag 12 Data")
#print(paste("The variance after differencing twice is: ", var(milk.diff12)))#Variance goes down which
#Pattern of spikes is gone which suggests seasonality is removed
#Also notice that variance is significantly lower
acf(milk.diff12 ,lag.max=60, main="ACF of Differened @ Lag 12 Data")
#Acf suggests that data is stationary
#Possible cut off at 12, so Q= 1, cut off at 1 so q = 1
pacf(milk.diff12 ,lag.max=60, main="PACFof Differened @ Lag 12 Data")
#Pacf suggests that data is stationary
#Cut off at 12, ,24 ,36 so P = 1, 2 or 3
#Cut off at 1 so p = 1
#d = 1, D = 1, P = (1,2,3), Q=1, p=1, q = 1

#fit1: SARIMA (1,1,1) x (1,1,1)_{12}$
#Sarima (1,1,1) x (1,1,1) 12
#Sarima (1,1,1) x (2,1,1) 12
#Sarima (1,1,1) x (3,1,1) 12

fit1 <- arima(x=prod, order=c(1,1,1),seasonal=list(order=c(1,1,1),period=12),method="ML")
fit1 #AIC = 991.93

```

```

##fit2: SARIMA (1,1,1) x (2,1,1)_{12}$
fit2 <- arima(x=prod, order=c(1,1,1),seasonal=list(order=c(2,1,1),period=12),method="ML")
fit2 #AIC = 992.83

##fit3: SARIMA (1,1,1) x (3,1,1)_{12}$
fit3 <- arima(x=prod, order=c(1,1,1),seasonal=list(order=c(3,1,1),period=12),method="ML")
fit3 #AIC = 994.87

# Model 1 has lowest AIC value of 991.93

#Before continuing make sure to check causality and invertibility of the models
par(mfrow=c(3,2))

plot.roots(NULL,polyroot(c(1,-.0937,-0.0468)),main="roots of AR part - fit1")
plot.roots(NULL,polyroot(c(1,-.1723,-.5845)),main="roots of MA part - fit1")
#Roots are outside unit circle. Process is causal and invertible.

plot.roots(NULL,polyroot(c(1,-.0718,-0.1612,-.1317)),main="roots of AR part - fit2")
plot.roots(NULL,polyroot(c(1,-.1883,-0.4741)),main="roots of MA part - fit2")
#Roots are outside unit circle. Process is causal and invertible.

plot.roots(NULL,polyroot(c(1,-.1258,-0.7035,-.4769,-.3451)),main="roots of AR part - fit3")
plot.roots(NULL,polyroot(c(1,-.0970,0.055)),main="roots of MA part - fit3")
#Roots are outside unit circle. Process is causal and invertible.

aic1<-AICc(arima(prod, order = c(1,1,1), seasonal=list(order=c(1,1,1), period=12), method="ML"))
#print(paste("The AICc value for model 1 is : ", aic1))
#Model 1 AICc value is 992.1965
aic2<-AICc(arima(prod, order = c(1,1,1), seasonal=list(order=c(2,1,1), period=12), method="ML"))
#print(paste("The AICc value for model 2 is : ", aic2))
#Model 2 AICc value is 993.2303
aic3<-AICc(arima(prod, order = c(1,1,1), seasonal=list(order=c(3,1,1), period=12), method="ML"))
#print(paste("The AICc value for model 3 is : ", aic3))
#Model 3 AICc value is 995.4352

#Model 1 had the lowest AIC value, and also the lowest AICc value
#Will go with model 1

### Analyze residuals for our model
par(mfrow= c(2,2))
ts.plot(residuals(fit1), main = "Plot of Final Model Residuals")

hist(residuals(fit1), main = "Histogram of Final Model Residuals") #Looks normally distributed.
#Slight skew on histogram towards the right.
#However for the most part looks normally distributed
qqnorm(residuals(fit1), main = "QQ-Plot of Final Model Residuals")
qqline(residuals(fit1),col="blue")
#The qq plot looks to be a little off especially around the right tail of the plot.
#However it looks like these are just outliers considering it is only a couple points.

shapiro.test(residuals(fit1))
#Shapiro Wilk test p-value is less than .05 which indicates a rejection of assumption of normality

```



```

#However this is expected considering our histogram and qq-plot were a bit skewed

####Diagnostic Checking
#lag = h = sqrt(n), n = 156, sqrt(156) = approx 12, p = 1, q = 1, fitdf = 1+1 = 2
Box.test(residuals(fit1), lag = 12, type = c("Box-Pierce"), fitdf = 2)
#Box.test(residuals(fit1), lag = 12, type = c("Ljung-Box"), fitdf = 2)
#Box.test(residuals(fit1), lag = 12, type = c("Ljung-Box"), fitdf = 0)
#P- values for all the box tests are greater than .05
#This indicates that autocorrelations are within the lags 1-11.
#Also indicates that residuals are independent

par(mfrow= c(2,2))
acf(residuals(fit1), main= "ACF of Final Model Residuals")
pacf(residuals(fit1), main = "PACF of Final Model Residuals")
acf((residuals(fit1))^2, main= "ACF of Squared Final Model Residuals")
pacf((residuals(fit1))^2, main = "PACF of Squared Final Model Residuals")

####Forecasting
par(mfrow=c(1,1))
fitfinal<- arima(x=prod, order=c(1,1,1),seasonal=list(order=c(1,1,1),period=12),
                method =c("ML"), xreg=1:length(prod))
pred.prod<- predict(fitfinal, n.ahead = 12,
                   newxreg=(length(prod)+1) : length(prod)+12)
U.tr= pred.prod$pred + 2*pred.prod$se # upper bound for the C.I. for data
L.tr= pred.prod$pred - 2*pred.prod$se #Lower bound
#Forecasted data for months of January 1975: December 1975
ts.plot(prod,xlim=c(0, 170), ylim=c(550,1000),
        ylab = "Monthly Milk Production",
        main = "Forecasted Values for Milk Production")
lines(U.tr, col="red", lty="dashed")
lines(L.tr, col="red", lty="dashed")
points((length(prod)+1):(length(prod)+12), pred.prod$pred, col="blue")

#Compare forecasted data to actual data for months of January to Decemeber 1975
ts.plot(full.prod,xlim=c(0, 170), ylim=c(550,1000),
        ylab = "Monthly Milk Production",
        main = "Forecasted Values for Milk Production vs Actual Values")
lines(U.tr, col="red", lty="dashed")
lines(L.tr, col="red", lty="dashed")
points((length(prod)+1):(length(prod)+12), pred.prod$pred, col="blue")
#Looks like our forecasted data slightly overestimated production
#For the most part follows our actual data

```