# Midterm Exam 2 Report

*Submitted By:*
Sayantan Majumdar
smxnv@mst.edu

*Course Instructor:*
Dr. Akim Adekpedjou
akima@mst.edu

November 21, 2019

# Contents

# List of Figures

# 1 Problem 1

The monthly values of the average hourly for US apparel and textile workers for July 1981 to July 1987 are the file wages of the library TSA.

(a) Plot the series and write a few lines on what you observe.

(b) Fit a linear trend model using least squares. Give the plot of the linear trend and superpose it with that of the data. Give the estimated regression equation.

(c) Plot the standardized residuals from the linear regression versus time. Comments?

(d) Fit a quadratic time trend model using least squares. Give the plot of quadratic trend and superpose it Give the estimated regression equation. Plot the residuals and comment on any patterns.

(e) Perform a diagnostic check of the residuals. Comments.

(f) Plot the autocorrelation function for the standardized residuals from the quadratic regression.

(g) Investigate the normality of the standardized residuals from the quadratic regression. Comment.

## 1.1 R Code

```
1  ### STAT 5814 MIDTERM 2/PROBLEM 1
2  ### AUTHOR: SAYANTAN MAJUMDAR
3  ### EMAIL: smxnv@mst.edu
4  ### SNO: 12566087
5
6  library(TSA)
7  library(snpar)
8
9  #get data
10 data(wages)
11 par(mfrow = c(1, 1))
12 plot(wages, type='o', ylab='Wages (USD/hr)', xlab='Time (Months)', main='
      Monthly Values of the Average Hourly Wages')
13
14 #linear model
15 wages.lm = lm(wages~time(wages))
16 summary(wages.lm)
17 model = ts(wages.lm$fitted.values, frequency=12, start=c(1981, 7), end=c
      (1987, 6))
18 lines(model, col='red')
19 plot(y=rstandard(wages.lm), x=as.vector(time(wages)), type = 'o', ylab = '
      Standardized Residuals', xlab = 'Time (Months)', main='Standardized
      Residual Plot')
20 par(mfrow = c(2, 2))
21 plot(wages.lm)
```

```
22
23  #Quadratic model trend
24  wages.qm = lm(wages ~ time(wages) + I(time(wages)^2))
25  summary(wages.qm)
26  model = ts(wages.qm$fitted.values, frequency=12, start=c(1981, 7), end=c
        (1987, 6))
27  par(mfrow=c(2, 2))
28  lines(model, col='red')
29  plot(y=rstandard(wages.qm), x=as.vector(time(wages)), type = 'o', ylab = '
        Standardized Residuals', xlab = 'Time (Months)', main='Standardized
        Residual Plot')
30  par(mfrow = c(2, 2))
31  plot(wages.qm)
32
33  #Residual Diagnostics
34  par(mfrow=c(1, 3))
35  hist(wages.qm$residuals, main="Residuals Histogram", xlab='Residuals')
36  qqnorm(wages.qm$residuals, main="QQ Plot for Residuals")
37  qqline(wages.qm$residuals, col="Red")
38  boxplot(wages.qm$residuals, main='Residual Boxplot', ylab='Residuals')
39
40  #ACF
41  acf(wages.qm$residuals, main='ACF of the Residuals')
42
43  #Residual Normality
44  shapiro.test(wages.qm$residuals)
45  runs.test(wages.qm$residuals, exact=TRUE)
```

## 1.2   Results

**Monthly Values of the Average Hourly Wages**

(a)



Figure 1.1: Time series plot of the wages data. This shows that the wages are mostly increasing and only deviating from this trend in small proportions.

(b) Linear trend model fit summary:

```
Call:
lm(formula = wages ~ time(wages))

Residuals:
     Min       1Q   Median       3Q      Max
-0.23828 -0.04981  0.01942  0.05845  0.13136

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) -5.490e+02  1.115e+01  -49.24   <2e-16 ***
time(wages)  2.811e-01  5.618e-03   50.03   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.08257 on 70 degrees of freedom
Multiple R-squared:  0.9728,Adjusted R-squared:  0.9724
F-statistic:  2503 on 1 and 70 DF,  p-value: < 2.2e-16
```

The $R^2$ value of this fit is very good ($\approx 0.97$) which signifies that the linear fit is working well. The estimated regression equation is given by Eq (1).

$$\widehat{X_t} = -549 + 0.2811X_t + \varepsilon \tag{1}$$

where $\varepsilon \sim \mathcal{N}(0, 0.08257)$



Figure 1.2: Wages plot showing the linear trend in red.

(c)



Figure 1.3: Standardized residuals from the linear regression versus time. Since the majority of the standardized residuals are within the [-2, 1] interval, so the model residuals strongly follow a normal distribution. Hence, we have a good fit.

(d) Quadratic time trend model summary:

```
Call:
lm(formula = wages ~ time(wages) + I(time(wages)^2))

Residuals:
     Min        1Q     Median        3Q        Max
-0.148318 -0.041440  0.001563  0.050089  0.139839

Coefficients:
                   Estimate Std. Error t value Pr(>|t|)
(Intercept)      -8.495e+04  1.019e+04  -8.336 4.87e-12 ***
time(wages)       8.534e+01  1.027e+01   8.309 5.44e-12 ***
I(time(wages)^2) -2.143e-02  2.588e-03  -8.282 6.10e-12 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.05889 on 69 degrees of freedom
Multiple R-squared:  0.9864,Adjusted R-squared:  0.986
F-statistic:  2494 on 2 and 69 DF,  p-value: < 2.2e-16
```

The $R^2$ value of this fit is even better ($\approx 0.98$) when compared to the linear fit earlier. This implies that the quadratic trend is more appropriate. The estimated regression equation is given by Eq (2).

$$\widehat{X_t} = -84950 + 85.34X_t - 0.02143X_t^2 + \varepsilon \tag{2}$$

where $\varepsilon \sim \mathcal{N}(0, 0.05889)$

**Monthly Values of the Average Hourly Wages**



Figure 1.4: Wages plot showing the quadratic fit in red.

**Standardized Residual Plot**



Figure 1.5: Standardized residual plot for the quadratic fit. Here, we observe that these most of these residuals are in the $\pm 2$ region. So, the quadratic fit is better than the linear one.

(e)



Figure 1.6: Residual diagnostics of the quadratic model showing residual histogram, qq-plot, and boxplot. These plots show that the residuals closely follow a normal distribution.

(f)



Figure 1.7: Autocorrelation function (ACF) for the standardized residuals from the quadratic regression. This shows that the ACF quickly degrades to zero as the lag increases.

(g) Residual normality checks:

```
Shapiro-Wilk normality test
data:  wages.qm$residuals
W = 0.98856, p-value = 0.7622

Exact runs test
data:  wages.qm$residuals
```

```
Runs = 15, p-value = 1.284e-07
alternative hypothesis: two.sided
```

The Shapiro-Wilk test shows that the residuals are normally distributed (p-value >
0.05 and W is close to 1). However, the runs test suggests that the order of the residuals
is not random.

# 2   Problem 2

A data set of 57 consecutive measurements from a machine tool are in the TSA package.

(a) Estimate the parameters of a (mean-centered) AR(1) model for this series. Use the
least squares method and maximum likelihood, and report the estimated parameters
from each of these methods. Comment on any similarities and differences. Give the
confidence intervals of your parameters.

(b) Estimate the parameters of a (mean-centered) AR(2) model for this series. Use the
least squares method and maximum likelihood, and report the estimated parameters
from each of these methods. Comment on any similarities and differences.

(c) Derive the confidence intervals for the parameters $\phi_1$ and $\phi_2$. Does the confidence
interval for $\phi_2$ suggests it should be included in the model?

(d) Compare the results of the maximum likelihood fits from parts (a) and (b). Which
model do you believe is preferable? Briefly explain your answer.

(e) Perform a diagnostic check of the residuals. Comments.

(f) Compare the two models using AIC, BIC. Which one do you prefer? Does any of the
quantities suggests one model is better than the other? If so, how much do you gain
in errors reduction if any. Is that in line with your answer in (d)?

## 2.1   R Code

```
1  ### STAT 5814 MIDTERM 2/PROBLEM 2
2  ### AUTHOR: SAYANTAN MAJUMDAR
3  ### EMAIL: smxnv@mst.edu
4  ### SNO: 12566087
5
6  library(TSA)
7  library(snpar)
8  library(stats)
9  library(forecast)
10
11 data("deere3")
12 par(mfrow = c(1, 1))
13 plot(deere3, type='o', ylab='Measurements', xlab='Time', main='Deere3 Data
       Plot')
```

```r
14
15  # Conditional least square estimate for AR(1) model
16  css_model_ar1 = arima(deere3, order=c(1, 0, 0), method='CSS')
17  css_model_ar1_coef = coef(css_model_ar1)
18  print(css_model_ar1_coef)
19  confint(css_model_ar1)
20
21
22  # Maximum Likelihood estimate for AR(1) model
23  ml_model_ar1 = arima(deere3, order=c(1, 0, 0), method='ML')
24  ml_model_ar1_coef = coef(ml_model_ar1)
25  print(ml_model_ar1_coef)
26  confint(ml_model_ar1)
27
28  # Conditional least square estimate for AR(2) model
29  css_model_ar2 = arima(deere3, order=c(2, 0, 0), method='CSS')
30  css_model_ar2_coef = coef(css_model_ar2)
31  print(css_model_ar2_coef)
32  confint(css_model_ar2)
33
34  # Maximum Likelihood estimate for AR(2) model
35  ml_model_ar2 = arima(deere3, order=c(2, 0, 0), method='ML')
36  ml_model_ar2_coef = coef(ml_model_ar2)
37  print(ml_model_ar2_coef)
38  confint(ml_model_ar2)
39
40  # Residual Diagnostics (manually repeated for all the models)
41  par(mfrow=c(3, 2))
42  model = ml_model_ar2
43  st_residuals = rstandard(model)
44  plot(y=st_residuals, x=as.vector(time(deere3)), type = 'o', ylab = '
        Standardized Residuals', xlab = 'Time', main='Standardized Residual
        Plot')
45  hist(st_residuals, main="Residuals Histogram", xlab='Residuals')
46  qqnorm(st_residuals, main="QQ Plot for Residuals")
47  qqline(st_residuals, col="Red")
48  boxplot(st_residuals, main='Residual Boxplot', ylab='Residuals')
49  acf(st_residuals, main='ACF Plot')
50  #Residual Normality
51  shapiro.test(st_residuals)
52  runs.test(st_residuals, exact=TRUE)
53
54  #AIC and BIC tests
55  AIC(model)
56  BIC(model)
57
58  #Confirmation
59  auto_arima = auto.arima(deere3)
60  auto_arima
```

## 2.2 Results



Figure 2.1: Time series plot of the Deere3 data.

(a) Least squares method for AR(1) parameter estimation:

```
> print(css_model_ar1_coef)
        ar1    intercept
  0.5332044 160.0797248
> confint(css_model_ar1)
                  2.5 %        97.5 %
ar1           0.3132122    0.7531966
intercept -646.0111173 966.1705670
```

Maximum Likelihood method for AR(1) parameter estimation:

```
> print(ml_model_ar1_coef)
        ar1    intercept
  0.5255778 124.3524257
> confint(ml_model_ar1)
                  2.5 %        97.5 %
ar1           0.3084104    0.7427452
intercept -648.3281895 897.0330409
```

The estimated $\hat{\phi}$ from both the least squares and maximum likelihood methods are very similar. However, the intercept ($\hat{\mu}$) is significantly different. In addition, the confidence intervals for $\hat{\phi}$ in both cases are similar and $0 \notin \text{CI}(\hat{\phi})$.

(b) Least squares method for AR(2) parameter estimation:

```
> print(css_model_ar2_coef)
        ar1           ar2     intercept
5.245848e-01 7.938728e-03 2.011876e+02
```

Maximum Likelihood method for AR(2) parameter estimation:

```
> print(ml_model_ar2_coef)
        ar1           ar2     intercept
  0.52111096    0.00830321 123.24182209
```

Like before, the estimated $\phi_1$ and $\phi_2$ in both cases are pretty similar, whereas, the intercepts are significantly different.

(c) Confidence intervals for $\phi_1$ and $\phi_2$:

```
Least Squares Method
> confint(css_model_ar2)
                  2.5 %        97.5 %
ar1           0.2660503     0.7831192
ar2          -0.2509164     0.2667939
intercept -607.0745845  1009.4498459


Maximum Likelihood Method
> confint(ml_model_ar2)
                  2.5 %        97.5 %
ar1           0.2642641     0.7779579
ar2          -0.2493401     0.2659465
intercept -656.0388502   902.5224944
```

Since $0 \in \text{CI}(\phi_2)$, therefore, $\phi_2 \approx 0$. So $\phi_2$ should not be included in the model.

(d) The maximum likelihood estimates of the intercepts for the AR(1) and AR(2) models are pretty similar. Also, the confidence interval of the intercept in case of AR(1) is a subset of that of AR(2). But, since our target is develop a parsimonious model and $0 \in \text{CI}(\phi_2)$ for the maximum likelihood method as well, so, AR(1) is the preferable model.

(e) Residual Diagnostics: The residual normality test results are given below.

```
Least Squares Estimation for AR(1):


Shapiro-Wilk normality test
data:  st_residuals
W = 0.98297, p-value = 0.6
```

Figure 2.2: Residual diagnostic plots for least squares estimation for AR(1). These plots shows that the residuals tend to closely follow a normal distribution.

```
Exact runs test
data:  st_residuals
Runs = 29, p-value = 1


Maximum Likelihood Estimation for AR(1):


Shapiro-Wilk normality test
data:  st_residuals
W = 0.98261, p-value = 0.5827


Exact runs test
data:  st_residuals
Runs = 29, p-value = 1


Least Squares Estimation for AR(2):


Shapiro-Wilk normality test
```

Figure 2.3: Residual diagnostic plots for maximum likelihood estimation for AR(1). Similar observations like that of Fig 2.2 can be made.

```
data:  st_residuals
W = 0.9809, p-value = 0.5028

Exact runs test
data:  st_residuals
Runs = 29, p-value = 1

Maximum Likelihood Estimation for AR(2):

Shapiro-Wilk normality test
data:  st_residuals
W = 0.98294, p-value = 0.599

Exact runs test
data:  st_residuals
Runs = 29, p-value = 1
```

Figure 2.4: Residual diagnostic plots for least squares estimation for AR(2). Like before, these also closely follow a normal distribution but less than AR(1) residuals.

The normality tests show that the residuals from the AR(1) tend to follow the normal distribution more closely than those of AR(2).

(f) AIC and BIC tests for the maximum likelihood method:

```
AR(1):
> AIC(model)
[1] 997.0189
> BIC(model)
[1] 1003.148

AR(2):
> AIC(model)
[1] 999.0148
> BIC(model)
[1] 1007.187
```

As observed, we obtain a reduction of $\approx 2$ and $\approx 4$ in the AIC and BIC values

Figure 2.5: Residual diagnostic plots for maximum likelihood estimation for AR(2). Similar observations like that of Fig 2.4 can be made.

respectively for the AR(1) model. Although, these are not significant reductions, our aim is always to develop a parsimonious model which aligns with the results obtained in (d). Therefore, AR(1) is the preferred model which is further confirmed using the auto.arima function in R.

```
> auto_arima
Series: deere3
ARIMA(1,0,0) with zero mean

Coefficients:
          ar1
       0.5291
s.e.   0.1103


sigma^2 estimated as 2109761:  log likelihood=-495.56
AIC=995.12    AICc=995.34    BIC=999.2
```

# 3 Problem 3

Are sales trends for two industrial products related? The two time series objects described below contain Sales of chemicals and allied products; and Sales of motor vehicles and parts, in the U.S. for each month from Jan. 1971 to Dec. 1991. File is petr.txt.

You should analyze the data in the "chemicals" time series and the "vehicles" time series and write a report to address such questions as:

(i) What trend model(s) best capture the trends in sales of chemicals over time?

(ii) What trend model(s) best capture the trends in sales of vehicles over time? Once the trend has been accounted for, what can you say about the behavior of the detrended data, for both models? Can either the original time series or the detrended series be described using any common models?

(iii) For the various models you tried, assess their fit. Can any transformations improve the fit? Is there any apparent association between chemical sales and vehicle sales over time? If so, describe the association.

(iv) Make conclusions that relate to how the sales amounts (for both chemicals and vehicles) change, both long-term over the observed period of years, and in terms of patterns of month- to-month variation.

## 3.1 R Code

```
1 ### STAT 5814 MIDTERM 2/PROBLEM 3
2 ### AUTHOR: SAYANTAN MAJUMDAR
3 ### EMAIL: smxnv@mst.edu
4 ### SNO: 12566087
5
6 library(TSA)
7 library(snpar)
8 library(stats)
9 library(forecast)
10 library(tseries)
11 library(fUnitRoots)
12
13 data = read.table('Data/ExamData/petr.txt', header=T)
14 chemicals.data = ts(data$Chemicals, frequency=12, start=c(1971, 1), end=c
      (1991, 12))
15 vehicles.data = ts(data$Vehicles, frequency=12, start=c(1971, 1), end=c
      (1991, 12))
16 plot(chemicals.data, main='Time Series Plot of Chemicals Data')
17 plot(vehicles.data, main='Time Series Plot of Vehicles Data')
18 chemicals.decompose = decompose(chemicals.data, type="mult")
19 plot(chemicals.decompose)
20 vehicles.decompose = decompose(vehicles.data, type="mult")
21 plot(vehicles.decompose)
22
```

```
23  # Stationarity Check
24  adf.test(chemicals.data)
25  adf.test(vehicles.data)
26
27  par(mfrow=c(2, 2))
28  # Differencing to remove linear trend
29  chemicals.diff= diff(chemicals.data)
30  plot(chemicals.diff, main='Differencing=1: Chemicals Data')
31
32  # Differencing at lag 12 (monthly data) to remove seasonality
33  chemicals.diff2 = diff(chemicals.diff, lag=12)
34  plot(chemicals.diff2, main='Lag=12: Chemicals Data')
35  # Stationarity Check
36  adf.test(chemicals.diff2)
37
38  # Differencing to remove linear trend
39  vehicles.diff= diff(vehicles.data)
40  plot(vehicles.diff, main='Differencing=1: Vehicles Data')
41
42  # Differencing at lag 12 (monthly data) to remove seasonality
43  vehicles.diff2 = diff(vehicles.diff, lag=12)
44  plot(vehicles.diff2, main='Lag=12: Vehicles Data')
45  # Stationarity Check
46  adf.test(vehicles.diff2)
47
48  par(mfrow=c(2, 1))
49  # Linear Trend and Harmonic (Seasonal) Trend for Chemicals Data
50  har = harmonic(chemicals.data, 1)
51  chemicals.fit = lm(chemicals.data ~ har + time(chemicals.data))
52  chemicals.residuals = rstudent(chemicals.fit)
53  summary(chemicals.fit)
54  chemicals.model = ts(chemicals.fit$fitted.values, frequency=12, start=c
        (1971, 1), end=c(1991, 12))
55
56  plot(chemicals.data, main='TS Plot: Chemicals Data')
57  lines(chemicals.model, col='red')
58  legend(1975, 25, legend=c("Actual", "Fitted"), col=c("Black", "Red"), lty
        =1:1, cex=0.5)
59
60  # Checking Transformation for Chemicals Data
61  chemicals.bc.lambda = BoxCox.lambda(chemicals.data)
62  chemicals.bc.tr = BoxCox(chemicals.data, lambda=chemicals.bc.lambda)
63  har = harmonic(chemicals.bc.tr, 1)
64  chemicals.bc.fit = lm(chemicals.bc.tr ~ (har + time(chemicals.bc.tr)))
65  chemicals.bc.residuals = rstudent(chemicals.bc.fit)
66  summary(chemicals.bc.fit)
67  chemicals.bc.model = ts(chemicals.bc.fit$fitted.values, frequency=12,
        start=c(1971, 1), end=c(1991, 12))
68
69  plot(chemicals.bc.tr, main='TS Plot: Transformed Chemicals Data')
70  lines(chemicals.bc.model, col='red')
71  legend(1975, 4.2, legend=c("Actual", "Fitted"), col=c("Black", "Red"), lty
        =1:1, cex=0.5)
72
```

```
73 # Linear Trend and Harmonic ( Seasonal ) Trend for Vehicles Data
74 har = harmonic ( vehicles . data , 1)
75 vehicles . fit = lm ( vehicles . data ~ har + time ( vehicles . data ))
76 vehicles . residuals = rstudent ( vehicles . fit )
77 summary ( vehicles . fit )
78 vehicles . model = ts ( vehicles . fit$fitted . values , frequency =12 , start =c
      (1971 , 1) , end =c (1991 , 12))
79
80 plot ( vehicles . data , main = 'TS Plot : Vehicles Data ')
81 lines ( vehicles . model , col = 'red ')
82 legend (1975 , 20 , legend =c ( " Actual " , " Fitted " ) , col =c ( " Black " , " Red " ) , lty
      =1:1 , cex =0.5)
83
84 # Checking Transformation for Vehicles Data
85 vehicles . bc . lambda = BoxCox . lambda ( vehicles . data )
86 vehicles . bc . tr = BoxCox ( vehicles . data , lambda = vehicles . bc . lambda )
87 har = harmonic ( vehicles . bc . tr , 1)
88 vehicles . bc . fit = lm ( vehicles . bc . tr ~ ( har + time ( vehicles . bc . tr )))
89 vehicles . bc . residuals = rstudent ( vehicles . bc . fit )
90 summary ( vehicles . bc . fit )
91 vehicles . bc . model = ts ( vehicles . bc . fit$fitted . values , frequency =12 , start =
      c (1971 , 1) , end =c (1991 , 12))
92
93 plot ( vehicles . bc . tr , main = 'TS Plot : Transformed Vehicles Data ')
94 lines ( vehicles . bc . model , col = 'red ')
95 legend (1975 , 3 , legend =c ( " Actual " , " Fitted " ) , col =c ( " Black " , " Red " ) , lty
      =1:1 , cex =0.5)
96
97 # Differencing after BoxCox Chemicals Data
98 chemicals . bc . diff = diff ( chemicals . bc . tr )
99 plot ( chemicals . bc . diff , main = 'Differencing =1: Transformed Chemicals Data ')
100
101 # Differencing at lag 12 ( monthly data ) to remove seasonality
102 chemicals . bc . diff2 = diff ( chemicals . bc . diff , lag =12)
103 plot ( chemicals . bc . diff2 , main = 'Lag =12: Transformed Chemicals Data ')
104 # Stationarity Check
105 adf . test ( chemicals . bc . diff2 )
106
107 # Differencing after BoxCox Vehicles Data
108 vehicles . bc . diff = diff ( vehicles . bc . tr )
109 plot ( vehicles . bc . diff , main = 'Differencing =1: Transformed Vehicles Data ')
110
111 # Differencing at lag 12 ( monthly data ) to remove seasonality
112 vehicles . bc . diff2 = diff ( vehicles . bc . diff , lag =12)
113 plot ( vehicles . bc . diff2 , main = 'Lag =12: Transformed Vehicles Data ')
114 # Stationarity Check
115 adf . test ( vehicles . bc . diff2 )
116
117 par ( mfrow =c (2 , 3))
118 # ACF , PACF , and EACF for Detrended Chemicals Data
119 acf ( chemicals . diff2 , main = " ACF : Detrended Chemicals Data " )
120 pacf ( chemicals . diff2 , main = " PACF : Detrended Chemicals Data " )
121 eacf ( chemicals . diff2 )
122 plot ( armasubsets ( chemicals . diff2 , 5 , 5))
```

```
123
124 # ACF, PACF, and EACF for Detrended Vehicles Data
125 acf(vehicles.diff2, main="ACF: Detrended Vehicles Data")
126 pacf(vehicles.diff2, main="PACF: Detrended Chemicals Data")
127 eacf(vehicles.diff2)
128 plot(armasubsets(vehicles.diff2, 5, 5))
129
130 # ACF, PACF, and EACF for Stabilized Chemicals Data
131 acf(chemicals.bc.diff2, main="ACF: Stabilized Chemicals Data")
132 pacf(chemicals.bc.diff2, main="PACF: Stabilized Chemicals Data")
133 eacf(chemicals.bc.diff2)
134 plot(armasubsets(chemicals.bc.diff2, 5, 5))
135
136 # ACF, PACF, and EACF for Stabilized Vehicles Data
137 acf(vehicles.bc.diff2, main="ACF: Stabilized Vehicles Data")
138 pacf(vehicles.bc.diff2, main="PACF: Stabilized Chemicals Data")
139 eacf(vehicles.bc.diff2)
140 plot(armasubsets(vehicles.bc.diff2, 5, 5))
141
142 par(mfrow=c(2, 1))
143 # ARIMA Model Fits for Chemicals Data
144 chemicals.arima = auto.arima(chemicals.diff2)
145 chemicals.bc.arima = auto.arima(chemicals.diff2, lambda='auto')
146 chemicals.bc.diff2.arima = auto.arima(chemicals.bc.diff2)
147
148 # ARIMA Model Fits for Vehicles Data
149 vehicles.arima = auto.arima(vehicles.diff2)
150 vehicles.bc.arima = auto.arima(vehicles.diff2, lambda='auto')
151 vehicles.bc.diff2.arima = auto.arima(vehicles.bc.diff2)
152
153 # Residual Diagnostics (manually repeated for all the models)
154 par(mfrow=c(3, 2))
155 dataset = chemicals.diff2
156 model = chemicals.arima
157 # st_residuals = chemicals.residuals
158 st_residuals = rstandard(model)
159 plot(y=st_residuals, x=as.vector(time(dataset)), type = 'o', ylab = '
        Standardized Residuals', xlab = 'Time', main='Standardized Residual
        Plot')
160 hist(st_residuals, main="Residuals Histogram", xlab='Residuals')
161 qqnorm(st_residuals, main="QQ Plot for Residuals")
162 qqline(st_residuals, col="Red")
163 boxplot(st_residuals, main='Residual Boxplot', ylab='Residuals')
164 acf(st_residuals, main='ACF Plot')
165 tsdiag(model)
166 # Residual Normality
167 shapiro.test(st_residuals)
168 snpar::runs.test(st_residuals, exact=TRUE)
169 # AIC and BIC tests
170 AIC(model)
171 BIC(model)
172
173 # Correlation between Chemicals and Vehicles Data
174 ccf_cv = ccf(chemicals.data, vehicles.data)
```

## 3.2   Results

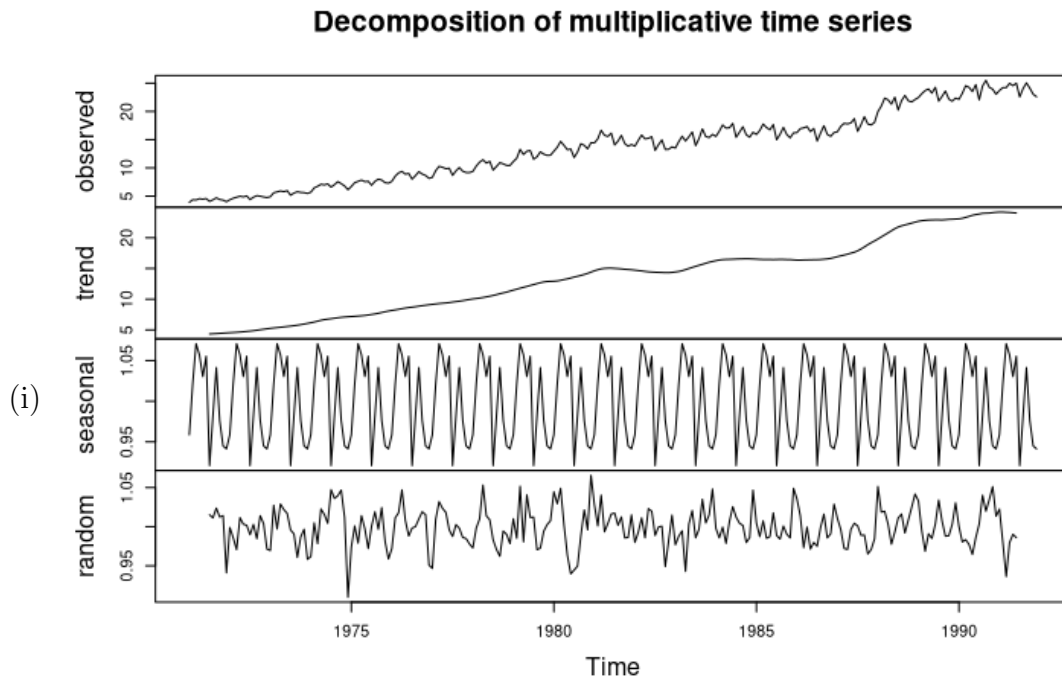**Decomposition of multiplicative time series**

(i)

Figure 3.1: Decomposition of the Chemicals data shows that both trend and seasonality are present.

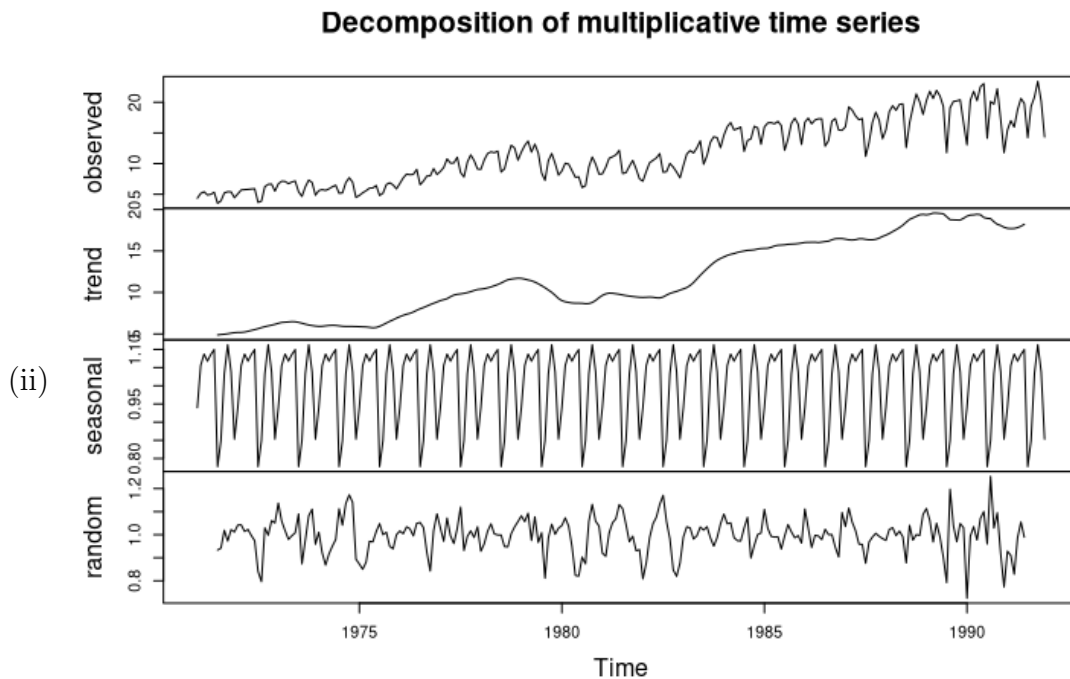**Decomposition of multiplicative time series**

(ii)

Figure 3.2: Decomposition of the Vehicles data.

Here, we again observe that the data exhibit both trend and seasonality. Nevertheless, in both the cases, we check their stationarity.

```
> adf.test(chemicals.data)


Augmented Dickey-Fuller Test
data:  chemicals.data
Dickey-Fuller = -3.0203, Lag order = 6, p-value = 0.1462
alternative hypothesis: stationary

> adf.test(vehicles.data)
Augmented Dickey-Fuller Test
data:  vehicles.data
Dickey-Fuller = -2.9112, Lag order = 6, p-value = 0.1921
alternative hypothesis: stationary
```

The augmented Dickey-Fuller tests for both the datasets fail to reject the null hypothesis, and hence, these time series data are not stationary. Therefore, we detrend (also remove seasonality) these, and the time series plots are shown in Fig 3.3.

The visual observation of the final detrended (Fig 3.3) data show that both the time series have become stationary. We further test this using the augmented Dickey-Fuller test which confirms that we have achieved stationarity.

```
> adf.test(chemicals.diff2)


Augmented Dickey-Fuller Test
data:  chemicals.diff2
Dickey-Fuller = -5.6096, Lag order = 6, p-value = 0.01
alternative hypothesis: stationary

> adf.test(vehicles.diff2)


Augmented Dickey-Fuller Test
data:  vehicles.diff2
Dickey-Fuller = -7.1672, Lag order = 6, p-value = 0.01
alternative hypothesis: stationary
```

Here, the original data can be described as a harmonic model with increasing trend. However, for the detrended data, we can calculate the ACF, PACF, and EACF to get an idea about the order of the ARIMA process.

(iii) Initally, we check the result of the harmonic model fit over the original data.

```
Call:
lm(formula = chemicals.data ~ har + time(chemicals.data))
```
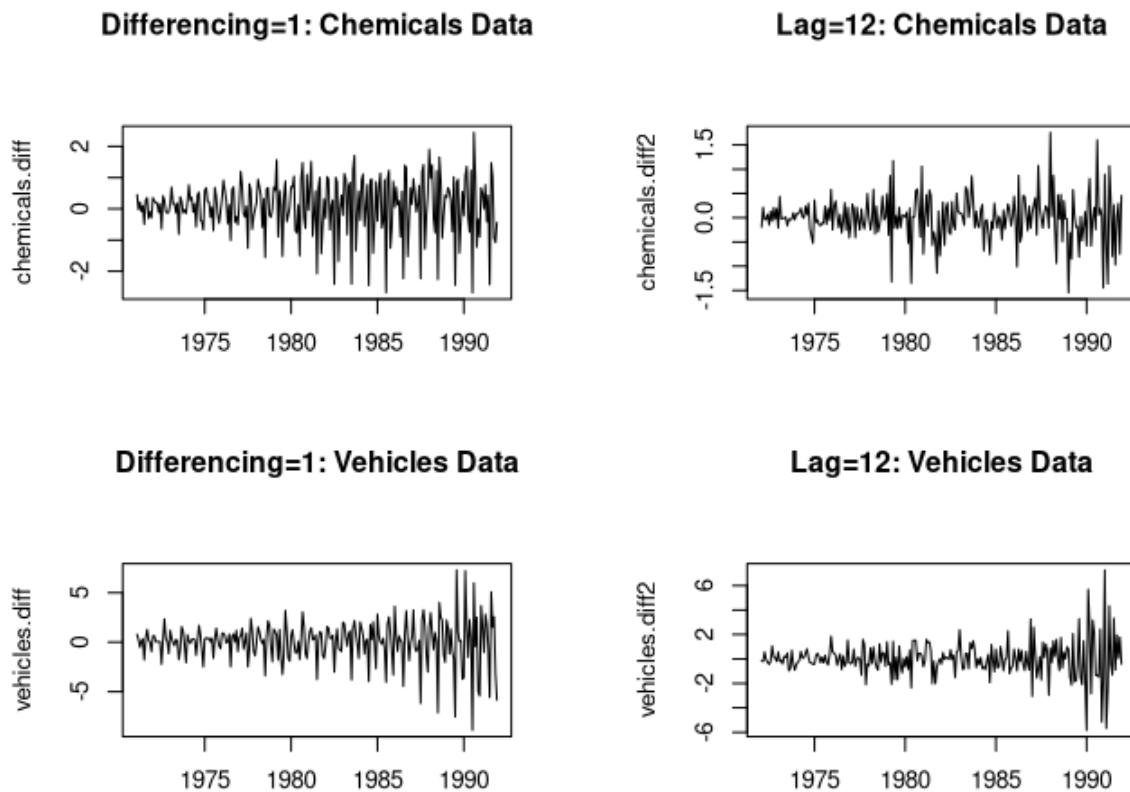
Figure 3.3: Time series plots of the detrended data. At first, the linear trend is removed using the first differences. Then the seasonality of the differenced data is removed by setting lag=12 as we deal with monthly data here.

```
Residuals:
    Min      1Q  Median      3Q     Max
-4.0447 -0.5332  0.0430  0.6677  3.0799


Coefficients:
                     Estimate Std. Error t value Pr(>|t|)
(Intercept)         -1.951e+03  2.365e+01 -82.478  < 2e-16 ***
harcos(2*pi*t)      -9.564e-02  1.023e-01  -0.935    0.351
harsin(2*pi*t)       5.420e-01  1.023e-01   5.297  2.6e-07 ***
time(chemicals.data) 9.913e-01  1.194e-02  83.058  < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1


Residual standard error: 1.148 on 248 degrees of freedom
Multiple R-squared:  0.9653,Adjusted R-squared:  0.9649
F-statistic:  2302 on 3 and 248 DF,  p-value: < 2.2e-16
```

```
Call:
lm(formula = vehicles.data ~ har + time(vehicles.data))

Residuals:
    Min      1Q  Median      3Q     Max
-7.0874 -1.3989  0.3886  1.5226  4.5814

Coefficients:
                     Estimate Std. Error t value Pr(>|t|)
(Intercept)        -1.499e+03  4.476e+01 -33.497  < 2e-16 ***
harcos(2*pi*t)      2.120e-01  1.935e-01   1.095  0.27449
harsin(2*pi*t)      5.378e-01  1.936e-01   2.777  0.00591 **
time(vehicles.data) 7.626e-01  2.259e-02  33.760  < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.172 on 248 degrees of freedom
Multiple R-squared:  0.8217,Adjusted R-squared:  0.8195
F-statistic: 380.9 on 3 and 248 DF,  p-value: < 2.2e-16
```

While the $R^2$ is good for both cases, the residual standard error is quite high along with the F-statistic. The model fits for both the original datasets are depicted in Fig 3.4. We now perform the residual diagnostics of the fit. Furthermore, we check the AIC
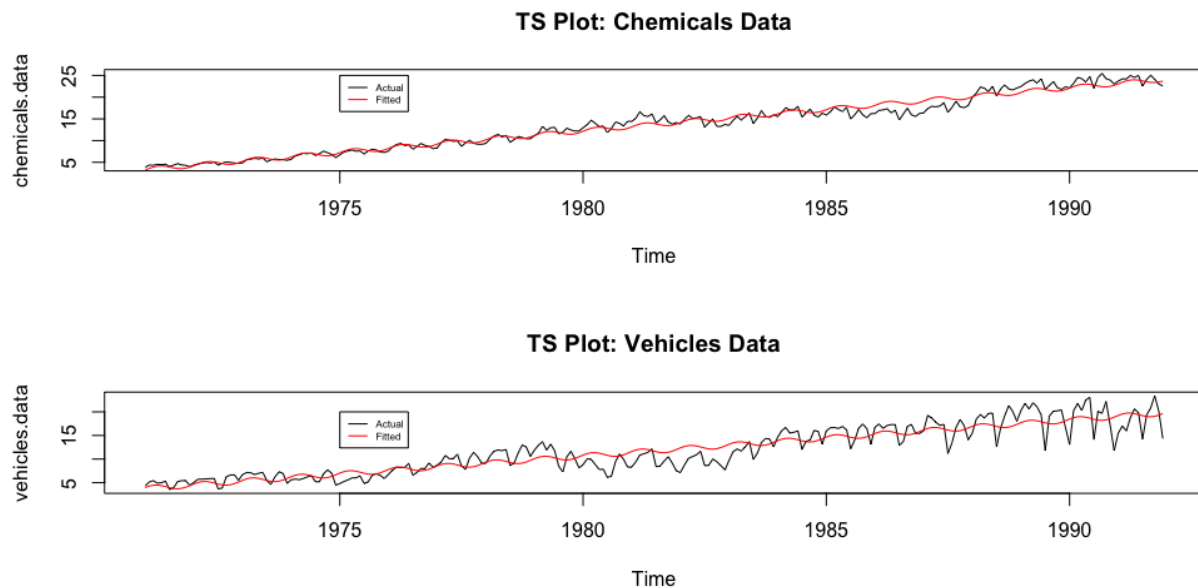


Figure 3.4: Fitting a harmonic model with linear trend over the original data. We see that this model somewhat matches the time series pattern but misses out a lot. Increasing the number of harmonic pairs better fits the model but could lead to possible overfitting and results in a non-parsimonious model.
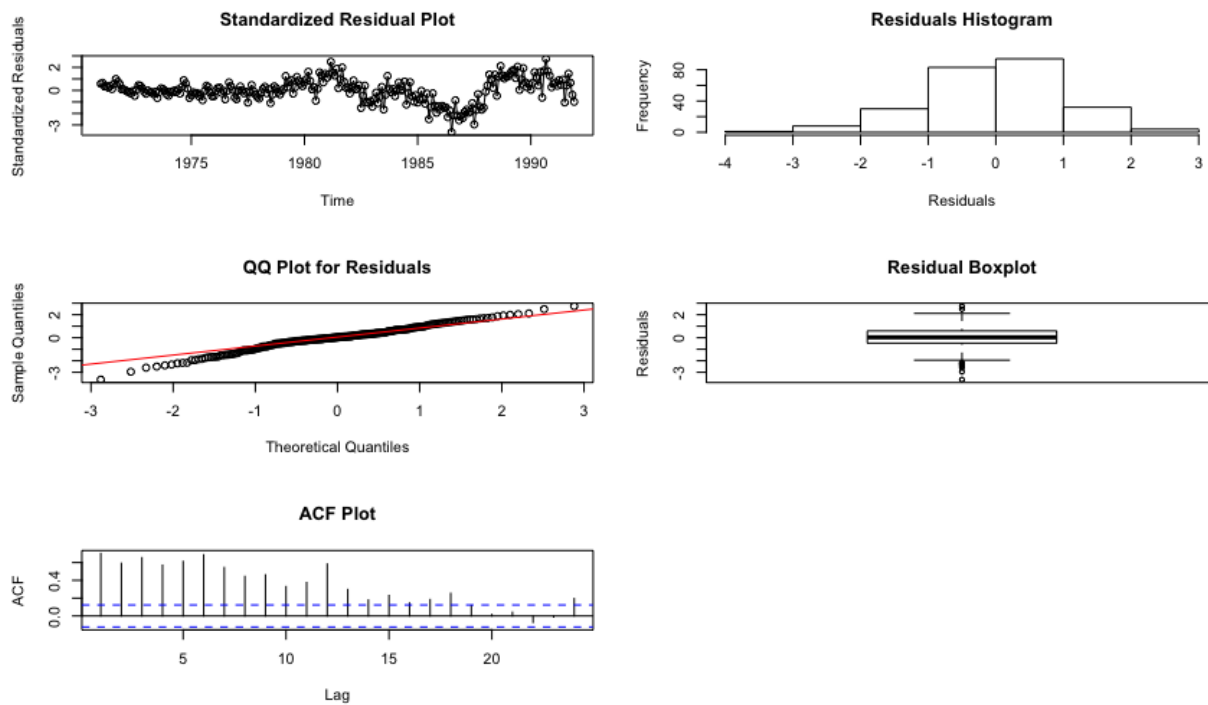
Figure 3.5: Residual diagnostics of the harmonic fit over the original Chemicals data. Here, we see the non-normality of the residuals.

and BIC values along with the p-values of the normality checks. Both the normality checks suggest that the residuals are from a non-normal distribution.

```
Original Chemicals Data

Shapiro-Wilk normality test
W = 0.98358, p-value = 0.005282

Exact runs test
Runs = 56, p-value < 2.2e-16
alternative hypothesis: two.sided

> AIC(chemicals.fit)
[1] 790.5797
> BIC(chemicals.fit)
[1] 808.2269

Original Vehicles Data

Shapiro-Wilk normality test
W = 0.98358, p-value = 0.005282

Exact runs test
```
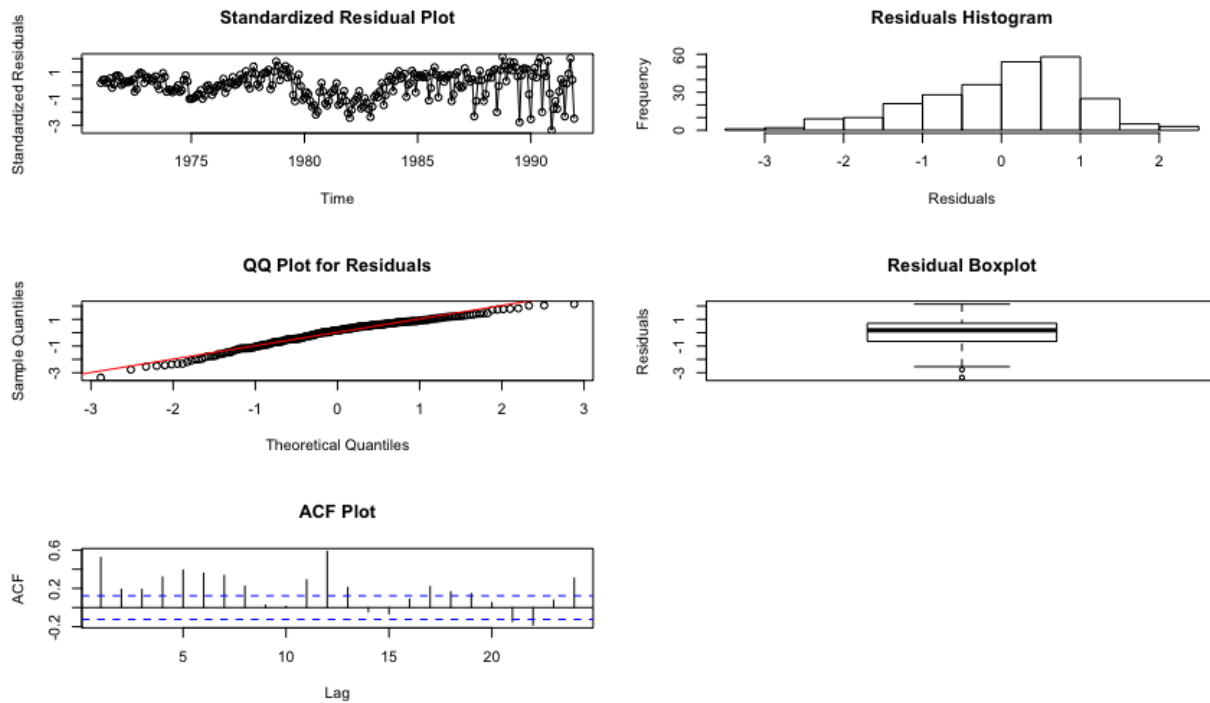
Figure 3.6: Residual diagnostics of the harmonic fit over the original Vehicles data. Here, we again see the non-normality of the residuals.

```
Runs = 56, p-value < 2.2e-16
alternative hypothesis: two.sided

> AIC(vehicles.fit)
[1] 1112.108
> BIC(vehicles.fit)
[1] 1129.755
```

Now, we observe the ACF, PACF, EACF, and ARMA Subset Plots for the detrended data. The EACF matrices for these data are given below.

```
> eacf(chemicals.diff2)
AR/MA
  0 1 2 3 4 5 6 7 8 9 10 11 12 13
0 x o x x o o o x o o o  x  o  o
1 o o x o o o o x o o o  x  o  o
2 x x o o o o o x o o o  x  o  x
3 x x x o o o o o o o o  x  o  x
4 x x o o o o o o o o o  x  o  x
5 x x o o o o o o o o o  x  x  o
6 x x o x x o o o o o o  x  o  o
7 x x x x o o o o o o o  x  o  o
```
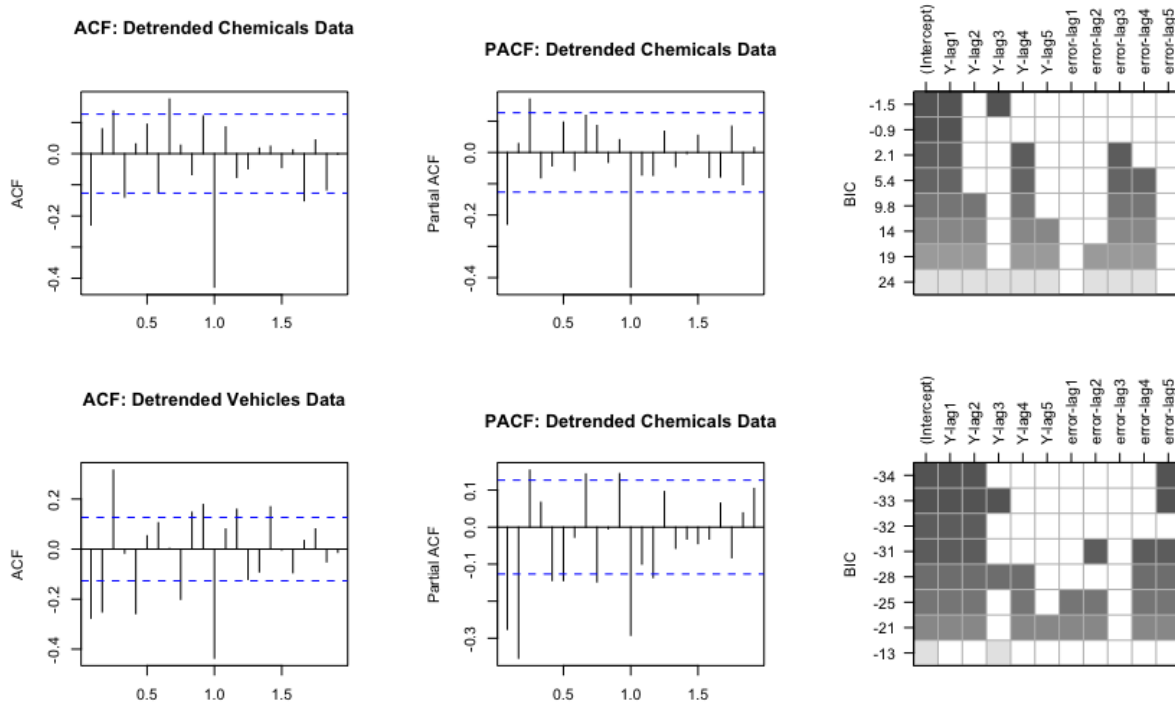
Figure 3.7: ACF, PACF, and ARMA subset plots for the detrended Chemicals and Vehicles data. These show that the detrended data could possible have both AR and MA parts.

```
> eacf(vehicles.diff2)
AR/MA
  0 1 2 3 4 5 6 7 8 9 10 11 12 13
0 x x x o x o o o x x x  x  o  x
1 x x x o x o o o x o o  x  x  x
2 x x o x x x o o o o o  o  x  o  x
3 x x x o x o o o o o o  x  x  x
4 x x o o o o o o o o o  x  x  x
5 x x x o x o o o o o o  x  x  o
6 x x o x x o o o o o o  x  x  o
7 x x o o o o o o o o o  x  x  o
```

Next, we fit the best ARIMA model using the auto.arima function in R and perform the necessary residual diagnostics (using tsdiag function). These are shown in Figures 3.8 and 3.9.

```
> chemicals.arima = auto.arima(chemicals.diff2)
> chemicals.arima
Series: chemicals.diff2
ARIMA(2,0,2)(0,0,1)[12] with zero mean

Coefficients:
```

```
          ar1      ar2      ma1      ma2     sma1
      -0.9289  -0.8776   0.7277   0.7503  -0.6089
s.e.   0.0644   0.0712   0.0901   0.0911   0.0583


sigma^2 estimated as 0.1409:  log likelihood=-104.97
AIC=221.95    AICc=222.31    BIC=242.8

> vehicles.arima = auto.arima(vehicles.diff2)
> vehicles.arima
Series: vehicles.diff2
ARIMA(2,0,2)(0,0,1)[12] with zero mean

Coefficients:
          ar1      ar2      ma1      ma2     sma1
      -0.4311  -0.7637   0.1581   0.5965  -0.5006
s.e.   0.0943   0.1128   0.1040   0.1725   0.0623


sigma^2 estimated as 1.259:  log likelihood=-366.3
AIC=744.61    AICc=744.97    BIC=765.47
```

Like before, we carry out the normality tests which align with our graphically inter-
preted results.

```
Detrended Chemicals Data


Shapiro-Wilk normality test
W = 0.95342, p-value = 5.817e-07

Exact runs test
Runs = 121, p-value = 0.8969
alternative hypothesis: two.sided

Detrended Vehicles Data


Shapiro-Wilk normality test
W = 0.97736, p-value = 0.0007176

Exact runs test
Runs = 109, p-value = 0.1537
alternative hypothesis: two.sided
```

Once these analyses are done, we apply BoxCox transformation (automatically chosen
$\lambda_C \approx 0.19$ and $\lambda_V \approx -0.03$ for Chemicals and Vehicles data respectively) to the original
data. We also do the same for the detrended data but in that case we first apply the
transformation and then detrend.

The model details of the harmonic fit over the transformed data are given below.
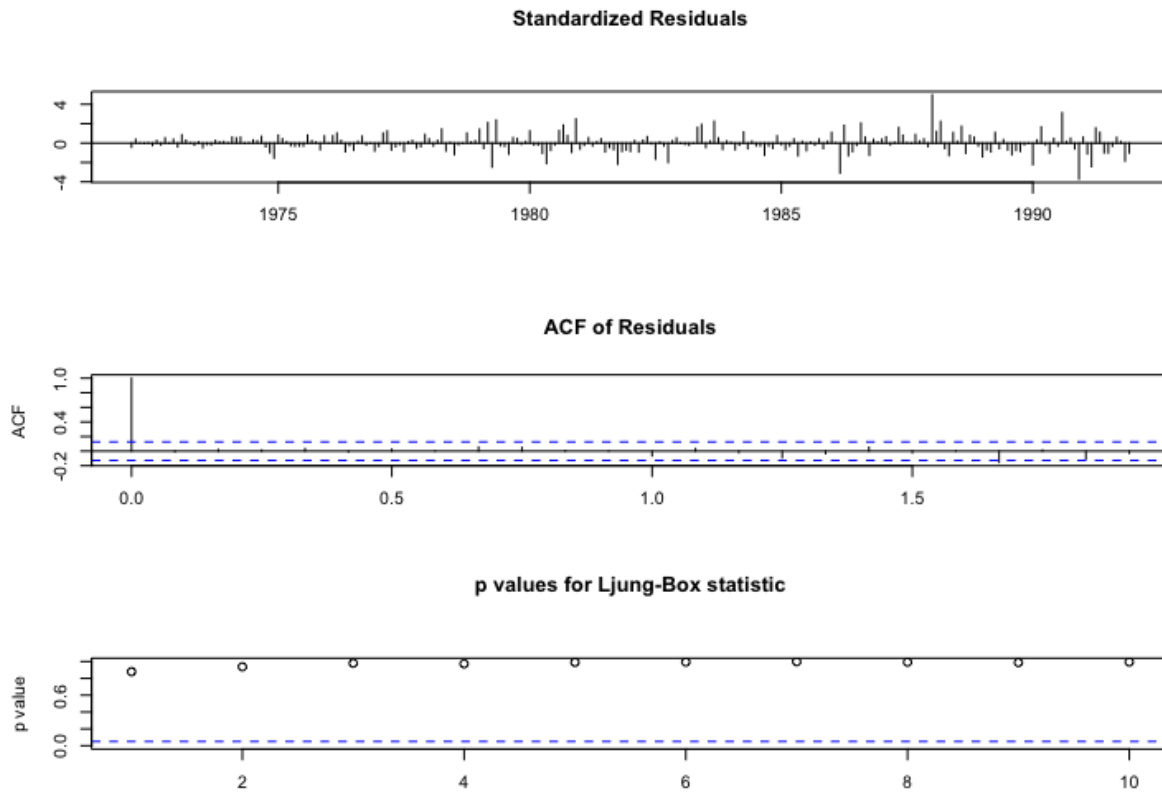
Figure 3.8: Residual diagnostics for the ARIMA fit over the detrended Chemicals data. These show that the residuals closely follow a normal distribution. Also, we observe that the AIC and BIC values have significantly reduced as compared to the harmonic fit over the same data.

```
Transformed Chemicals Data

Call:
lm(formula = chemicals.bc.tr ~ (har + time(chemicals.bc.tr)))

Residuals:
    Min       1Q    Median       3Q      Max
-0.40239 -0.13766 -0.00284  0.12268  0.47168

Coefficients:
                       Estimate Std. Error t value Pr(>|t|)
(Intercept)          -2.587e+02  3.716e+00 -69.611  < 2e-16 ***
harcos(2*pi*t)       -1.439e-02  1.607e-02  -0.895    0.371
harsin(2*pi*t)        6.795e-02  1.608e-02   4.226 3.34e-05 ***
time(chemicals.bc.tr)  1.322e-01  1.875e-03  70.490  < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

**Standardized Residuals**



**ACF of Residuals**
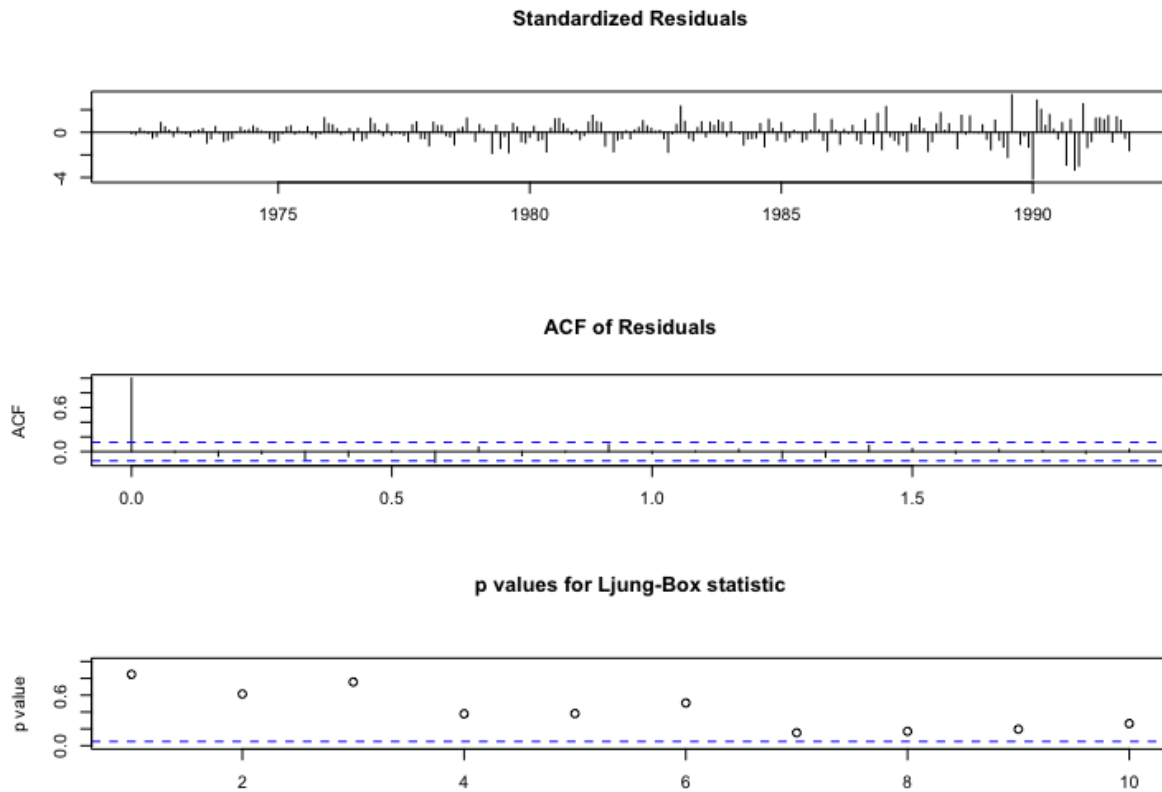


**p values for Ljung-Box statistic**



Figure 3.9: Residual diagnostics for the ARIMA fit over the detrended Vehicles data. These show that the residuals do not closely follow a normal distribution. But, we observe that the AIC and BIC values have significantly reduced as compared to the harmonic fit over the same data.

```
Residual standard error: 0.1803 on 248 degrees of freedom
Multiple R-squared:  0.9525,Adjusted R-squared:  0.9519
F-statistic:  1658 on 3 and 248 DF,  p-value: < 2.2e-16


Transformed Vehicles Data
Call:
lm(formula = vehicles.bc.tr ~ (har + time(vehicles.bc.tr)))


Residuals:
     Min       1Q   Median       3Q      Max
-0.51612 -0.12429  0.02448  0.12732  0.41170


Coefficients:
                    Estimate Std. Error t value Pr(>|t|)
(Intercept)        -1.272e+02  3.673e+00 -34.628   <2e-16 ***
harcos(2*pi*t)      2.391e-02  1.588e-02   1.506   0.1334
```
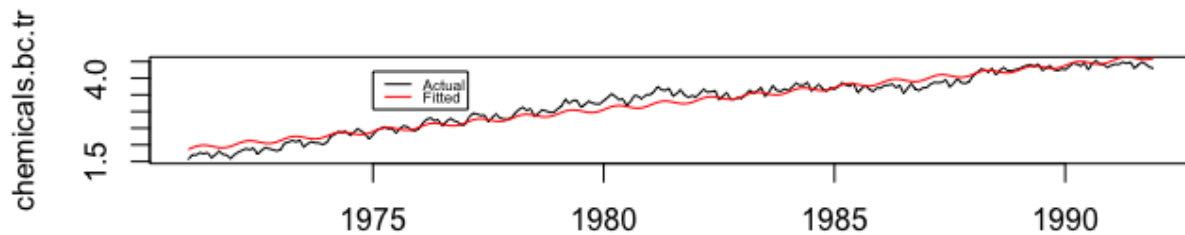
```
harsin(2*pi*t)          4.306e-02  1.589e-02   2.710    0.0072 **
time(vehicles.bc.tr)  6.534e-02  1.854e-03   35.250    <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1783 on 248 degrees of freedom
Multiple R-squared:  0.834,Adjusted R-squared:  0.832
F-statistic: 415.3 on 3 and 248 DF,  p-value: < 2.2e-16
```

## TS Plot: Transformed Chemicals Data
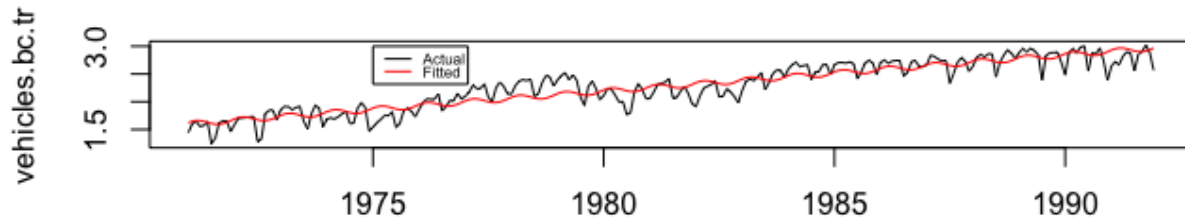
## TS Plot: Transformed Vehicles Data

Figure 3.10: Fitting a harmonic model with linear trend over the transformed data. There is not much improvement in the fit as compared to Fig 3.4. In fact, the fit over the Chemicals data is worse. Therefore, we do not perform any other diagnostics for this model and carry on with the detrended and stabilized datasets.

Here, in Fig 3.11 we observe the temporal patterns of the stabilized datasets (transformed, detrended, and seasonality removed).

Now, we check the different ACF plots depicted in Fig 3.12. The EACF matrices for these stabilized datasets are given here.

```
> eacf(chemicals.bc.diff2)
AR/MA
  0 1 2 3 4 5 6 7 8 9 10 11 12 13
0 x x o o o o o o o o o  x  x  o  o
```
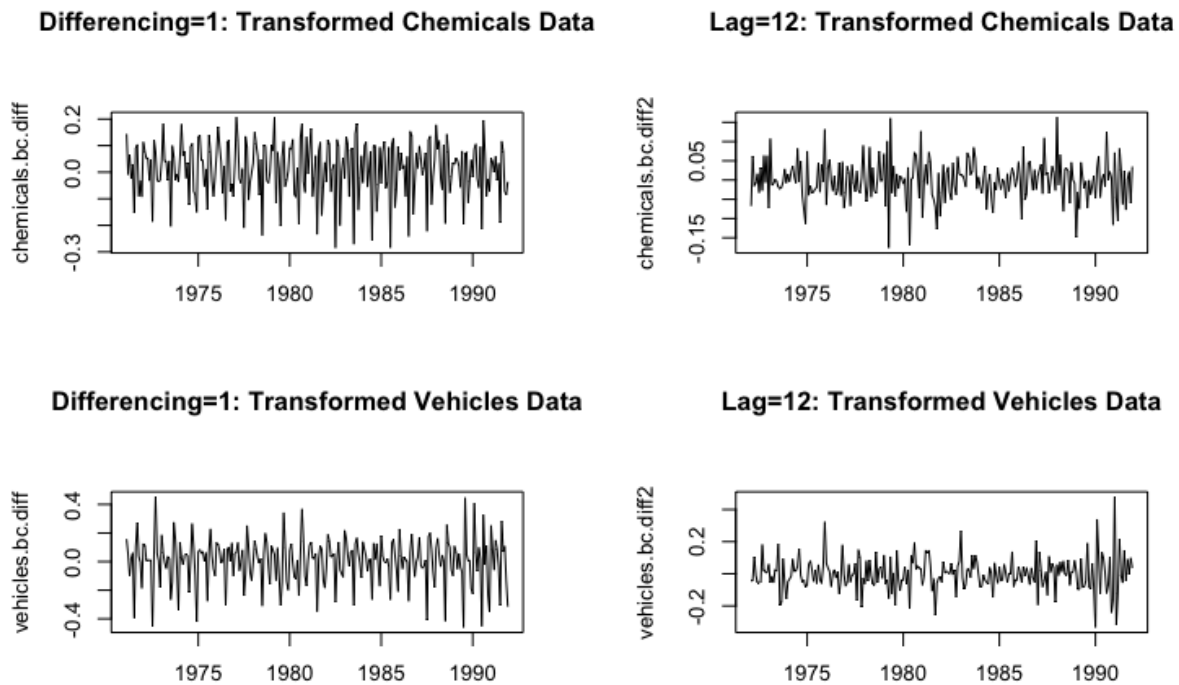
Figure 3.11: Time series plots of the stabilized data. The plots are constructed the same way as in Fig 3.3. However, we see that the stabilized data appears to be more stationary than before.

```
1 x x o o o o o o o o o o   x   x   o
2 x x o o o o o o o o o o   x   x   o
3 o x x o o o o o o o o o   x   x   o
4 x x o o o o o o o o o o   x   x   x
5 x x o o o o o o o o o o   x   o   x
6 x o o o x o o o o o o o   x   o   o
7 x o o x o o o o o o o o   x   x   o

> eacf(vehicles.bc.diff2)
AR/MA
  0 1 2 3 4 5 6 7 8 9 10 11 12 13
0 x x o o x o o o o o o x  x   o   o
1 x x o o x o o o o o o o  x   x   o
2 x o o o x o o o o o o o  x   x   x
3 o o x o x o o o o o o o  x   x   x
4 o x x x o o o o o o o o  x   x   o
5 x x x x o o o o o o o o  x   x   o
6 x x o x o o o o o o o o  x   x   o
7 x x o o o o x o o o o o  x   o   o
```

As expected from Fig 3.12, we have greatly reduced the AIC and BIC values using
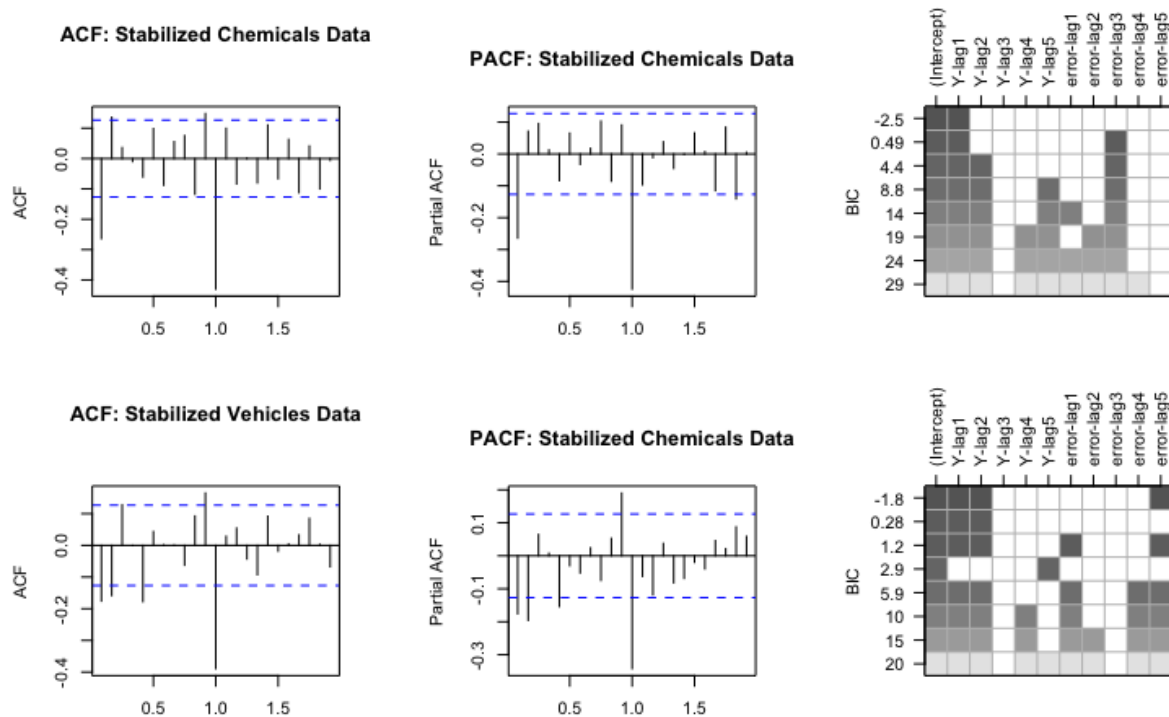
32

Figure 3.12: ACF, PACF, and ARMA subset plots for the stabilized Chemicals and Vehicles data. These show that the detrended data could possible have both AR and MA parts like before but the AIC and BIC values could be significantly reduced (since $\lambda_C \approx 0.19$ and $\lambda_V \approx -0.03$ for Chemicals and Vehicles data respectively) as compared to Fig 3.7.

transformation. The model details obtained using the auto.arima function are given here. It is noteworthy that, while the number of parameters for the Chemicals data remain the same as before, we have been able to reduce one parameter in case of the Vehicles data. Hence, we do not further check the residuals.

```
> chemicals.bc.diff2.arima = auto.arima(chemicals.bc.diff2)
> chemicals.bc.diff2.arima

Series: chemicals.bc.diff2
ARIMA(2,0,2)(0,0,1)[12] with zero mean

Coefficients:
          ar1      ar2     ma1      ma2      sma1
      -0.9330  -0.9102  0.7469  0.7854  -0.6996
s.e.   0.0548   0.0595  0.0850  0.0741   0.0641

sigma^2 estimated as 0.001787:  log likelihood=415.93
AIC=-819.86    AICc=-819.5    BIC=-799.01

> vehicles.bc.diff2.arima = auto.arima(vehicles.bc.diff2)
```

```
> vehicles.bc.diff2.arima

Series: vehicles.bc.diff2
ARIMA(0,0,2)(0,0,2)[12] with zero mean

Coefficients:
          ma1       ma2      sma1      sma2
      -0.2106   -0.1494   -0.6332   -0.1205
s.e.   0.0676    0.0740    0.0678    0.0723

sigma^2 estimated as 0.006902:  log likelihood=252.91
AIC=-495.82    AICc=-495.56    BIC=-478.44
```

Therefore, from the above results, it can be concluded that transformation improves the fit of these models. In order to check the cross-correlation between the Chemicals and Vehicles data, we use the ccf function in R (Fig 3.13). This shows that the sales of chemicals and vehicles are positively correlated.
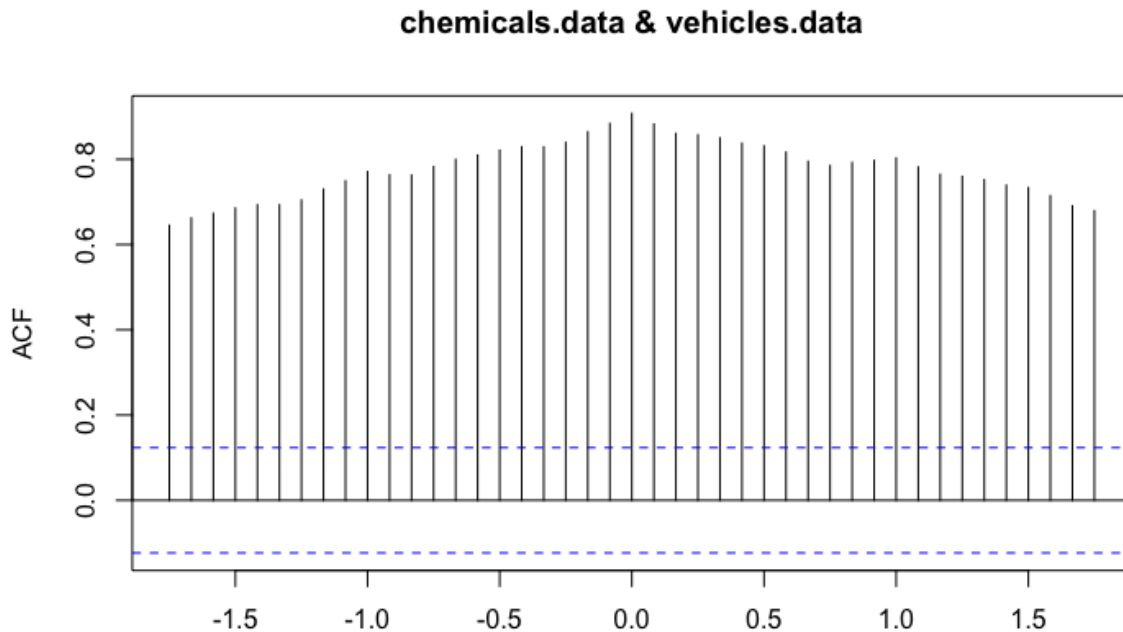
## chemicals.data & vehicles.data



Figure 3.13: Cross-correlation plot of Chemicals and Vehicles data.

(iv) The strong positive correlation between the Chemicals and Vehicles data could be attributed to the increase of vehicle usage over time thereby necessitating more use of chemicals not only for vehicle parts but also for additional utilities like stain removers etc. However, if the month-to-month variation is considered, there are bigger variations in the sales of vehicles. This is because Chemical sales on a monthly basis will not fluctuate much as chemicals are required in a variety of industries. But in case of vehicles, there are additional factors involved such as the economic slowdown which significantly affect the sales.

# 4    Problem 4

France has enjoyed a long and storied history as a world leader in culture, food, literature, film, and revolutions. Over its history, the average life expectancy of people in France has changed notably. The time series object described below contains the annual average life expectancy in France (measured in years lived), measured during the years 1816 to 2019. The data is in the file FLE.txt.

Analyze the data in the time series and write a report addressing the following:

(a) What trend model(s) best capture the trends in French life expectancy over time?

(b) Are there any noticeable patterns, or observations apparent in the series? Can you tell from your plot what could be the causes of your observations?

(c) Regardless of what you observe, suppose you apply some transformation to the data, how does the new data looks like? Any changes in trend?

(d) For the various models you tried, assess the fit of the models using any tools at your disposal. Perform any necessary test to make sure sure your model is appropriate. Write down the equation of your final model. Are any transformations of the data necessary? Perform all diagnostics tests.

Now conclude with brief summary: How the life expectancy changes over time, both long-term over the observed period of years, and in terms of patterns of year-to-year variation. Add any necessary graphs, tests and/or confidence interval etc...

## 4.1    R Code

```
1  ### STAT 5814 MIDTERM 2/PROBLEM 4
2  ### AUTHOR: SAYANTAN MAJUMDAR
3  ### EMAIL: smxnv@mst.edu
4  ### SNO: 12566087
5
6  library(TSA)
7  library(snpar)
8  library(stats)
9  library(forecast)
10 library(tseries)
11 library(fUnitRoots)
12
13 data = read.table('Data/ExamData/FLE.txt')
14 fle.data = ts(data$V2, frequency=1, start=1816, end=2019)
15 par(mfrow=c(1, 1))
16 plot(fle.data, main='Time Series Plot of French Life Expectancy', xlab='
     Time (Years)', ylab='Annual Average Life Expectancy (years)')
17 # Stationarity Check
18 adf.test(fle.data)
19
```

```
20 # Quadratic fit over original data
21 fle.fit = lm(fle.data ~ time(fle.data) + I(time(fle.data)^2))
22 fle.residuals = rstudent(fle.fit)
23 summary(fle.fit)
24 fle.model = ts(fle.fit$fitted.values, frequency=1, start=1816, end=2019)
25
26 plot(fle.data, main='Quadratic Fit: Life Expectancy Data')
27 lines(fle.model, col='red')
28 legend(1850, 70, legend=c("Actual", "Fitted"), col=c("Black", "Red"), lty
      =1:1, cex=0.7)
29
30 # Transformed Data
31 fle.bc.lambda = BoxCox.lambda(fle.data)
32 fle.bc.tr = BoxCox(fle.data, lambda=fle.bc.lambda)
33 plot(fle.bc.tr, main='Box Cox Transformed Data')
34
35 fle.bc.fit = lm(fle.bc.tr ~ time(fle.bc.tr) + I(time(fle.bc.tr)^2))
36 fle.bc.model = ts(fle.bc.fit$fitted.values, frequency=1, start=1816, end
      =2019)
37 summary(fle.bc.fit)
38
39 plot(fle.bc.tr, main='Box Cox Transformed Data')
40 lines(fle.bc.model, col='red')
41 legend(1850, 2500, legend=c("Actual", "Fitted"), col=c("Black", "Red"),
      lty=1:1, cex=0.7)
42
43 # Model assessment for the transformed data
44 fle.bc.residuals = rstudent(fle.bc.fit)
45
46 # Detrending of the original data
47 fle.detrend = diff(fle.data, differences=2)
48 plot(fle.detrend, main='Detrended FLE Data')
49 # Stationarity Check
50 adf.test(fle.detrend)
51
52 par(mfrow=c(3, 1))
53 # ACF, PACF, and EACF for Detrended FLE Data
54 acf(fle.detrend, main="ACF: Detrended FLE Data")
55 pacf(fle.detrend, main="PACF: Detrended FLE Data")
56 eacf(fle.detrend)
57 plot(armasubsets(fle.detrend, 5, 5))
58
59 # ARIMA Model Fits for Detrended FLE Data
60 fle.detrend.arima = arima(fle.detrend, order=c(1, 0, 2))
61 fle.auto.arima = auto.arima(fle.detrend)
62 fle.bc.auto.arima = auto.arima(fle.detrend, lambda="auto")
63
64 # Residual Diagnostics (manually repeated for all the models)
65 par(mfrow=c(1, 3))
66 dataset = fle.detrend
67 model = fle.auto.arima
68 # st_residuals = fle.bc.residuals
69 st_residuals = rstandard(model)
70 plot(y=st_residuals, x=as.vector(time(dataset)), type = 'o', ylab = '
```

```
      Standardized Residuals', xlab = 'Time', main='Standardized Residual
      Plot')
71 hist(st_residuals, main="Residuals Histogram", xlab='Residuals')
72 qqnorm(st_residuals, main="QQ Plot for Residuals")
73 qqline(st_residuals, col="Red")
74 boxplot(st_residuals, main='Residual Boxplot', ylab='Residuals')
75 acf(st_residuals, main='ACF Plot')
76 tsdiag(model)
77 # Residual Normality
78 shapiro.test(st_residuals)
79 snpar::runs.test(st_residuals, exact=TRUE)
80 # AIC and BIC tests
81 AIC(model)
82 BIC(model)
83 confint(fle.bc.auto.arima)
```
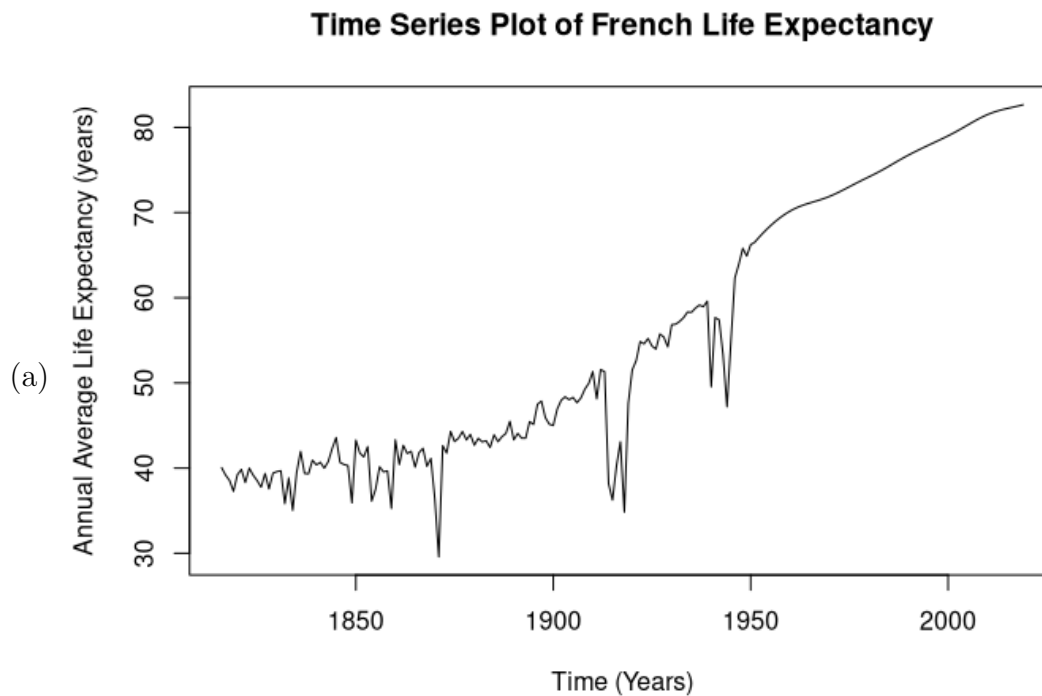
## 4.2   Results



(a)

Figure 4.1: Time series plot of the French life expectancy data.

From Fig 4.1 it can be said that a quadratic trend is suitable for this dataset. The model summary is shown here.

```
Call:
lm(formula = fle.data ~ time(fle.data) + I(time(fle.data)^2))


Residuals:
     Min       1Q   Median       3Q      Max
-17.2140  -1.3665   0.4999   1.8464   5.6716
```

```
Coefficients:
                       Estimate Std. Error t value Pr(>|t|)
(Intercept)            3.499e+03  3.042e+02   11.50   <2e-16 ***
time(fle.data)        -3.849e+00  3.175e-01  -12.12   <2e-16 ***
I(time(fle.data)^2)   1.069e-03  8.279e-05   12.92   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.668 on 201 degrees of freedom
Multiple R-squared:  0.946,Adjusted R-squared:  0.9455
F-statistic:  1762 on 2 and 201 DF,  p-value: < 2.2e-16
```

Here, we observe that the F-statistic and the residual standard are quite high. Moreover, the augmented Dickey-Fuller stationarity test (Dickey-Fuller = -2.9107, Lag order = 5, p-value = 0.1945) confirms that the original data is not stationary. However, by only visually observing the original data (Fig 4.1), we can say that the quadratic fit (Fig 4.2) is the best trend for this dataset.
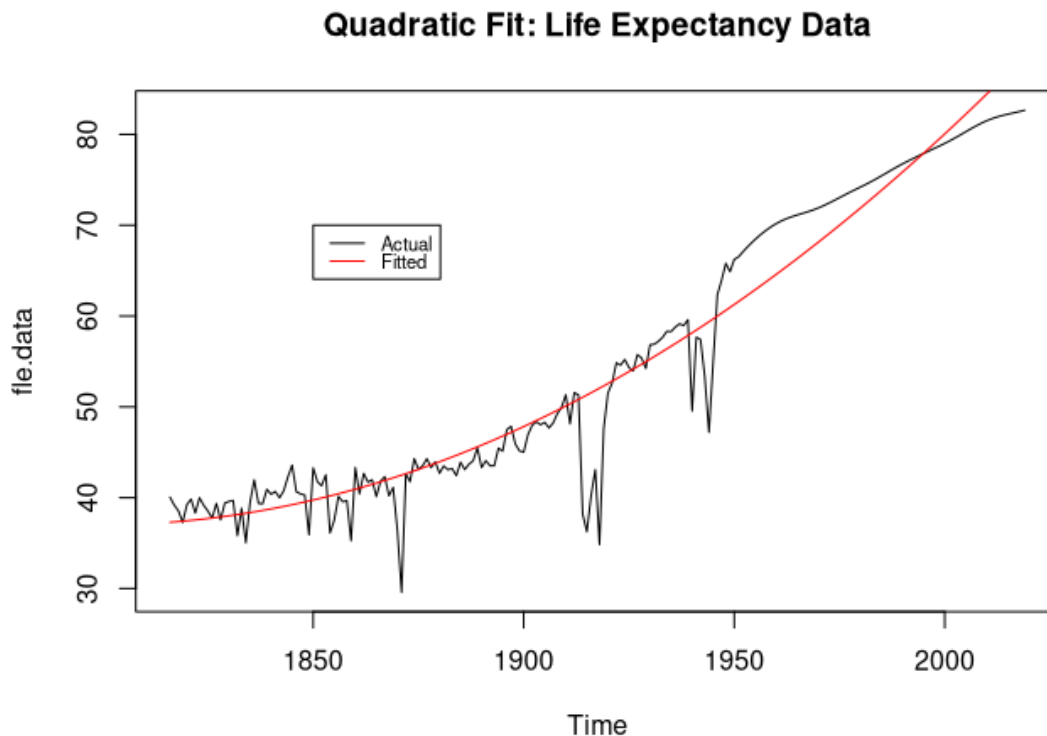
**Quadratic Fit: Life Expectancy Data**



Figure 4.2: Quadratic fit over the original French life expectancy data. The fit performs poorly from 1950 onward as we observe a linear trend from 1950-2019.

(b) The fluctuations in the dataset are mostly present before 1950. The increase in the life expectancy values between 1816-1850 could be attributed to the advancements

in medicine. However, the sharp drops in the values, particularly during mid 1915-1920 and mid 1940s could be caused by the two world wars. Interestingly, the life expectancy increases post 1950 which could be primarily attributed to the invention of modern healthcare and economic growth.
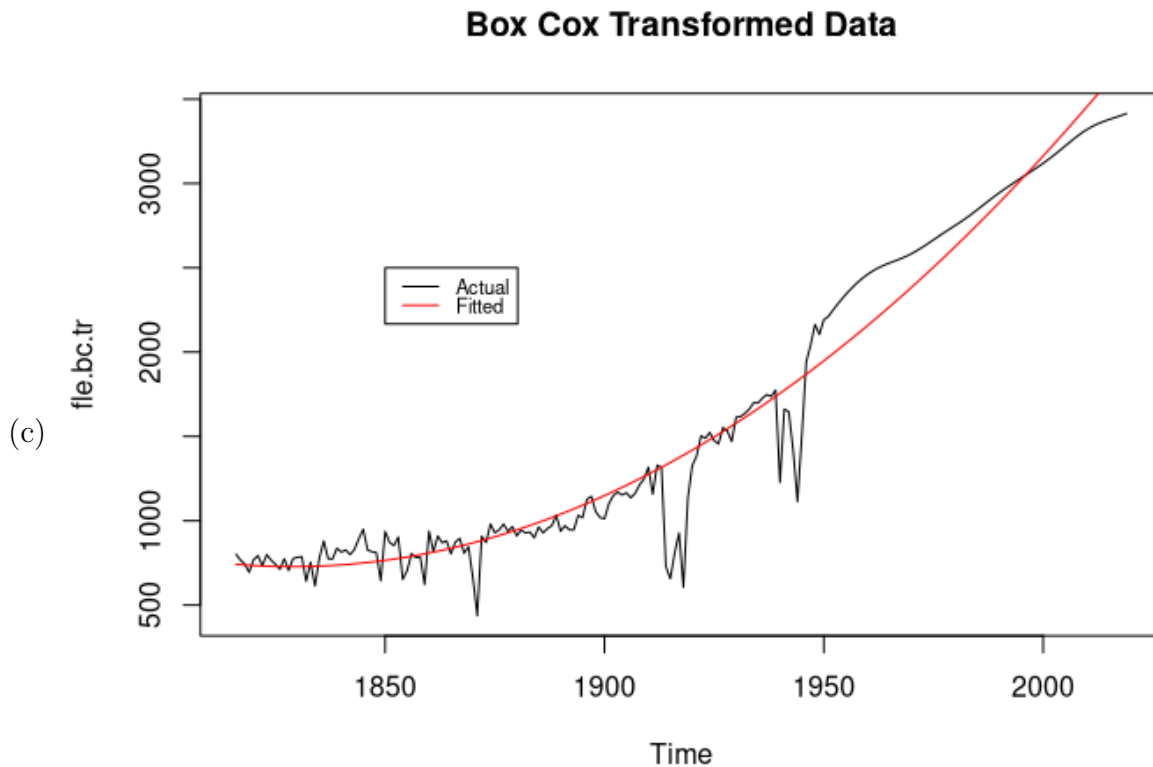
(c)



**Box Cox Transformed Data**

Figure 4.3: Transformed FLE Data using Box Cox method. Here, this transformation is obtained by auto-setting $\lambda \approx 1.99$. The quadratic trend still remains.

The model fit summary for this transformed dataset is given here.

```
Call:
lm(formula = fle.bc.tr ~ time(fle.bc.tr) + I(time(fle.bc.tr)^2))

Residuals:
    Min      1Q  Median      3Q     Max
-782.46  -70.39   23.94   83.28  305.68

Coefficients:
                       Estimate Std. Error t value Pr(>|t|)
(Intercept)           2.787e+05  1.494e+04   18.66   <2e-16 ***
time(fle.bc.tr)      -3.040e+02  1.559e+01  -19.50   <2e-16 ***
I(time(fle.bc.tr)^2)  8.311e-02  4.064e-03   20.45   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 180.1 on 201 degrees of freedom
Multiple R-squared:  0.9625,Adjusted R-squared:  0.9621
F-statistic:  2579 on 2 and 201 DF,  p-value: < 2.2e-16
```

(d) The residual diagnostics for the quadratic model are shown in Fig 4.4. We further
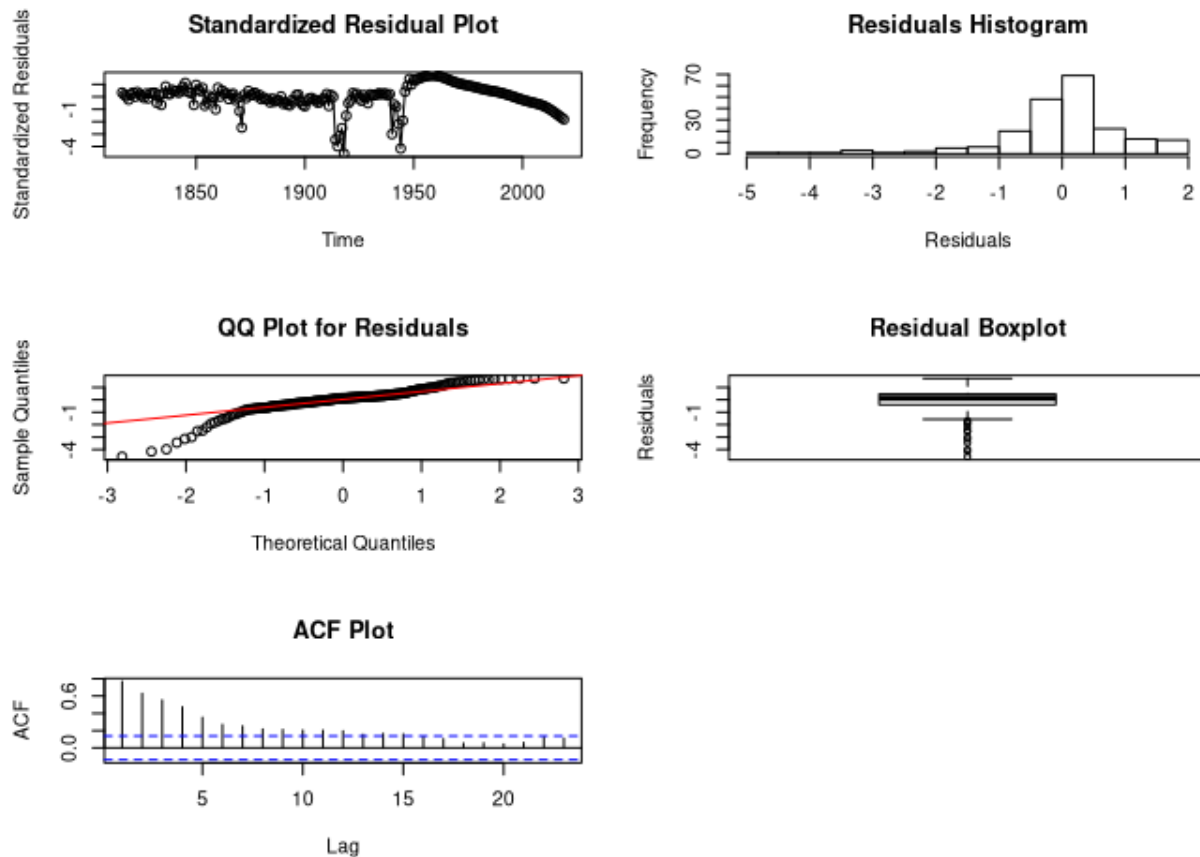


Figure 4.4: Residual diagnostics for the quadratic model. Here, we observe that this model fits poorly as the residuals follow a non-normal distribution.

verify the non-normality of the residuals using the Shapiro-Wilk and Exact Runs tests.

```
Quadratic Model

Shapiro-Wilk normality test
W = 0.88032, p-value = 1.21e-11

Exact runs test
Runs = 46, p-value = 3.374e-16
alternative hypothesis: two.sided
```

Next, we detrend the transformed dataset (Fig 4.5) using two differences because of the quadratic trend (i.e. the transformed dataset is not stationary). Since no seasonality is present, we set lag=1. Once the detrending is done, we check the ACF, PACF, EACF,
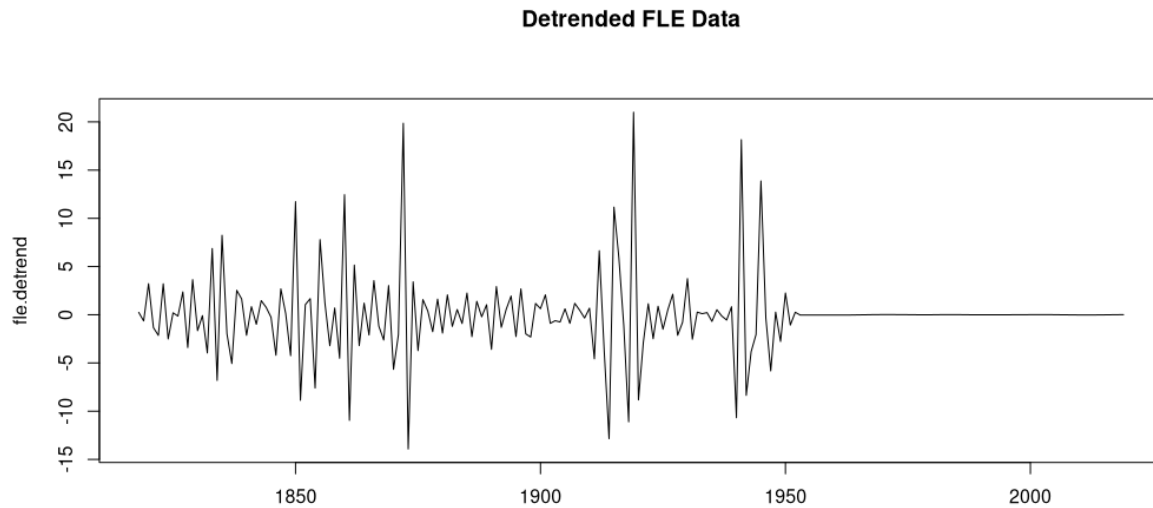
**Detrended FLE Data**



Figure 4.5: Plot of the detrended FLE data. This shows that we have achieved stationarity and can now proceed with various ARIMA modeling. We also verify this using the augmented Dickey-Fuller test (Dickey-Fuller = -10.009, Lag order = 5, p-value = 0.01) where we accept the alternative hypothesis of stationarity.

and the ARMA subset plots as given in Fig. 4.6.

```
> eacf(fle.detrend)

AR/MA
  0 1 2 3 4 5 6 7 8 9 10 11 12 13
0 x o o o o o o o o o o  x  x  o
1 x o o o o o o o o o o  o  o  o
2 x x o o o o o o o o o  o  o  o
3 x o x o o o o o o o o  o  o  o
4 x o x x x o o o o o o  o  o  o
5 x x o o x o o o o o o  o  o  o
6 x x x o o o x o o o o  o  o  o
7 x x x o o o o x o o o  o  o  o
```

Initially, we consider ARMA(1, 1) model or ARIMA(1, 0, 1). The model fit summary is given here.

```
>arima(x = fle.detrend, order = c(1, 0, 1))

Coefficients:
```
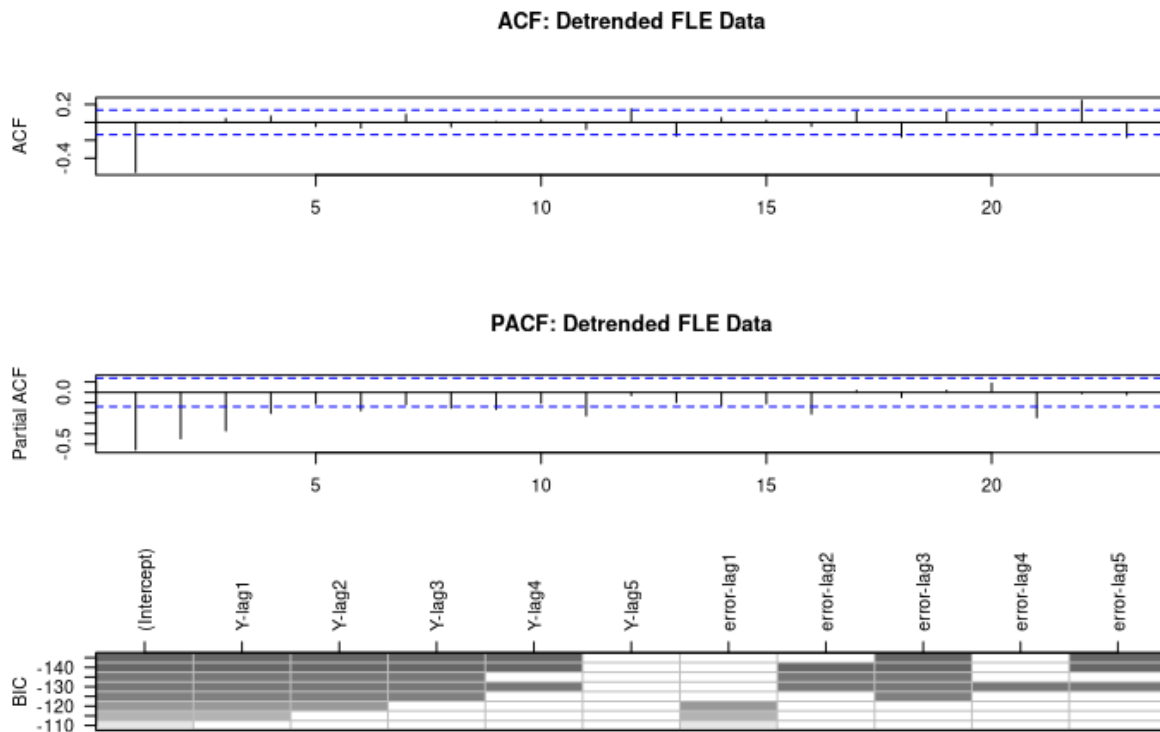
**ACF: Detrended FLE Data**

**PACF: Detrended FLE Data**

Figure 4.6: ACF, PACF, and ARMA subset plots of the detrended FLE data. From these we can say that ARMA(p, q) is a suitable model.

```
          ar1       ma1   intercept
       -0.2498    -1.000      0.0017
s.e.    0.0680     0.013      0.0025


sigma^2 estimated as 6.765:  log likelihood = -482.63,  aic = 971.26
```

We also perform the residual diagnostics (Figures 4.7 and 4.8)and apply the normality checks. The residual diagnostics show that the ARMA(1, 1) residuals do not closely follow a normal distribution. This is further confirmed by the Shapiro-Wilk and Runs tests.

```
ARMA(1, 1) for the detrended FLE data


Shapiro-Wilk normality test
W = 0.76725, p-value < 2.2e-16


Exact runs test
Runs = 86, p-value = 0.02855
alternative hypothesis: two.sided
```

**Standardized Residuals**



**ACF of Residuals**



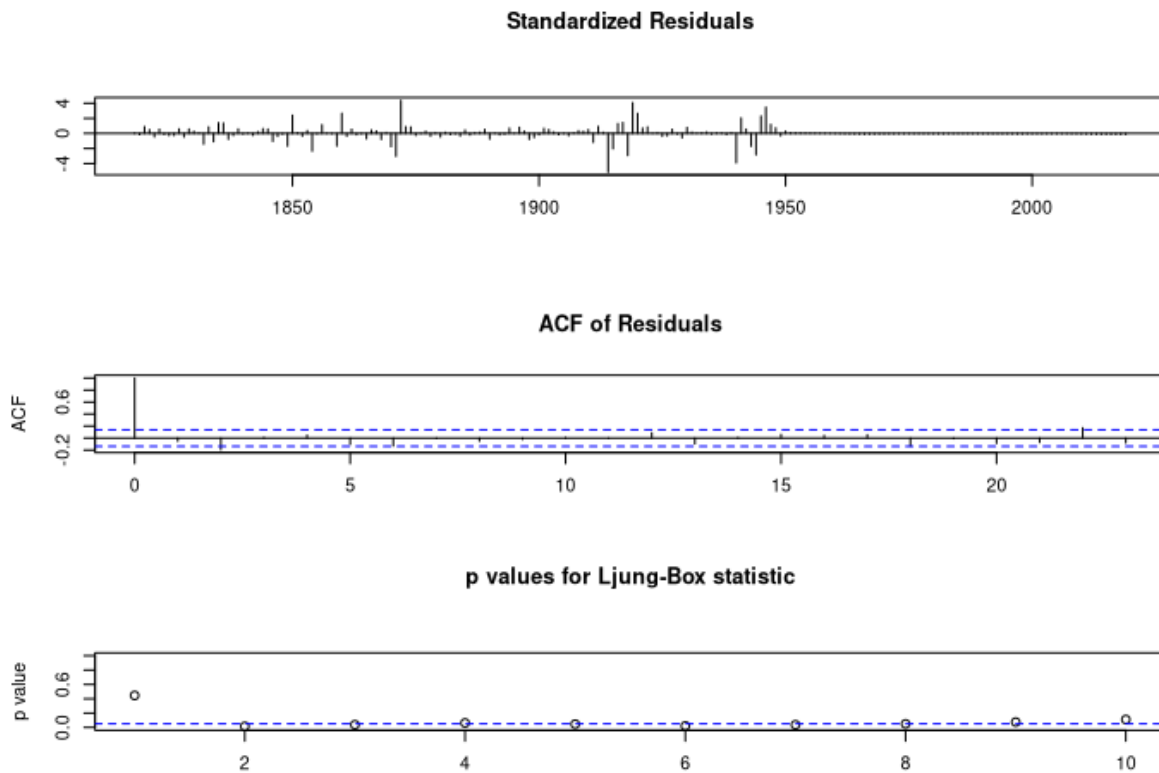**p values for Ljung-Box statistic**



Figure 4.7: Residual diagnostics for the ARMA(1, 1) model for the FLE data.

```
> AIC(model)
[1] 973.2626
> BIC(model)
[1] 986.4957
```

Next, we try with ARMA(1, 2) to see if we have a better fit.

```
Call:
arima(x = fle.detrend, order = c(1, 0, 2))

Coefficients:
         ar1      ma1     ma2  intercept
      0.5658  -1.8831  0.8831     0.0017
s.e.  0.0946   0.0600  0.0594     0.0009

sigma^2 estimated as 6.266:  log likelihood = -476.12,  aic = 960.25
```

The residual diagnostics (Figures 4.9 and 4.10) and normality checks are carried out like before. The residual diagnostics show that the ARMA(1, 1) residuals do not closely follow a normal distribution. This is further confirmed by the Shapiro-Wilk and Runs tests.
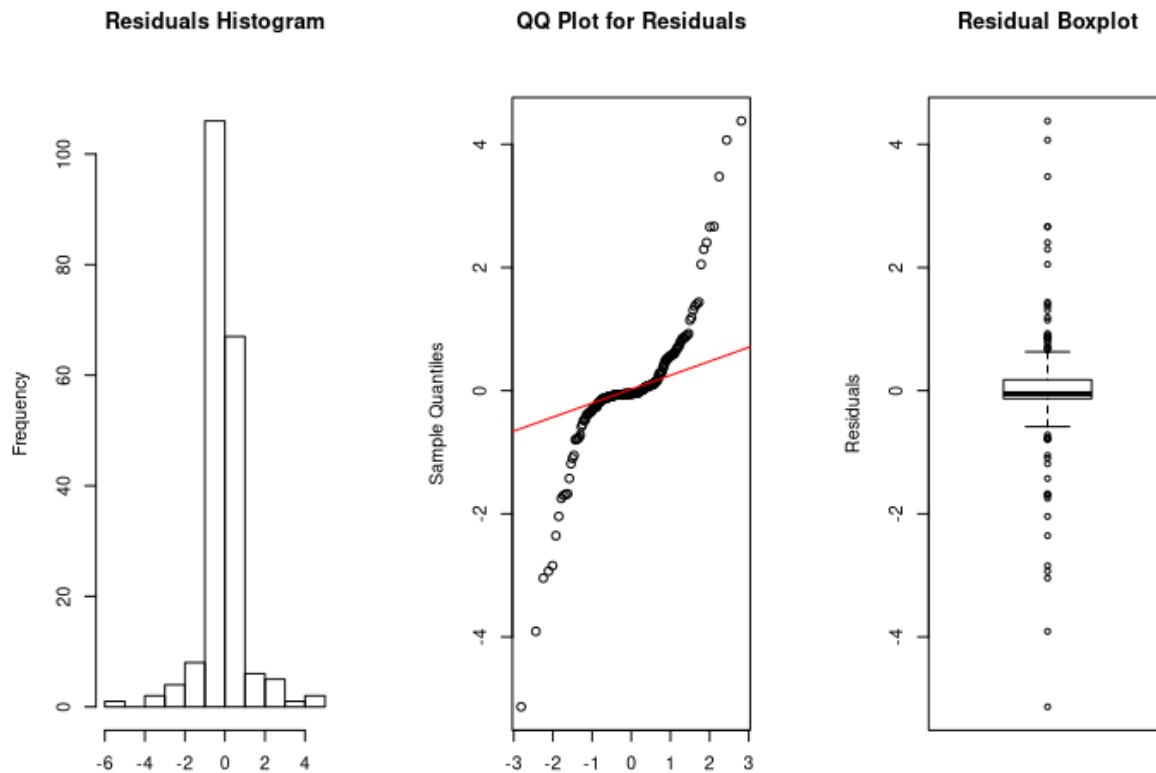
Figure 4.8: Residual plots for the ARMA(1, 1) model for the FLE data.

```
ARMA(1, 2) for the detrended FLE data

Shapiro-Wilk normality test
W = 0.81151, p-value = 6.773e-15

Exact runs test
Runs = 66, p-value = 4.192e-07
alternative hypothesis: two.sided

> AIC(model)
[1] 962.246
> BIC(model)
[1] 978.7874
```

From these we can conclude that ARMA(1, 1) is preferred model. However, to obtain a better fit, we use the auto.arima function and check the effect of transformation. The model summary including the AIC and BIC values and the confidence intervals of the parameters are given here.
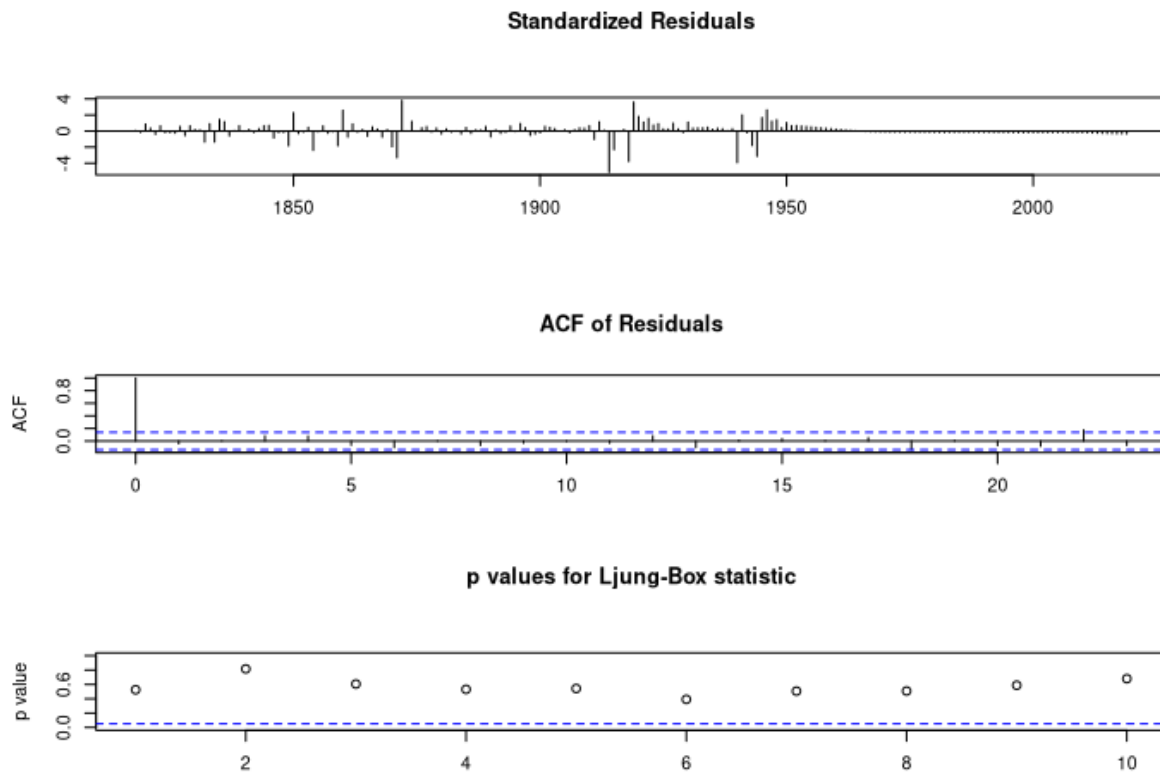
Figure 4.9: Residual diagnostics for the ARMA(1, 2) model for the FLE data.

```
> fle.auto.arima = auto.arima(fle.detrend)
> fle.auto.arima
Series: fle.detrend
ARIMA(1,0,3) with zero mean

Coefficients:
         ar1      ma1     ma2      ma3
      0.7382  -2.0855  1.2361  -0.1467
s.e.  0.0960   0.1212  0.2218   0.1026

sigma^2 estimated as 6.436:  log likelihood=-475.51
AIC=961.01    AICc=961.32    BIC=977.55

> fle.bc.auto.arima = auto.arima(fle.detrend, lambda="auto")
> fle.bc.auto.arima
Series: fle.detrend
ARIMA(1,0,2) with non-zero mean
Box Cox transformation: lambda= 0.7566293

Coefficients:
         ar1      ma1     ma2      mean
```
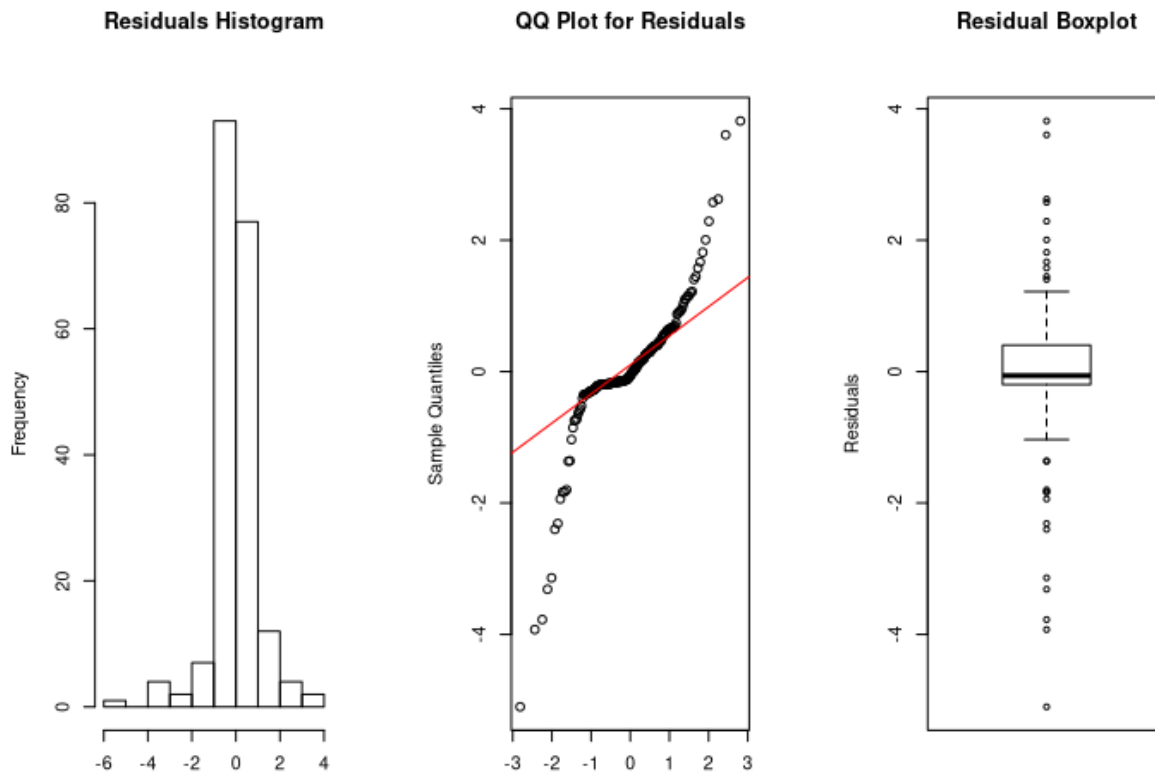
Figure 4.10: Residual plots for the ARMA(1, 2) model for the FLE data.

```
        0.5021  -1.6531   0.7389   -1.3864
s.e.    0.2453   0.1990   0.1701    0.0253


sigma^2 estimated as 4.409:  log likelihood=-435.52
AIC=881.04    AICc=881.34    BIC=897.58


> confint(fle.bc.auto.arima)
                  2.5 %       97.5 %
ar1          0.02135014   0.9827585
ma1         -2.04318152  -1.2630966
ma2          0.40548839   1.0722244
intercept   -1.43600896  -1.3367935
```

Here, we see that the transformed dataset has lower AIC and BIC values along with lesser parameters wherein the confidence intervals do not contain zero. Therefore, data transformation is required.

So from the above analyses, the final model chosen is the ARIMA(1, 2, 2) model (ARMA(1, 2) for the second order differences) where BoxCox transformation ($\lambda \approx 0.75$) is used. The model equation is given by Eq. (3).

$$Z_t = -1.3864 + 0.5021Z_{t-1} + 1.6531\varepsilon_{t-1} - 0.7389\varepsilon_{t-2} + \varepsilon_t \tag{3}$$

where $\varepsilon \sim \mathcal{N}(0, 4.409)$ and $Z_t = \text{BoxCox}(\nabla^2 X_t, \lambda)$, $\lambda \approx 0.75$, $t = \{1, 2, 3, \dots\}$, and $X_t$ being the original FLE data.

If the entire duration of the time series is considered, then the life expectancy has greatly varied before 1950 but there is an overall increase in the life expectancy. From 1950 onward, the life expectancy keeps on increasing linearly.

When individual years are taken into account, then the sharp increments and decrements in the values can be attributed to the improved healthcare systems and degrading economy or wars respectively.