

## Problem 2

### a) Explain why the MA(q) models are regarded as finite memory models

A MA(q) process has finite memory since a shock only affects the process for periods of q ahead. Also, all autocorrelation functions for lags larger than q are equal to zero. Lastly, forecasts for MA(q) process converge to the mean after q steps ahead.

### b) Discuss the ACF behavior of an ARMA(p,q) process

ACF of an Arma() decays geometrically towards zero

## Problem 3

```
library(tseries)
library(fBasics)

## Loading required package: timeDate
## Loading required package: timeSeries
##
## Rmetrics Package fBasics
## Analysing Markets and calculating Basic Statistics
## Copyright (C) 2005-2014 Rmetrics Association Zurich
## Educational Software for Financial Engineering and Computational Science
## Rmetrics is free software and comes with ABSOLUTELY NO WARRANTY.
## https://www.rmetrics.org --- Mail to: info@rmetrics.org
library(forecast)
library(lmtest)

## Loading required package: zoo
##
## Attaching package: 'zoo'
## The following object is masked from 'package:timeSeries':
##
##     time<-
## The following objects are masked from 'package:base':
##
##     as.Date, as.Date.numeric
```

a) Import the data. In R you need to create a time object for index using the ts() function where the starting date is the first month of 1998, and frequency is set equal to 12.

```
setwd("~/Desktop/CSC425/hwork3/")
myd = read.table("INDPRO.csv", header = T, sep = ',')
head(myd)
```

```
##      date      rate
## 1 1/1/1998 0.005042245
## 2 2/1/1998 0.001123778
## 3 3/1/1998 0.000800165
## 4 4/1/1998 0.003449379
## 5 5/1/1998 0.006472101
## 6 6/1/1998 -0.006410126
```

```
tail(myd)
```

```
##      date      rate
## 223 7/1/2016 0.002881344
## 224 8/1/2016 -0.000836184
## 225 9/1/2016 -0.002270310
## 226 10/1/2016 0.001949174
## 227 11/1/2016 -0.006613915
## 228 12/1/2016 0.008298091
```

```
date = myd[,1]
rate = myd[,2] # INDPRO index growth rate series
#create time series object
ratets = ts(rate, start = c(1998,1), freq = 12)
ratets
```

```
##      Jan      Feb      Mar      Apr      May
## 1998 0.005042245 0.001123778 0.000800165 0.003449379 0.006472101
## 1999 0.004601433 0.005267962 0.001684204 0.002550556 0.007486796
## 2000 0.000091600 0.002855146 0.004021726 0.007226003 0.001890342
## 2001 -0.006897773 -0.006088877 -0.002405733 -0.002480285 -0.006611099
## 2002 0.005691094 0.000202648 0.007909293 0.004224599 0.004083065
## 2003 0.005595023 0.003056432 -0.002191305 -0.007151676 0.000163176
## 2004 0.001919137 0.005997042 -0.004703005 0.004165085 0.008357944
## 2005 0.004513475 0.006824789 -0.001670405 0.001460047 0.001670760
## 2006 0.001197194 0.000739109 0.001655129 0.004349440 -0.001389117
## 2007 -0.004881160 0.010230117 0.001709672 0.007162427 0.000438640
## 2008 -0.003014407 -0.003347978 -0.002430179 -0.007423756 -0.005062454
## 2009 -0.023607983 -0.006248784 -0.015756483 -0.008887032 -0.010484095
## 2010 0.011189411 0.003340351 0.006504471 0.004149190 0.015566791
## 2011 -0.000700548 -0.004326531 0.008640668 -0.004202778 0.001959896
## 2012 0.007142720 0.002808449 -0.007023033 0.008909850 0.001648644
## 2013 -0.000760109 0.005117797 0.001498846 -0.000469349 0.000188025
## 2014 -0.004751761 0.008098604 0.007722774 0.001453555 0.003963895
## 2015 -0.004953144 -0.001254828 -0.003236423 -0.002291617 -0.002349126
## 2016 0.004846932 -0.001296037 -0.009467121 0.003993213 -0.001394473
##      Jun      Jul      Aug      Sep      Oct
## 1998 -0.006410126 -0.003570624 0.020470151 -0.001731702 0.007976057
## 1999 -0.001609289 0.006224716 0.004048526 -0.003833174 0.013147521
## 2000 0.000742898 -0.001275552 -0.003327721 0.004060686 -0.002886096
## 2001 -0.006150227 -0.005642703 -0.001899855 -0.003437936 -0.004419474
## 2002 0.009546644 -0.002379221 0.000247961 0.001466118 -0.003492072
## 2003 0.001582439 0.004085057 -0.001838590 0.006355563 0.001214951
```

```

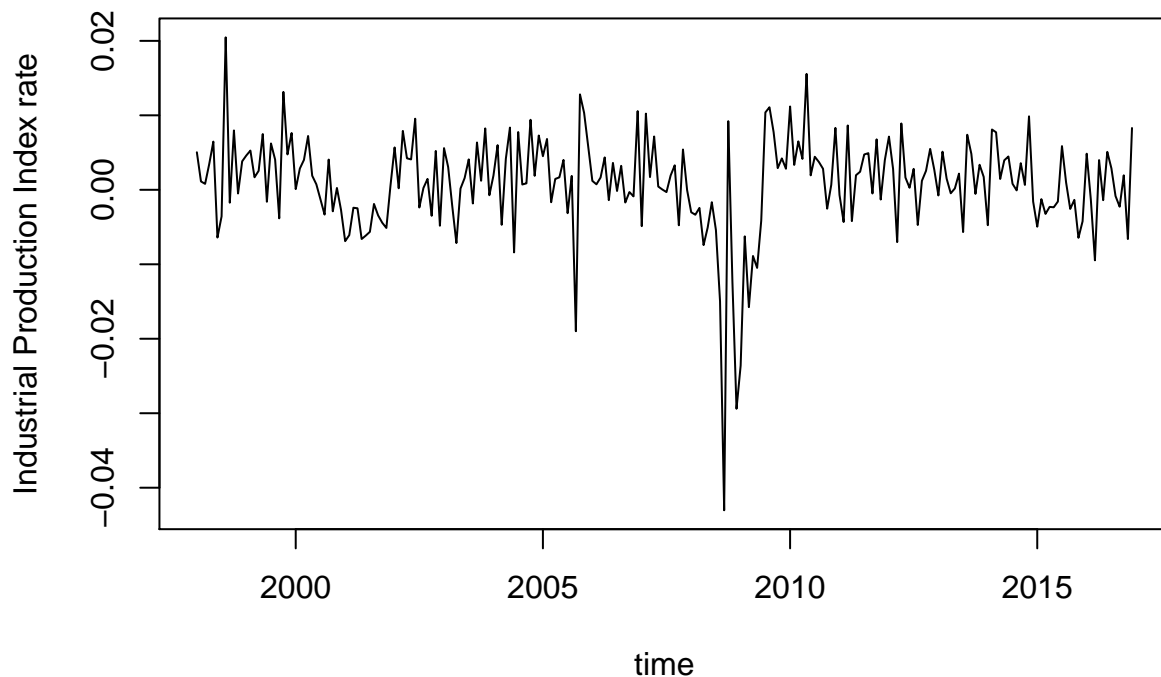
## 2004 -0.0084071360  0.0077605390  0.0007123480  0.0008673950  0.0094033850
## 2005  0.0040133120 -0.0031387960  0.0018668800 -0.0189901820  0.0128105410
## 2006  0.0036112460 -0.0001589900  0.0032193330 -0.0017027160 -0.0003009920
## 2007  0.0000361410 -0.0003090900  0.0019188530  0.0032881810 -0.0047642080
## 2008 -0.0016667620 -0.0054112600 -0.0149665490 -0.0430291090  0.0092165800
## 2009 -0.0041389540  0.0103921060  0.0110811820  0.0077402360  0.0029102970
## 2010  0.0019397950  0.0044404470  0.0037300310  0.0028260880 -0.0025365200
## 2011  0.0024280010  0.0047451250  0.0049189810 -0.0004928650  0.0067746050
## 2012  0.0002759880  0.0028131050 -0.0047281910  0.0012119520  0.0025480220
## 2013  0.0021702410 -0.0056746030  0.0073930090  0.0048228890 -0.0005464230
## 2014  0.0044737860  0.0008085820 -0.0000798424  0.0035751320  0.0006545110
## 2015 -0.0015805630  0.0058706900  0.0009746340 -0.0025677630 -0.0013512850
## 2016  0.0050967420  0.0028813440 -0.0008361840 -0.0022703100  0.0019491740
##
##           Nov           Dec
## 1998 -0.0004994950  0.0038232170
## 1999  0.0047530830  0.0076279430
## 2000  0.0002462260 -0.0027046990
## 2001 -0.0051251340  0.0005470560
## 2002  0.0052207490 -0.0048266760
## 2003  0.0082476380 -0.0007309140
## 2004  0.0018762090  0.0072892630
## 2005  0.0103324840  0.0058066590
## 2006 -0.0008837600  0.0105706120
## 2007  0.0054155830 -0.0000321577
## 2008 -0.0123207760 -0.0294012250
## 2009  0.0042148320  0.0028289890
## 2010  0.0006692030  0.0083145070
## 2011 -0.0013067830  0.0039977210
## 2012  0.0054932120  0.0027102640
## 2013  0.0033633150  0.0017329420
## 2014  0.0098699310 -0.0015803270
## 2015 -0.0064451160 -0.0042292300
## 2016 -0.0066139150  0.0082980910

```

b) Create the time plot of the index growth rate  $X_t$  and analyze trends displayed by the plot

```
plot(ratets, type = 'l', xlab = 'time', ylab = 'Industrial Production Index rate', main = 'Time Plot')
```

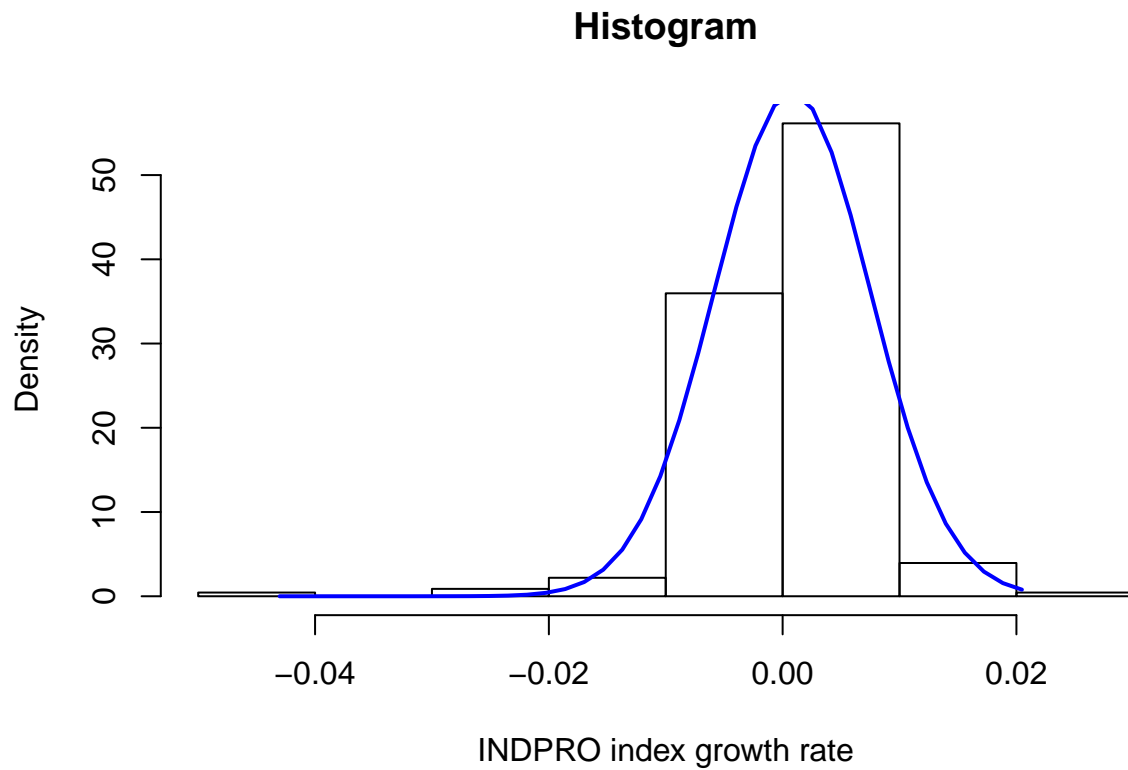
## Time Plot



By looking at the plot, it fluctuates over times and has no clear trends. It seems like its mean and variance over time are more likely constant rather than not constant. However, I see that there is a huge fluctuation between 2008 and 2009.

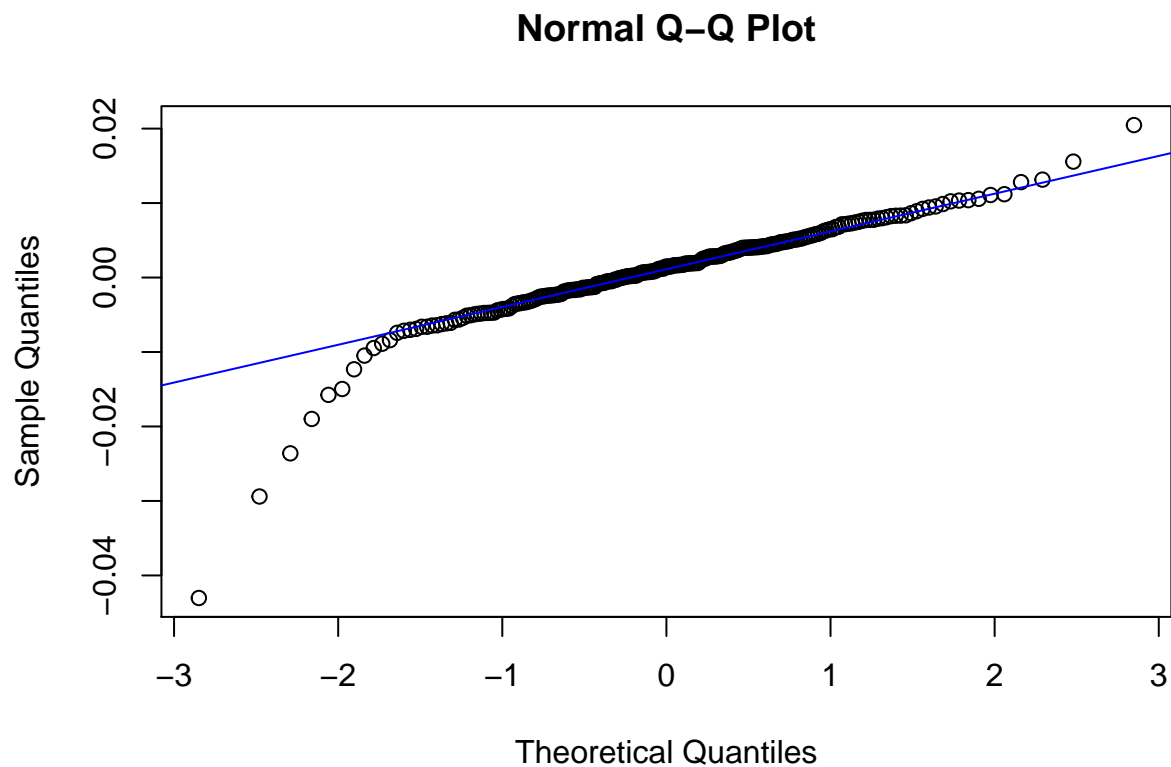
c) Analyze the distribution of  $X_t$ . Can you assume that  $X_t$  is normally distributed?

```
#Histogram
par(mfcol = c(1,1))
hist(rate, xlab = 'INDPRO index growth rate', main = "Histogram", prob = T)
xfit<-seq(min(rate), max(rate), length = 40)
yfit<-dnorm(xfit, mean = mean(rate), sd = sd(rate))
lines(xfit, yfit, col = "blue", lwd = 2)
```



based on the histogram, it is skewed to the left, but I could say it is close to normal distribution

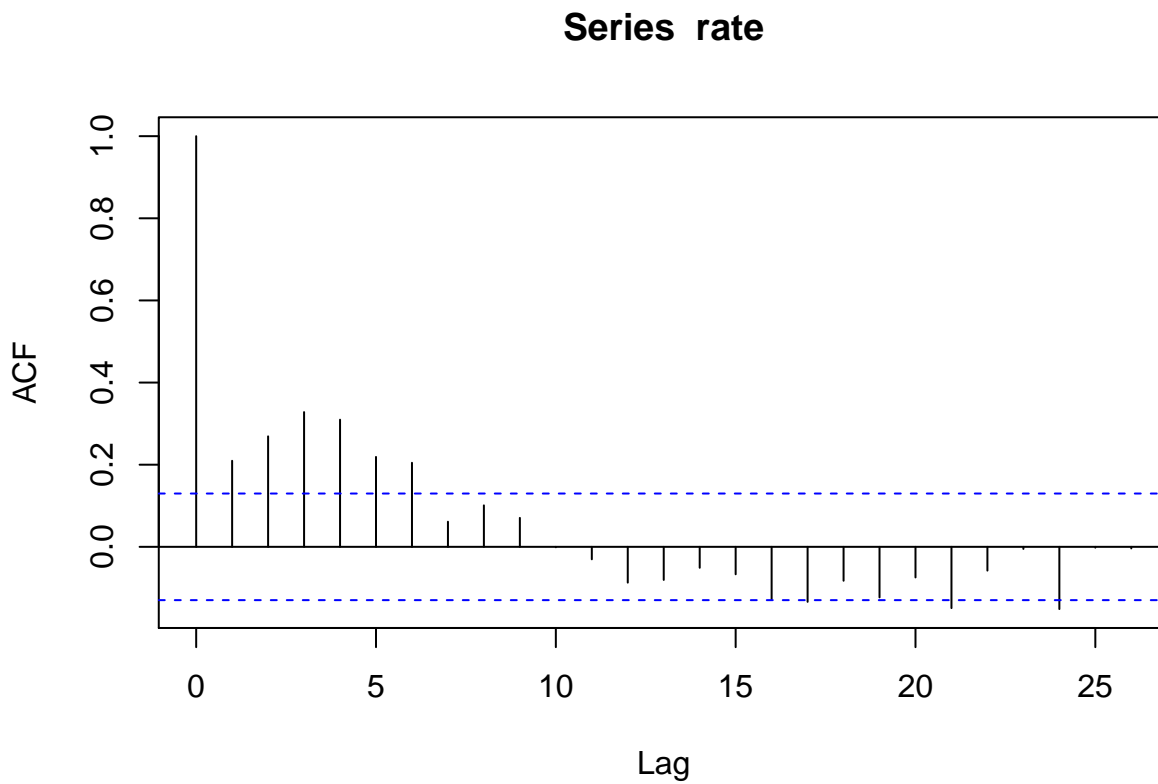
```
#qq-plot  
qqnorm(rate)  
qqline(rate, col = "blue")
```



based on the qq-plot, it is very close to normal distribution.

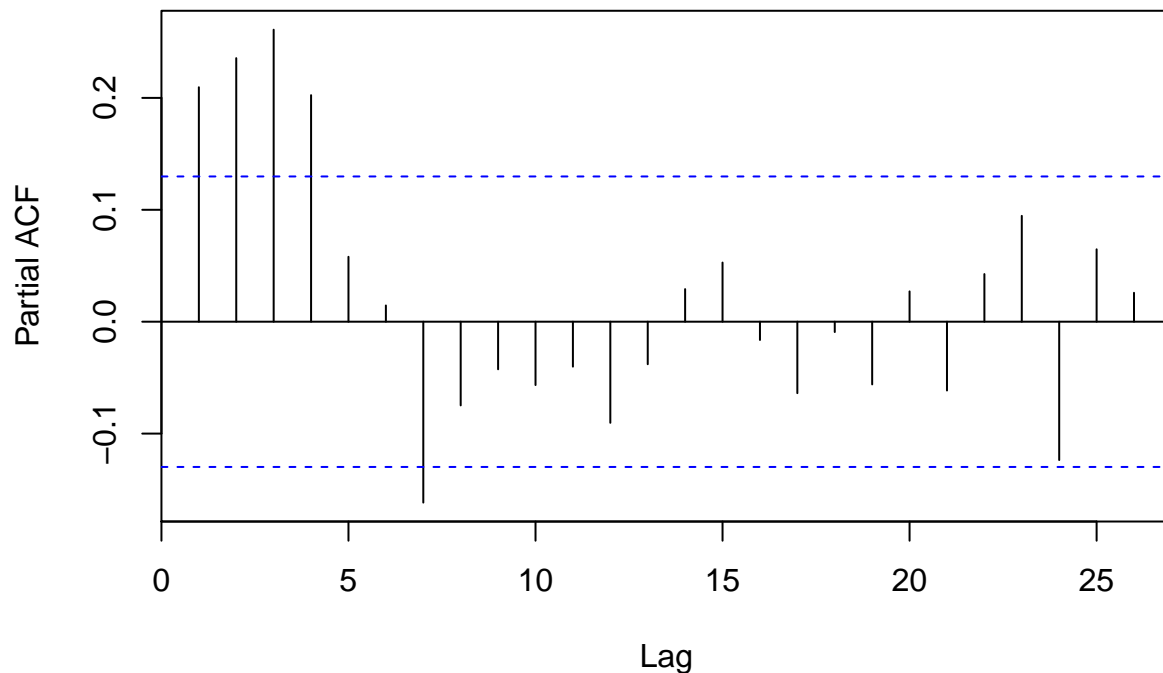
d) Analyze the ACF and the PACF functions and discuss the following questions

```
acf(rate, plot = T, lag = 26)
```



```
pacf(rate, lag = 26)
```

## Series rate



d-1) Does the TS exhibits a stationary behavior?

It is not a perfect stationary, but close to stationary. There are just few significances to be observed over time. Other points are close to zeros.

d-2) Does the process show more of an AR behavior or MA behavior

PACF cuts off after the 4th lag, while the ACF tail off to zero. In this reason, I decide that a model of AR(4) could work for this analysis.

e) Apply the BIC criteria using the `auto.arima(xvar, ic = c("bic"))`function, and fit the model suggested by the BIC criterion. (M1)

```
m1 = auto.arima(ratets, max.P=8, max.Q=8, ic = ("bic"), trace = T, stationary = T)
```

```
##
## Fitting models using approximations to speed things up...
##
## ARIMA(2,0,2)(1,0,1)[12] with non-zero mean : -1657.383
## ARIMA(0,0,0) with non-zero mean : -1627.824
## ARIMA(1,0,0)(1,0,0)[12] with non-zero mean : -1632.173
## ARIMA(0,0,1)(0,0,1)[12] with non-zero mean : -1626.184
## ARIMA(0,0,0) with zero mean : -1629.715
## ARIMA(2,0,2)(0,0,1)[12] with non-zero mean : -1655.759
## ARIMA(2,0,2)(2,0,1)[12] with non-zero mean : -1649.472
## ARIMA(2,0,2)(1,0,0)[12] with non-zero mean : -1659.219
```

```
## ARIMA(1,0,2)(1,0,0)[12] with non-zero mean : -1659.184
## ARIMA(3,0,2)(1,0,0)[12] with non-zero mean : -1653.089
## ARIMA(2,0,1)(1,0,0)[12] with non-zero mean : -1653.734
## ARIMA(2,0,3)(1,0,0)[12] with non-zero mean : -1654.633
## ARIMA(1,0,1)(1,0,0)[12] with non-zero mean : -1654.316
## ARIMA(3,0,3)(1,0,0)[12] with non-zero mean : -1648.008
## ARIMA(2,0,2)(1,0,0)[12] with zero mean : -1664.188
## ARIMA(2,0,2) with zero mean : -1663.718
## ARIMA(2,0,2)(2,0,0)[12] with zero mean : -1659.891
## ARIMA(2,0,2)(1,0,1)[12] with zero mean : -1662.4
## ARIMA(2,0,2)(2,0,1)[12] with zero mean : -1654.537
## ARIMA(1,0,2)(1,0,0)[12] with zero mean : -1664.324
## ARIMA(1,0,1)(1,0,0)[12] with zero mean : -1659.476
## ARIMA(1,0,3)(1,0,0)[12] with zero mean : -1665.736
## ARIMA(0,0,2)(1,0,0)[12] with zero mean : -1633.876
## ARIMA(2,0,4)(1,0,0)[12] with zero mean : -1656.216
## ARIMA(1,0,3)(1,0,0)[12] with non-zero mean : -1660.681
## ARIMA(1,0,3) with zero mean : -1663.921
## ARIMA(1,0,3)(2,0,0)[12] with zero mean : -1659.129
## ARIMA(1,0,3)(1,0,1)[12] with zero mean : -1664.133
## ARIMA(1,0,3)(2,0,1)[12] with zero mean : -1653.7
## ARIMA(0,0,3)(1,0,0)[12] with zero mean : -1640.814
## ARIMA(2,0,3)(1,0,0)[12] with zero mean : -1659.646
## ARIMA(1,0,4)(1,0,0)[12] with zero mean : -1660.895
##
## Now re-fitting the best model(s) without approximations...
##
## ARIMA(1,0,3)(1,0,0)[12] with zero mean : -1661.046
##
## Best model: ARIMA(1,0,3)(1,0,0)[12] with zero mean
```

e-1) Examine the significance of the model coefficients and analyze the residuals to check adequacy of the model

BIC() suggests that ARIMA(1,0,3)(1,0,0)[12] with zero mean is the best model.

```
coeftest(m1)
```

```
##
## z test of coefficients:
##
##      Estimate Std. Error z value Pr(>|z|)
## ar1    0.821877   0.061570 13.3486 < 2.2e-16 ***
## ma1   -0.822426   0.083459 -9.8542 < 2.2e-16 ***
## ma2    0.138905   0.081621  1.7018  0.088787 .
## ma3    0.182975   0.068459  2.6728  0.007523 **
## sar1  -0.106368   0.068149 -1.5608  0.118566
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

By coefficient test, ma2 and sar1 have pretty high p-values for t-test and we could drop those two values. However, I am not able to see Intercept by using ARIMA(1,0,3)(1,0,0)

Thus, I will check arima(1,0,3) instead if this function shows the Intercept.



```

m11 = Arima(ratets, c(1,0,3), method = 'ML')
coeftest(m11)

##
## z test of coefficients:
##
##           Estimate Std. Error z value Pr(>|z|)
## ar1      0.81457954  0.06140278 13.2662 < 2.2e-16 ***
## ma1     -0.81045702  0.08385289 -9.6652 < 2.2e-16 ***
## ma2      0.12436064  0.07992196  1.5560  0.119702
## ma3      0.18883096  0.06627749  2.8491  0.004384 **
## intercept 0.00086275  0.00104399  0.8264  0.408581
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

I am now able to check intercept, but its p-value is so high as well as the p-value for MA2. Hence, I could drop these two, but I will put 0 by using fixed() in Arima.

Instead of drop them, I put 0 by using fixed option.

```

m111 = Arima(ratets, c(1,0,3), fixed = c(NA,NA,0,NA,0), method = 'ML')
coeftest(m111)

##
## z test of coefficients:
##
##           Estimate Std. Error z value Pr(>|z|)
## ar1    0.819617    0.059892  13.6849 < 2.2e-16 ***
## ma1   -0.753919    0.067651 -11.1442 < 2.2e-16 ***
## ma3    0.251288    0.055065   4.5634 5.032e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Finally, the t-tests are good.

Now, I am going to check adequacy of the models.

```

# Box test for m1$residuals
Box.test(m1$residuals, lag = 4, type = 'Ljung')

```

```

##
## Box-Ljung test
##
## data:  m1$residuals
## X-squared = 0.48873, df = 4, p-value = 0.9746
Box.test(m1$residuals, lag = 6, type = 'Ljung')

```

```

##
## Box-Ljung test
##
## data:  m1$residuals
## X-squared = 1.095, df = 6, p-value = 0.9818

```

Ljung-Box test on m1\$residuals cannot reject hypothesis of white noise time series

```

# Box test for m11$residuals
Box.test(m11$residuals, lag = 4, type = 'Ljung')

```

```
##
## Box-Ljung test
##
## data: m1$residuals
## X-squared = 0.7197, df = 4, p-value = 0.9489
Box.test(m1$residuals, lag = 6, type = 'Ljung')
```

```
##
## Box-Ljung test
##
## data: m1$residuals
## X-squared = 1.2721, df = 6, p-value = 0.9732
```

The p-values for Ljung-Box test on m1\$residuals are still high, so cannot reject hypothesis of white noise time series

```
# Box test for m11$residuals
Box.test(m11$residuals, lag = 4, type = 'Ljung')
```

```
##
## Box-Ljung test
##
## data: m11$residuals
## X-squared = 4.1008, df = 4, p-value = 0.3925
Box.test(m11$residuals, lag = 6, type = 'Ljung')
```

```
##
## Box-Ljung test
##
## data: m11$residuals
## X-squared = 4.8653, df = 6, p-value = 0.5612
```

Ljung-Box test on m1\_1\$residuals get improved but still cannot reject hypothesis of white noise time series. Therefore,

## e-2) Write down the model expression and discuss if this is a good model for the data

For m1 model, I have not learned expressing the models with (sar1), also this model does not have an intercept. so I conclude with this model cannot be expressed based on lecture we have covered.

For m11 model, the model expression could be  $X_t - 0.00086275 = 0.81457954(X_{t-1} - 0.00086275) + a_t - 0.81045702(a_{t-1}) + 0.12436064(a_{t-2}) + 0.18883096(a_{t-3})$

This model is not a good model for this analysis. Since it does not have a white noise, and I had some bad coefficients results.

f) Find an alternative model(M2) for the index growth rate time series, using either an MA(q) or an AR(p) model depending on your analysis in d). Make sure the model coefficients are significant and the residuals are white noise. Write down the model expression.

```
# alternative model AR(4) since the PACF cuts off after lag 4
m2<-Arima(ratets, c(4,0,0), method = 'ML')
coefest(m2)
```

```
##
## z test of coefficients:
##
##           Estimate Std. Error z value Pr(>|z|)
## ar1      0.04422845 0.06497972  0.6807 0.4960930
## ar2      0.15370358 0.06297554  2.4407 0.0146594 *
## ar3      0.23860823 0.06299084  3.7880 0.0001519 ***
## ar4      0.20093968 0.06480786  3.1005 0.0019316 **
## intercept 0.00087412 0.00106987  0.8170 0.4139121
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

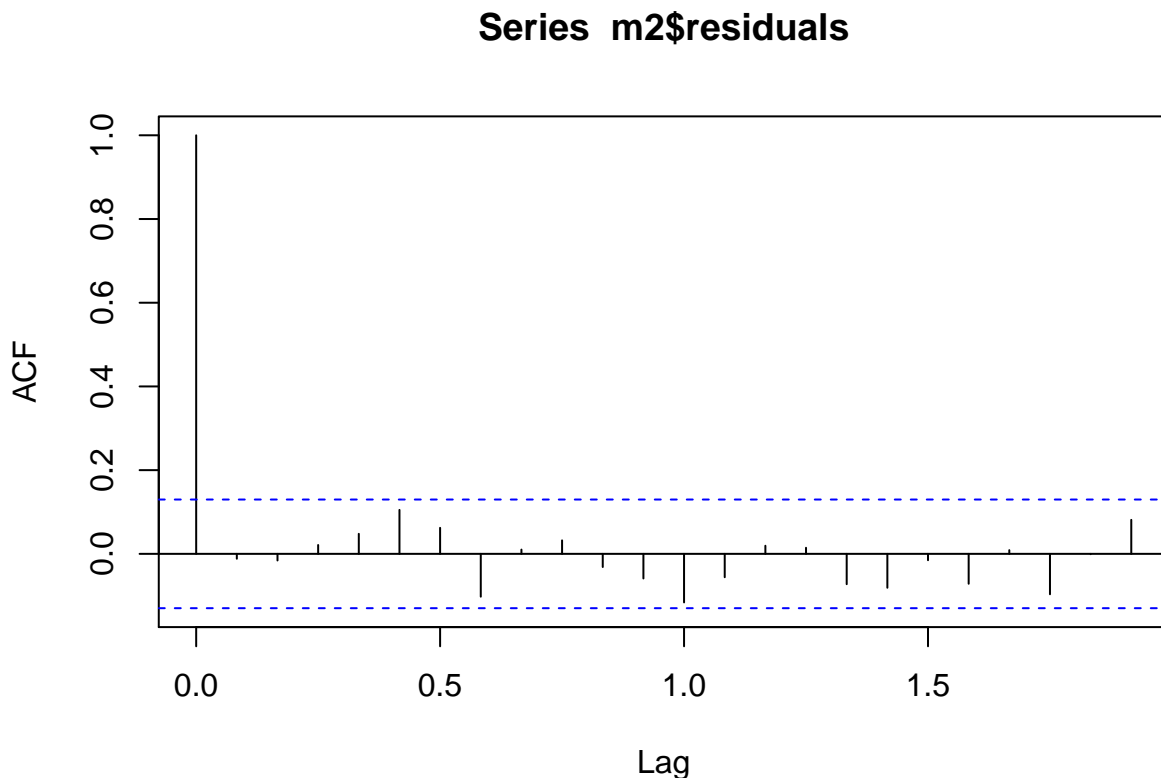
coefest() shows that t-test for ar2,ar3,ar4 are very good, but failed for ar1 and intercept. so we could drop them

```
Box.test(m2$residuals, lag = 4, fitdf = 4, type = 'Ljung')
```

```
##
## Box-Ljung test
##
## data:  m2$residuals
## X-squared = 0.73156, df = 0, p-value < 2.2e-16
```

I got 2.2e-16 for p-value which is good.

```
acf(m2$residuals)
```



acf the residuals looks pretty good, thus white noise.

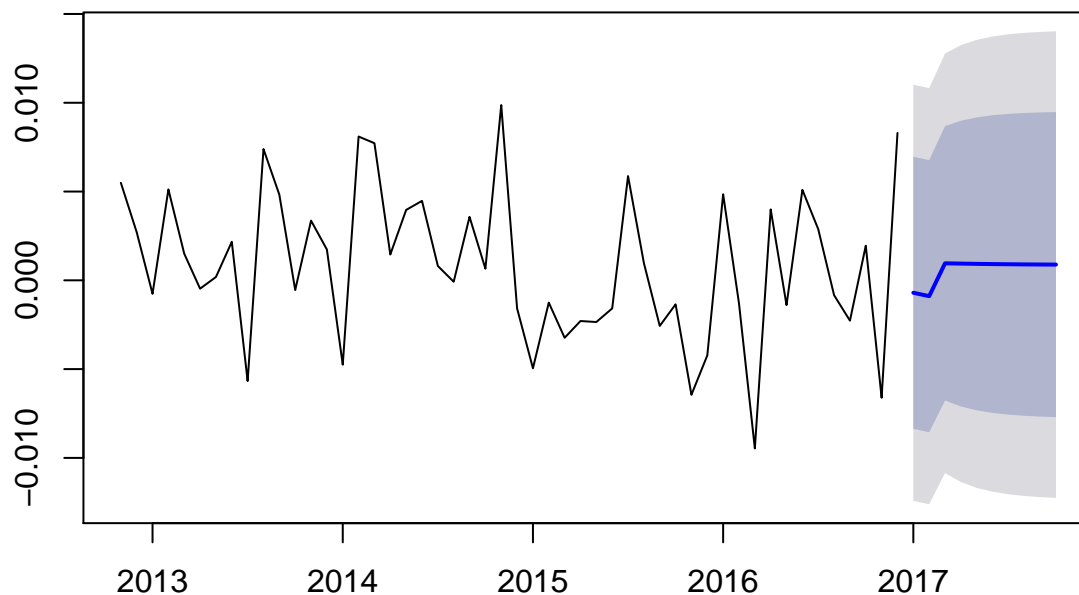
the model expression :

$$x_t - 0.00087412 = 0.04422845(x_{t-1} - 0.00087412) + 0.15370358(x_{t-2} - 0.00087412) + 0.23860823(x_{t-3} - 0.00087412) + 0.20093968(x_{t-4} - 0.00087412) + a_t$$

g) Plot the forecasts for the model M1 and M2 and compare their behavior. Are the forecasts similar in trends? Do you notice any major difference?

```
pred1 = forecast(m1)
plot(forecast(m1, h = 10), include = 50)
```

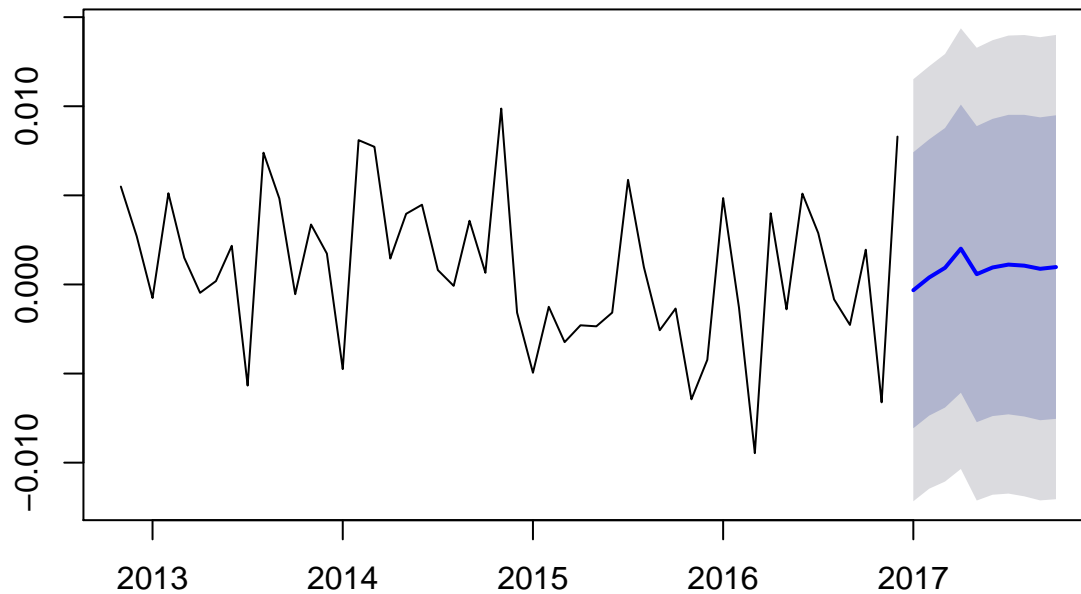
### Forecasts from ARIMA(1,0,3) with non-zero mean



This trend shows there would be a flat trend after 2017 and this forecasting answers me with no interesting patterns.

```
#alternative model
plot(forecast(m2, h = 10), include = 50)
```

## Forecasts from ARIMA(4,0,0) with non-zero mean



Based on the alternative model forecasting, there will be some gentle fluctuations after 2017 with a soft-positive trend. Also the changes around values of 0.000.

h) Apply the backtesting procedure using 85% (=194 values) of the data for training and 15% for testing to evaluate the forecasting power for both models.

```
source('backtest.R')
pm1 = backtest(m11, ratets, 194, 1)

## [1] "RMSE of out-of-sample forecasts"
## [1] 0.004749266
## [1] "Mean absolute error of out-of-sample forecasts"
## [1] 0.003861951
## [1] "Mean Absolute Percentage error"
## [1] 2.573491
## [1] "Symmetric Mean Absolute Percentage error"
## [1] 1.478262

pm2 = backtest(m2, ratets, 194, 1)

## [1] "RMSE of out-of-sample forecasts"
## [1] 0.004635424
## [1] "Mean absolute error of out-of-sample forecasts"
## [1] 0.003664983
## [1] "Mean Absolute Percentage error"
## [1] 2.00443
## [1] "Symmetric Mean Absolute Percentage error"
## [1] 1.411719
```

Based on the results, both models are pretty good since all RMSE, MAE, MAPE, SMAPE are very low.

RMSE for m1 is 0.004749266 and m2 is 0.004635424. Thus m2 has more accuracy MAE for m1 is 0.003861951 and m2 is 0.003664983. Thus m2 has more accuracy MAPE for m1 is 2.573491 and m2 is 2.00443. Thus m2

has more accuracy SMAPE for m1 is 1.478262 and m2 is 1.411719 Thus m2 has more accuracy

**i) At the end of this analysis, which model would you recommend for this analysis? Discuss the trend captured by it.**

I think model 2 is more adequate for this analysis since it is more accurate based on the result from backtest. Also, it has a white noise, and shows a good result from `coeftest()`. Lastly, based on `forecast()`, the trend for model 2 shows some positive trend. In these reasons, model 2 is a better model in this analysis.