# Data Analysis – Lab 6

Victor Lebrun – Fanny Streiff

In this report the source code will be written in blue and the results in black
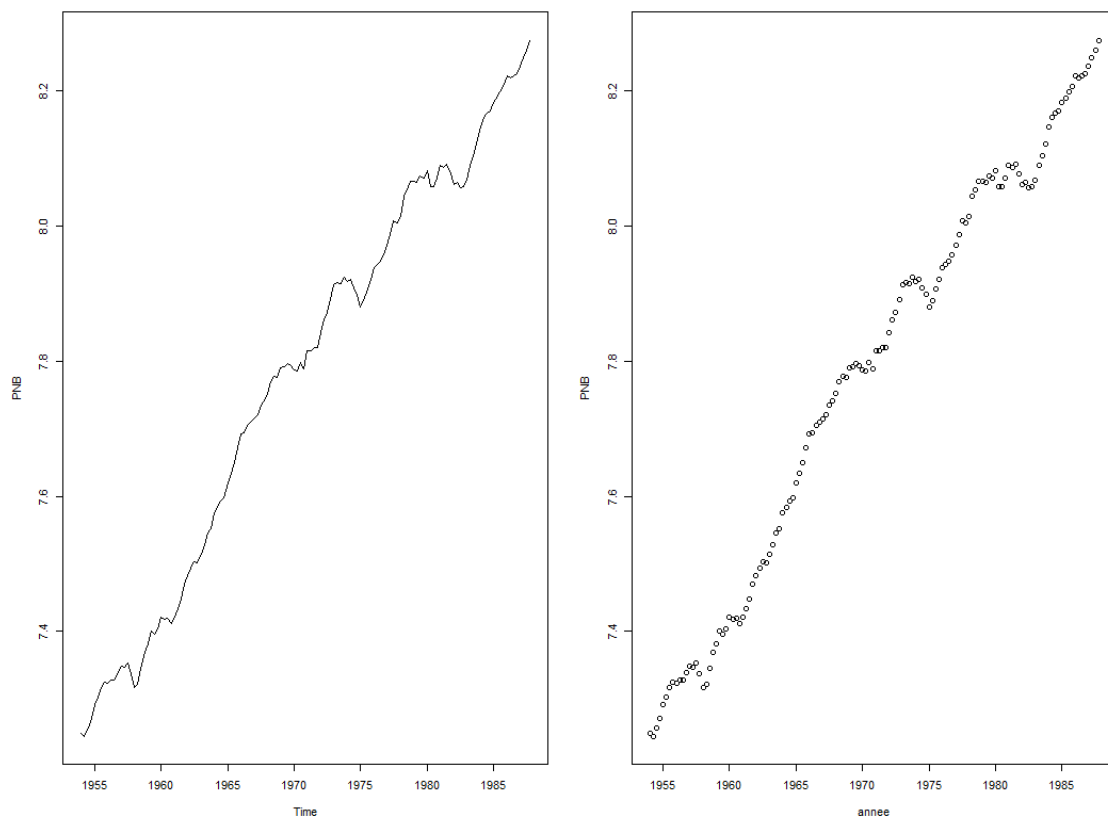
## A Stationarity analysis

1. Install and load the tseries library, then load the USeconomic data using the command line data(USeconomic).

```
library("tseries")
data(USeconomic)
```

2. We will now create the variables that we are going to study:

• The log(GNP) data are in the 2nd column of the USeconomic dataset. Retrieve this column and store it in a variable of your choice.

• The log(GNP) data were acquired on a trimestrial basis. Use the command seq to create a second variable year containing all years and semesters from 1954 to 1987.75 .

• Display the log(GNP) evolution between 1954 and the 3rd trimester of 1987.

```
PNB=USeconomic[,2]
plot(PNB)
annee=seq(from=1954 ,to=1987.75, by=0.25)
plot(annee,PNB)
```

**3. Remind what "stationarity" means for a time series. What can you visually say about the stationarity of the log(GNP) time series.**
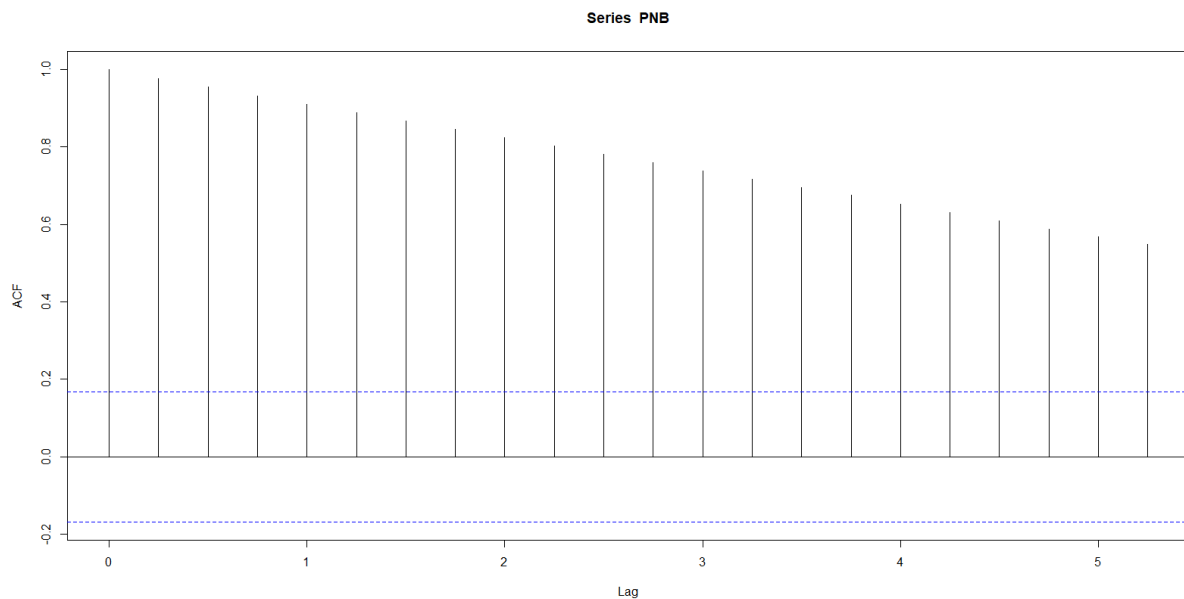
"Strict Stationarity : A time series is said to be strictly stationary if all its observations are drawn from the same distribution: the joint probability does not change in time."

"Weak Stationarity : We do not require that each draw comes from the exact same distribution, only that the distributions have the same mean and variance (all of them not a function of time).

- Constant mean: $E(x_t) = \mu$
- Constant variance: $Var(x_t) = \gamma_0$
- Constant co-variance: $Cov(x_t, x_{t-h}) = \gamma_h \ \forall h \in [1..T]$"

The log(GNP) function does not look to be stationary as it takes increasing values, its mean value and its variance are not constant.

## 4. Use the acf command to draw this series autocorrelogram. Comment.

Series PNB



First of all, just as the exemple we saw in class : the ACF does not go down exponentially which means that the series may need to be differentiated. And so this is what we are going to do in the following questions.

## 5. Use the Box.test function on your series and interpret the result.

We use the Box-Pierce test in order to assess if the series is mostly white noise.

```
> Box.test(PNB)

        Box-Pierce test

data:  PNB
X-squared = 129.98, df = 1, p-value < 2.2e-16
```

We observe that the p value for this test is very close to zero. Therefore we can reject the null hypothesis and conclude that this series is not composed white noise. This implies that this time series is not stationary.

## 6. What can you conclude on this time series stationarity ?

According to the previous results we conclude that this time series is not stationary.
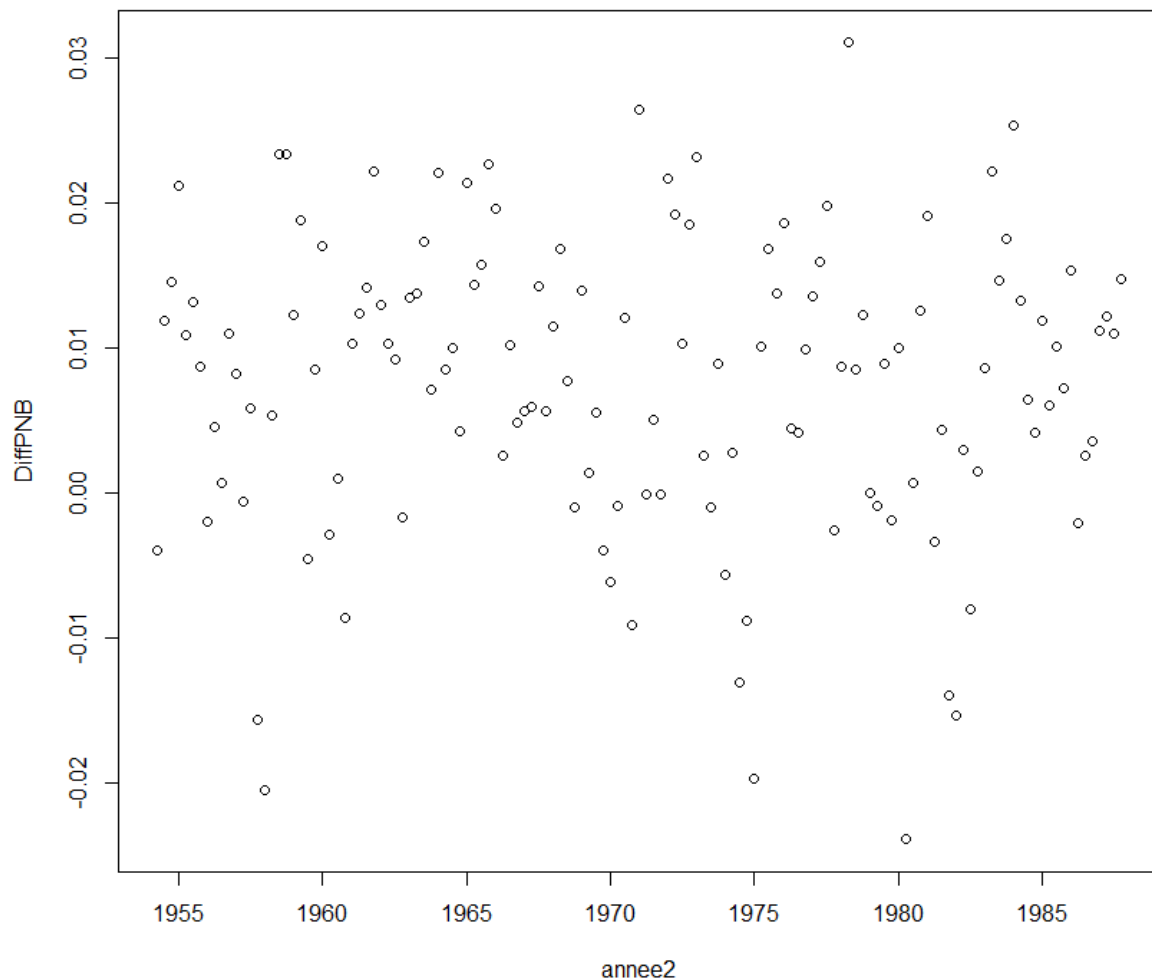
## B Study of diGNP

1. Create a variable DiGNP so that: ($Y_t = X_t - X_{t-1}$)$2 \leq t \leq T$ where $X_t$ is the logGNP at time t. Explain what this time series represents.

```
DiffPNB=diff(PNB)
```

This new time serie represents the difference of GDP between two quarters and t-1 for any t between Q2 of 1954 and Q3 of 1987.

2. Plot the evolution of this series between 1954 and the 3rd semester of 1987.

```
annee2=seq(from=1954.25 ,to=1987.75, by=0.25)
plot(annee2,DiffPNB)
```



We cannot see any clear tendency from this plot.

3. Is this series centered ? You can use the empirical mean value of the series and a Student test to justify your answer.

In order to say if this series is centered or not we compute a Student Test that will also provide us with the mean of the series.
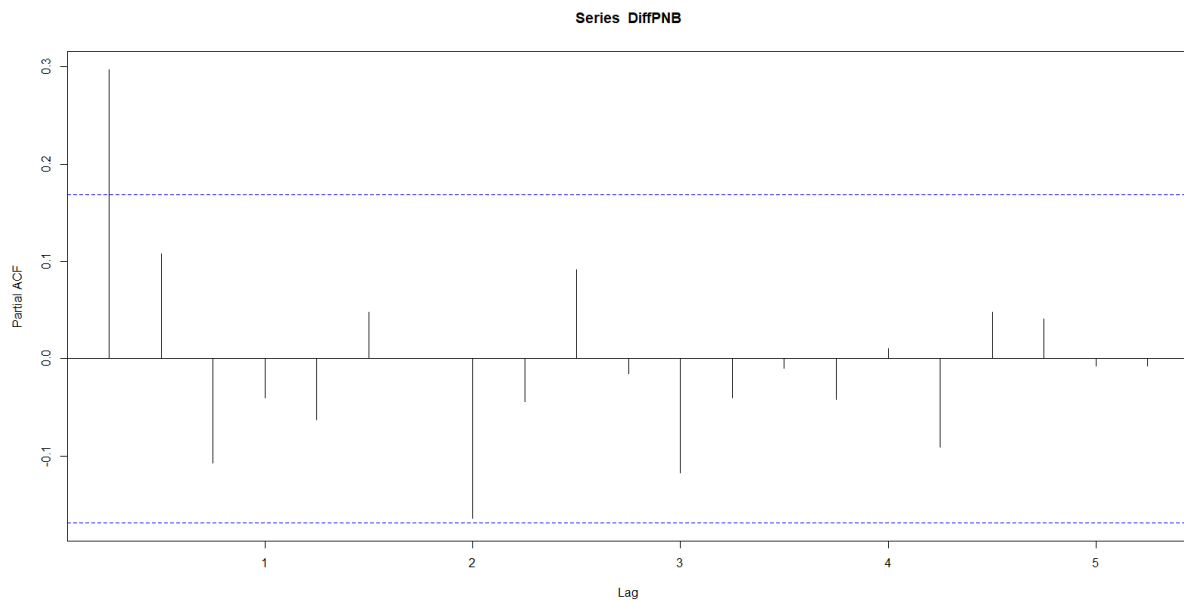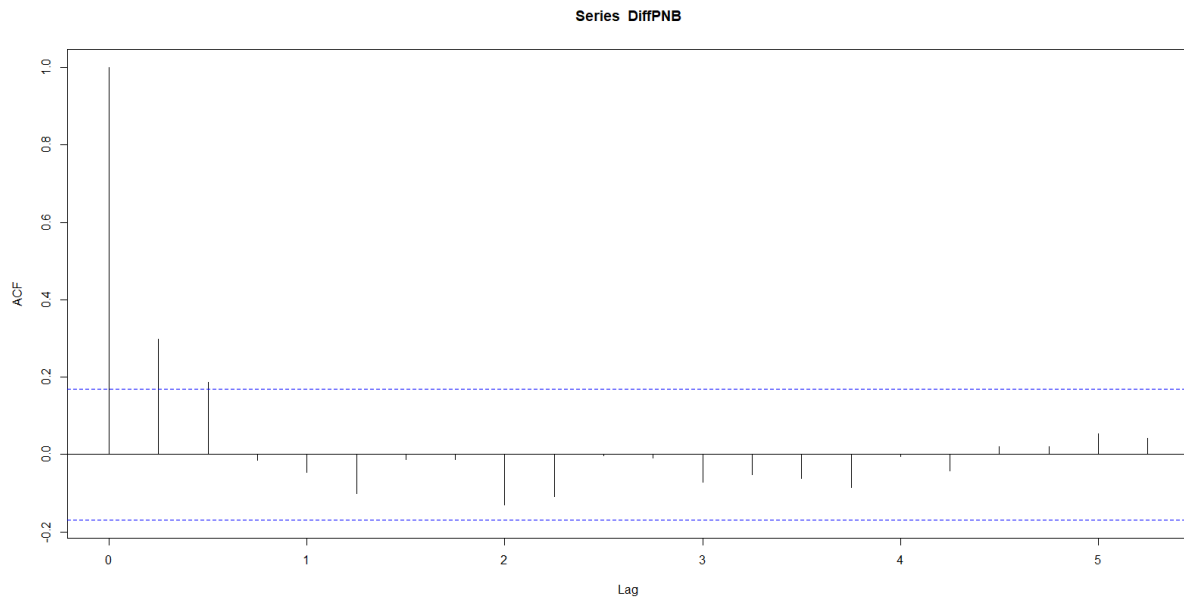
```
> t.test(DiffPNB)

        One Sample t-test

data:  DiffPNB
t = 8.6739, df = 134, p-value = 1.223e-14
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 0.005864407 0.009328764
sample estimates:
  mean of x
0.007596586
```

The mean value of this time series is very close to zero. But the p value of the Student Test is is almost zero, this means that we can reject the null hypothesis. Therefore the true mean value is different from zero and the series is not centered.

4. Use the acf and pacf functions to draw the autocorrelogram and partial autocorrelogram of this time series. From there, deduce the most likely parameter(s) p and q for an ARMA(p,q) model to modelize DiGNP.

```
acf(DiffPNB)
pacf(DiffPNB)
```

Series DiffPNB



Series DiffPNB

From the acf, it appears that q=2, or maybe q=1.

From the pacf, it appears that p=0 or maybe p=8.

5. Test all the couples (p,q) that seemed relevant to you from the previous question:

• The function arima(DiGNP,c(p,0,q)) will help you to evaluate each model.

• Explain the different values returned by the function arima. In particular, you will have to retrieve the fitted parameters and explain the meaning of the two following parameters: log likelihood" and aic" (Akaike's Information Criterion), and how you can you them to rate your model.

• Which model seems to be the best one if you only use the AIC criterion ? We are now going to study 3 models in particular: ARMA(0,1), ARMA(0,2) and ARMA(8,2). We remind you that the Shapiro-Wilk test assesses the null hypothesis that a sample follows a normal distribution.

We will test 3 couples of values for (p,q) : (0,2), (0,1), (8,2).

```
A1=arima(DiffPNB,c(0,0,2))
A2=arima(DiffPNB,c(0,0,1))
A3=arima(DiffPNB,c(8,0,2))

> A1

Call:
arima(x = DiffPNB, order = c(0, 0, 2))

Coefficients:
         ma1      ma2   intercept
      0.2681   0.1976     0.0076
s.e.  0.0851   0.0790     0.0012

sigma^2 estimated as 9.178e-05:  log likelihood = 435.86,  aic = -863.73

> A2

Call:
arima(x = DiffPNB, order = c(0, 0, 1))

Coefficients:
         ma1   intercept
      0.2278     0.0076
s.e.  0.0707     0.0010

sigma^2 estimated as 9.578e-05:  log likelihood = 433.03,  aic = -860.05

> A3

Call:
arima(x = DiffPNB, order = c(8, 0, 2))

Coefficients:
        ar1      ar2     ar3     ar4      ar5     ar6      ar7      ar8      ma1     ma2  intercept
     0.3857  -0.4613  0.0457  0.0523  -0.1347  0.0604  -0.0039  -0.1393  -0.1172  0.6047    0.0075
s.e. 0.3313   0.2883  0.1180  0.1094   0.1005  0.1047   0.1195   0.0975   0.3312  0.2354    0.0010

sigma^2 estimated as 8.646e-05:  log likelihood = 439.66,  aic = -855.32
```

Call : remind the parameters

Coefficients : test the model for the p-ar and the q-ma.

log likelihood : it is used to assess the quality of fit of two statistical models. The test is based on the likelihood ratio, which expresses how many times more likely the data are under one model than the other. The higher this coefficient is, the better it is.

aic : the AIC (Akaike information criterion) coefficient is an estimator of the relative quality of statistical models for a given set of data and can be interpreted as an estimate of how

much information would be lost if a given model is chosen. When comparing models, one wants to minimize AIC. So the lower it is, the better it is. The smallest the number is, the "simpler" is the model so the better it is.

Looking at the log likelihood results, we select first (8,2) and (0,2) because (8,2) has a greater log likelihood than (0,2). But if now look at aic, we will select (0,2) because of its small coefficient.

Finally, we conclude that (0,2) seems to be the better model due to its highest log likelihood correlation coefficient and to its small aic coefficient.

6. Use the Box-Pierce test and the Shapiro-Wilk test (shapiro.test) on the residuals of all 3 models applied to the log(GNP) data and display their autocorrelogram. From there, what can you say on the stationarity of the residuals ? How do you justify which model is the best ?

At first we will test the residuals of the (p=0,q=2) model.

```
> Box.test(A1$residuals)

        Box-Pierce test

data:  A1$residuals
X-squared = 0.0043454, df = 1, p-value = 0.9474
```

We test the hypothesis that the mean of the series of the residual is null. We observe that the p-value is higher than 0.05 so we cannot reject our initial hypothesis.
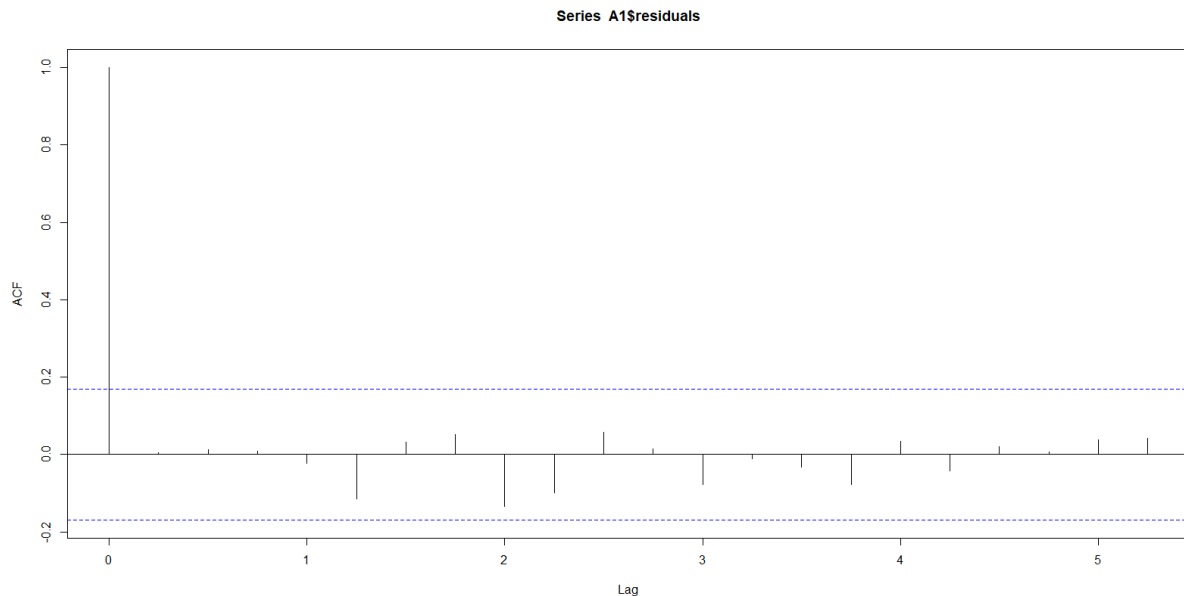
```
> shapiro.test(A1$residuals)

        Shapiro-Wilk normality test

data:  A1$residuals
W = 0.98847, p-value = 0.3231

> acf(A1$residuals)
> mean(A1$residuals)
[1] 3.271741e-05
~ |
```

The Shapiro–Wilk test assesses the null hypothesis that the serie of the residuals came from abnormally distributed population. In that case p > 0.05, so we cannot reject the null hypothesis.

Series A1$residuals

We are now testing the residuals of the (p=0,q=1) model.

```
> Box.test(A2$residuals)

        Box-Pierce test

data:  A2$residuals
X-squared = 0.28287, df = 1, p-value = 0.5948
```

We test the hypothesis that the mean of the series of the residual is null. We observe that the p-value is higher than 0.05 so we cannot reject our initial hypothesis.

```
> shapiro.test(A2$residuals)

        Shapiro-Wilk normality test

data:  A2$residuals
W = 0.97756, p-value = 0.02493

> acf(A2$residuals)
> mean(A2$residuals)
[1] 1.80625e-05
```
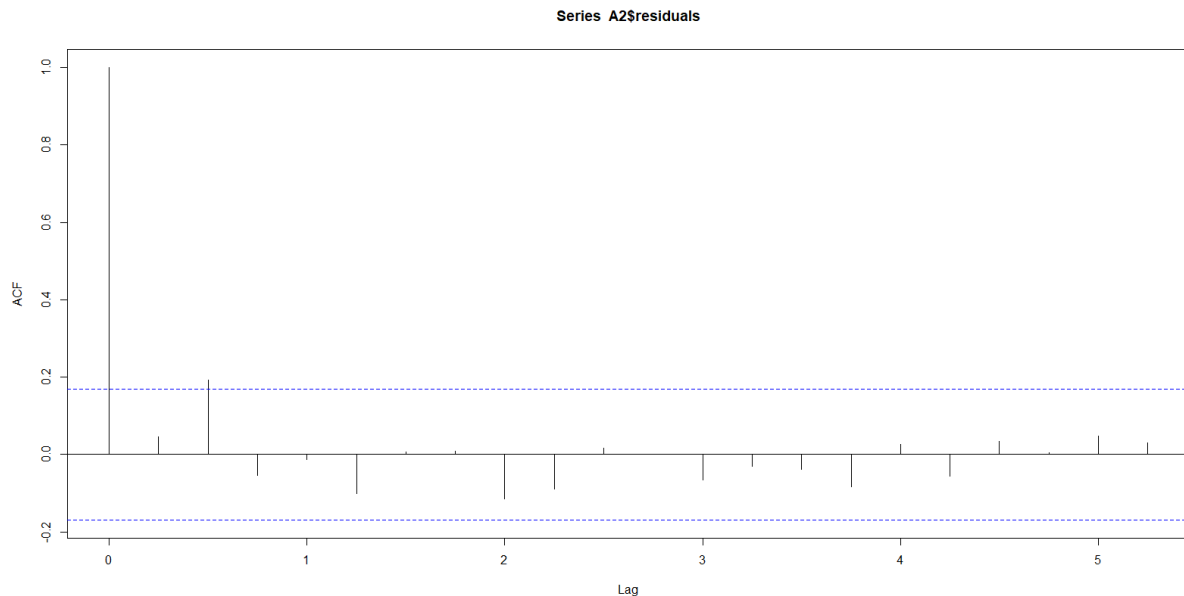
The Shapiro–Wilk test assesses the null hypothesis that the serie of the residuals came from abnormally distributed population. In that case $p < 0.05$, so we can reject the null hypothesis.

Series A2$residuals

Lastly we test the residuals of the (p=8,q=2) model.

```
> Box.test(A3$residuals)

        Box-Pierce test

data:  A3$residuals
X-squared = 1.3959e-07, df = 1, p-value = 0.9997
```

We test the hypothesis that the mean of the series of the residual is null. We observe that the p-value is higher than 0.05 so we cannot reject our initial hypothesis.
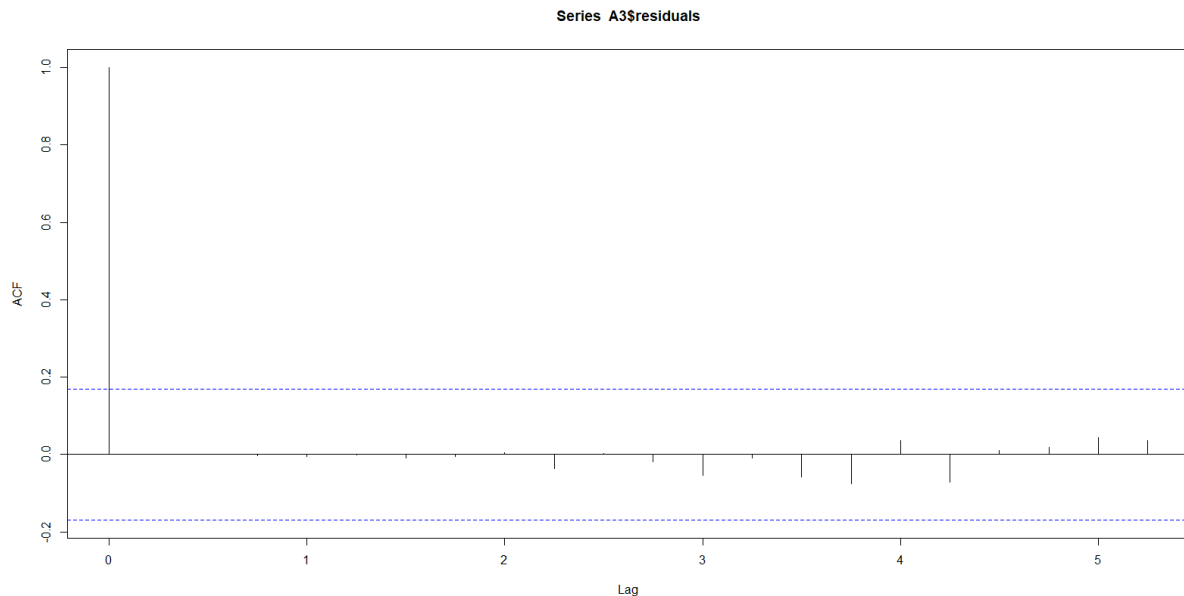
```
> shapiro.test(A3$residuals)

        Shapiro-Wilk normality test

data:  A3$residuals
W = 0.98606, p-value = 0.1872

> acf(A3$residuals)
> mean(A3$residuals)
[1] 3.89988e-05
```

The Shapiro–Wilk test assesses the null hypothesis that the serie of the residuals came from abnormally distributed population. In that case p > 0.05, so we cannot reject the null hypothesis.

Series A3$residuals

To define if a model is good or not looking at its residuals, we need this residuals to follow 4 rules :

- "The residuals are uncorrelated. If there are correlations between residuals, then there is information left in the residuals which should be used in computing forecasts.
- The residuals have zero mean. If the residuals have a mean other than zero, then the forecasts are biased.
- The residuals have constant variance.
- The residuals are normally distributed."

*(source : https://otexts.org/fpp2/residuals.html)*

Looking at the acf, we can see that (0,2) and (8,2) seems to be uncorrelated because "the model goes into the box and doesn't go outside after". Which is not the case of (0,1) : we can see a line going out of the box at q=2 while q=1 is in the box.

For the 3 residuals, the mean is very close to 0, which is very good but doesn't help to choose one.

For the 3 Box Pierce's test, the p-value is far beyond 0.05 which is very good because we want the residuals to be white noise. However, once again, (0,2) and (0,8) is better than (0,1) because theirs p-values is near to 1 when (0,1) p-values is close to 0.6 only.

For the Shapiro test, we need to have a high p-value (>0.05) to conclude that the residual is following a normal distribution, which is only the case for the case (0,2) and (0,8).

Finally, we can keep the model (0,2) which seems the best (especially because of Q5). But the model (8,2) could also be a less correct, but still coherent answer too as far as we can go in our analysis.
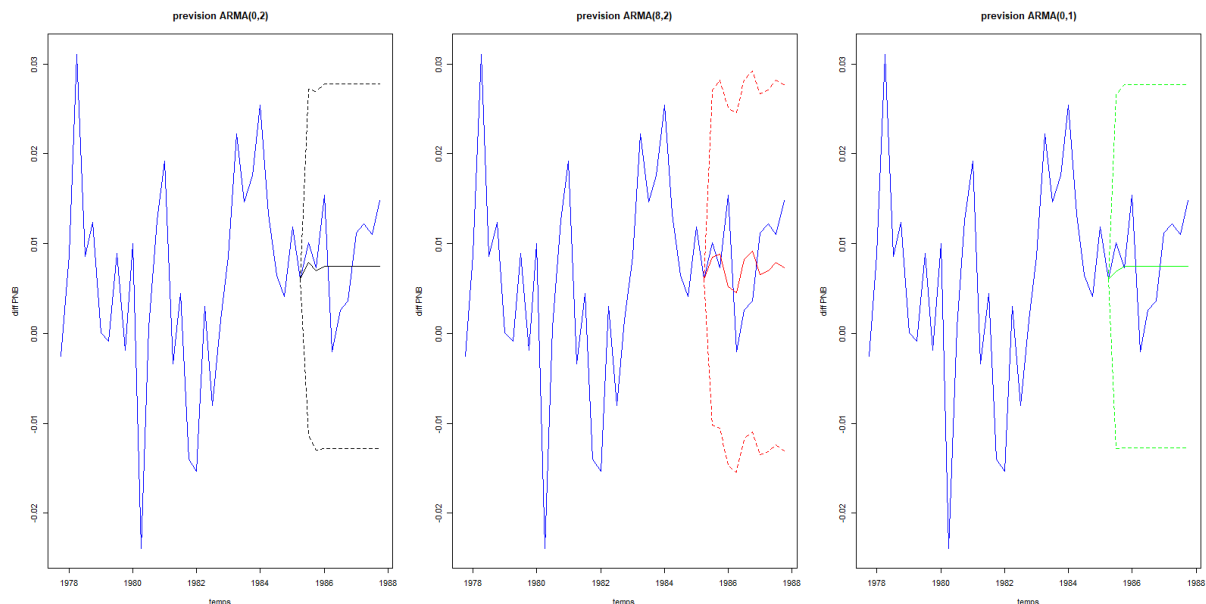
# C Predictions using ARMA

1. For each model, draw on the same graph the 3 following elements (you may want to zoom on the end of the time series):

• The original GNP series

• The mean value predicted by each model.

• The  95% confidence interval on your prediction under the hypothesis that the residuals are normally distributed (regardless of question B.6 result). See the code below to answer this question (you may have to adapt it to your variables).


We used the code that was given in the subject in order to make the following plot.

```
n <- 10
T=length(annee)
index <- 1:(T - n - 1)
res01 <- predict(arima(DiffPNB[index], c(0, 0, 1)), n)
res02 <- predict(arima(DiffPNB[index], c(0, 0, 2)), n)
res82 <- predict(arima(DiffPNB[index], c(8, 0, 2)), n)

par(mfrow = c(1, 3))
plot(annee[(T - 4 * n):T], DiffPNB[(T - 4 * n):T - 1], main = "prevision ARMA(0,2)", t = "l", col = "blue", xlab = "temps", ylab = "diff PNB")
lines(annee[(T - n):T], c(DiffPNB[T - n - 1], res02$pred))
lines(annee[(T - n):T], c(DiffPNB[T - n - 1], res02$pred) + c(0,res02$se) * 1.96, lty = 2)
lines(annee[(T - n):T], c(DiffPNB[T - n - 1], res02$pred) - c(0,res02$se) * 1.96, lty = 2)
plot(annee[(T - 4 * n):T], DiffPNB[(T - 4 * n):T - 1], main = "prevision ARMA(8,2)",t = "l", col = "blue", xlab = "temps", ylab = "diff PNB")
lines(annee[(T - n):T], c(DiffPNB[T - n - 1], res82$pred), col = "red")
lines(annee[(T - n):T], c(DiffPNB[T - n - 1], res82$pred) + c(0,res82$se) * 1.96, lty = 2, col = "red")
lines(annee[(T - n):T], c(DiffPNB[T - n - 1], res82$pred) - c(0,res82$se) * 1.96, lty = 2, col = "red")
plot(annee[(T - 4 * n):T], DiffPNB[(T - 4 * n):T - 1], main = "prevision ARMA(0,1)",t = "l", col = "blue", xlab = "temps", ylab = "diff PNB")
lines(annee[(T - n):T], c(DiffPNB[T - n - 1], res01$pred), col = "green")
lines(annee[(T - n):T], c(DiffPNB[T - n - 1], res01$pred) + c(0,res01$se) * 1.96, lty = 2, col = "green")
lines(annee[(T - n):T], c(DiffPNB[T - n - 1], res01$pred) - c(0,res01$se) * 1.96, lty = 2, col = "green")
```



We can observe that the predicted values are relatively far from the real ones. The first and third give results that almost do not evolve and that are very far from the real values. The second model gives a result that contains more evolving values but they are not really closer to the expected ones.

2. Using the previous question, which model seem to give the best results ?

The best model seems to be the ARMA(0,2) as it has the highest log likelihood and the lowest AIC coefficient of the three models we tested. It has a good p-value to the Box-Pierce test and the Shapiro-Wilk test compared to the two other models, so its residuals came from a normally distributed population and they are stationary as the mean value of the series is null. The predicted values are better than the ones of the two other models as the ARMA(8,2) model often predicts trends that are the opposite of the real ones. The residuals of the ARMA(0,1) model were not normalised and the other tests did not give good results.

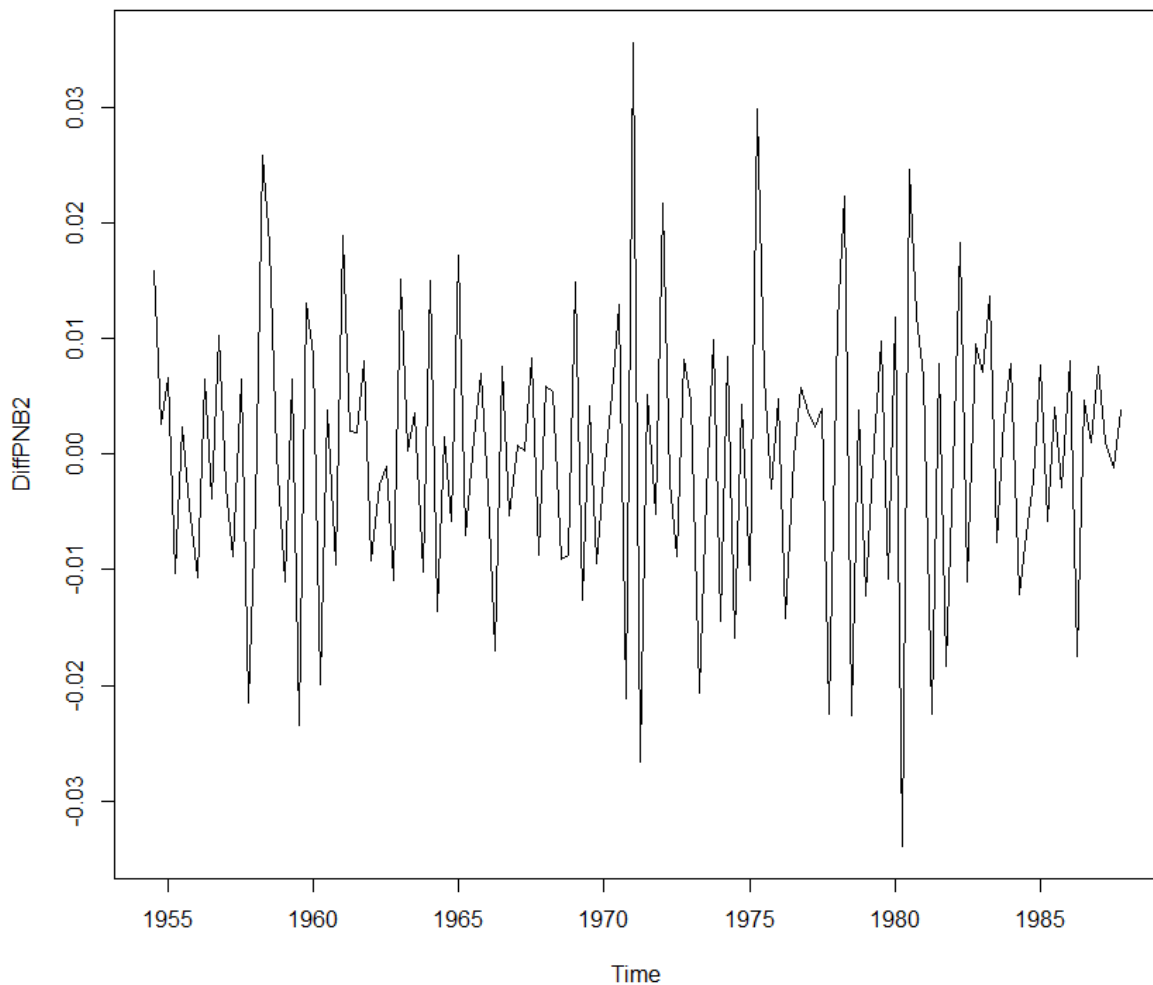## 3. Translate these models into ARIMA(p,d,q) for the original logGNP time series.

We already used the ARIMA model to compute the previous results. We used the ARIMA model on the differentiated values of logGNP with the parameter d=0. This is the same as using the ARIMA model with d=1 on the logGNP time series.

## D ARIMA Model

Differentiate the GNP serie a second time and use the analysis of sections B) and C) to find the best possible ARIMA(p,2,q) possible for the GNP series. Remark : You may directly use the function ARIMA(p,2,q) on the GNP series, it will be more convenient to draw the graphics and make analysis.

We begin by creating the variable DiffPNB2 which is the derivative of DiffPNB.

```
DiffPNB2=diff(DiffPNB)

> plot(DiffPNB2)
```



We compute the Student Test for this time series :
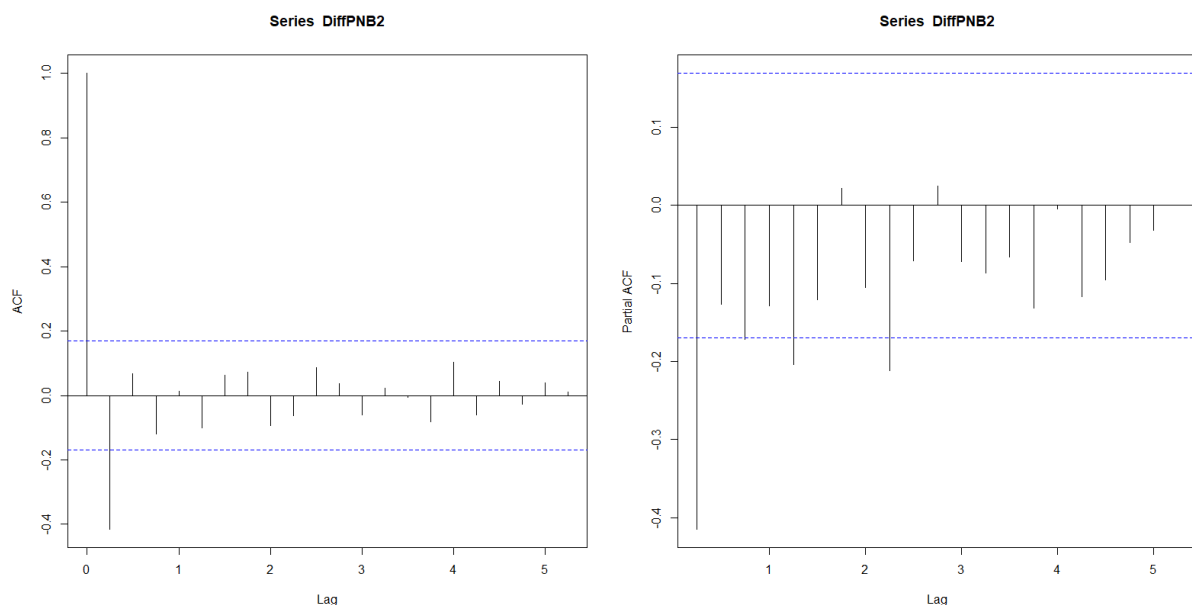
```
> t.test(DiffPNB2)

        One Sample t-test

data:  DiffPNB2
t = 0.13481, df = 133, p-value = 0.893
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 -0.001918837  0.002199523
sample estimates:
    mean of x
0.0001403431
```

The mean value of the time series is very close to zero. The p-value is higher than 0.05 so we cannot reject the null hypothesis, so this time series is centered.

We now take a look at the acf and pacf of DiffPNB2 to try to guess the correct p and q values for the arima model.

```
acf(DiffPNB2)
pacf(DiffPNB2)
```



From acf, we can guess that we will have q=1

From pacf, we can guess that we will have p=0,2,4,8

So now let's try all this 4 models and see which one is the best.

```
A4=arima(PNB,c(0,2,1))
A5=arima(PNB,c(2,2,1))
A6=arima(PNB,c(4,2,1))
A7=arima(PNB,c(8,2,1))
```

We will first take a look at the models themselves :

```
> A4

Call:
arima(x = PNB, order = c(0, 2, 1))

Coefficients:
          ma1
      -1.0000
s.e.   0.0378

sigma^2 estimated as 0.0001035:  log likelihood = 422.17,  aic = -840.34

> A5

Call:
arima(x = PNB, order = c(2, 2, 1))

Coefficients:
         ar1     ar2      ma1
      0.2747  0.1149  -1.0000
s.e.  0.0866  0.0864   0.0216

sigma^2 estimated as 9.32e-05:  log likelihood = 429.64,  aic = -851.28

> A6

Call:
arima(x = PNB, order = c(4, 2, 1))

Coefficients:
         ar1     ar2      ar3      ar4      ma1
      0.2817  0.1464  -0.0933  -0.0357  -1.0000
s.e.  0.0868  0.0894   0.0891   0.0867   0.0233

sigma^2 estimated as 9.19e-05:  log likelihood = 430.43,  aic = -848.87

> A7

Call:
arima(x = PNB, order = c(8, 2, 1))

Coefficients:
         ar1     ar2      ar3      ar4      ar5     ar6     ar7      ar8      ma1
      0.2786  0.1511  -0.0925  -0.0320  -0.0846  0.0730  0.0469  -0.1529  -1.0000
s.e.  0.0856  0.0886   0.0888   0.0895   0.0888  0.0886  0.0880   0.0853   0.0266

sigma^2 estimated as 8.887e-05:  log likelihood = 432.43,  aic = -844.87
```

Looking at the log likelihood results, we select first (8,1), (4,1) and (2,1) because (8,1) has a greater log likelihood than (4,1) and (2,1). But if now look at aic, we will select (2,1) because of its small coefficient.

Finally, we conclude that (4,1) seems to be the best model due to its high log likelihood correlation coefficient and to its small aic coefficient. The (2,1) model also provides very good results and it is hard to say which one of the two is the best.

To be sure of our choice, we will take a look at the residuals.

```
A4=arima(PNB,c(0,2,1))
Box.test(A4$residuals)
shapiro.test(A4$residuals)
acf(A4$residuals)
pacf(A4$residuals)
mean(A4$residuals)

> Box.test(A4$residuals)

        Box-Pierce test

data:   A4$residuals
X-squared = 11.792, df = 1, p-value = 0.0005949

> shapiro.test(A4$residuals)

        Shapiro-Wilk normality test

data:   A4$residuals
W = 0.98388, p-value = 0.1093

> acf(A4$residuals)
> pacf(A4$residuals)
> mean(A4$residuals)
[1] -0.00019578
```
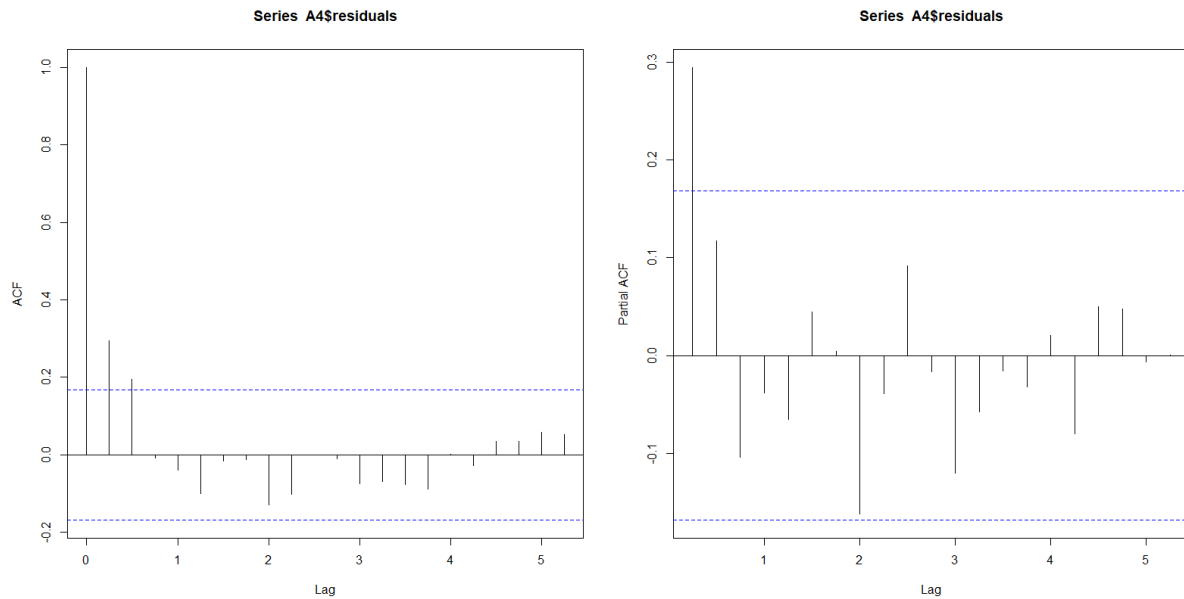
```
A5=arima(PNB,c(2,2,1))
Box.test(A5$residuals)
shapiro.test(A5$residuals)
acf(A5$residuals)
pacf(A5$residuals)
mean(A5$residuals)
```
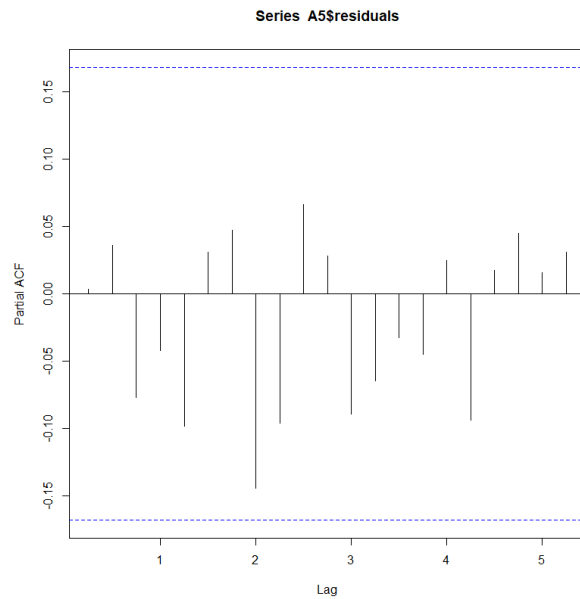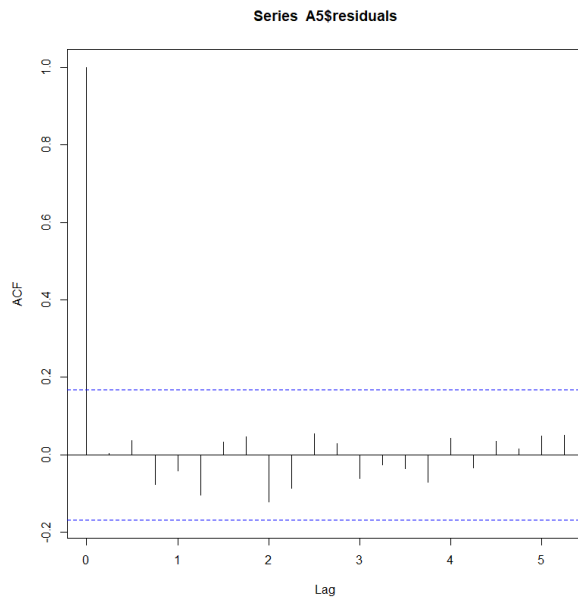
```
> Box.test(A5$residuals)

        Box-Pierce test

data:  A5$residuals
X-squared = 0.0014778, df = 1, p-value = 0.9693

> shapiro.test(A5$residuals)

        Shapiro-Wilk normality test

data:  A5$residuals
W = 0.98934, p-value = 0.3829

> acf(A5$residuals)
> pacf(A5$residuals)
> mean(A5$residuals)
[1] 3.385963e-05
```

ACF

Partial ACF

Lag

Lag

```
A6=arima(PNB,c(4,2,1))
Box.test(A6$residuals)
shapiro.test(A6$residuals)
acf(A6$residuals)
pacf(A6$residuals)
mean(A6$residuals)

> Box.test(A6$residuals)

        Box-Pierce test

data:  A6$residuals
X-squared = 0.018541, df = 1, p-value = 0.8917

> shapiro.test(A6$residuals)

        Shapiro-Wilk normality test

data:  A6$residuals
W = 0.99043, p-value = 0.4782

> acf(A6$residuals)
> pacf(A6$residuals)
> mean(A6$residuals)
[1] 1.186934e-05
```
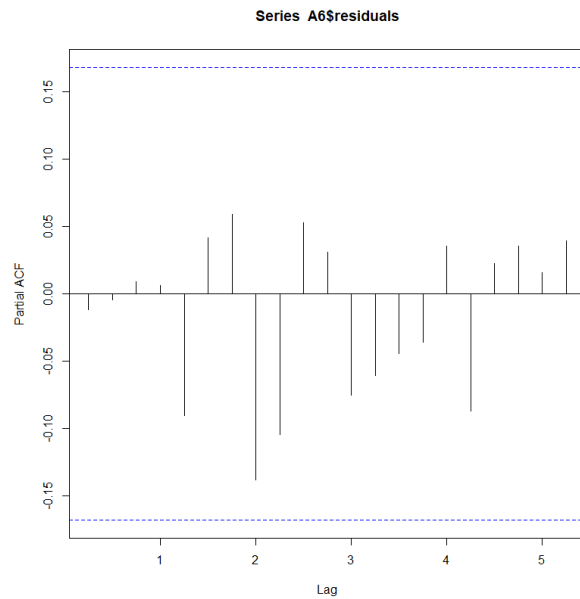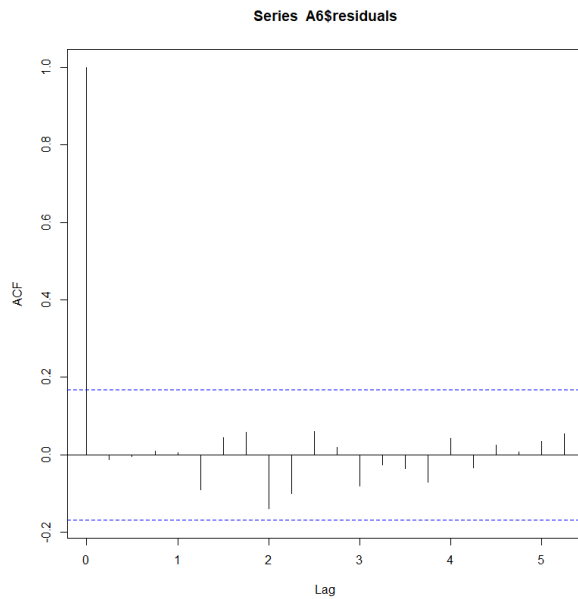
Series A6$residuals | Series A6$residuals

```
A7=arima(PNB,c(8,2,1))
Box.test(A7$residuals)
shapiro.test(A7$residuals)
acf(A7$residuals)
pacf(A7$residuals)
mean(A7$residuals)

> Box.test(A7$residuals)

        Box-Pierce test

data:  A7$residuals
X-squared = 0.023344, df = 1, p-value = 0.8786

> shapiro.test(A7$residuals)

        Shapiro-Wilk normality test

data:  A7$residuals
W = 0.98848, p-value = 0.3188

> acf(A7$residuals)
> pacf(A7$residuals)
> mean(A7$residuals)
[1] -5.09822e-06
```
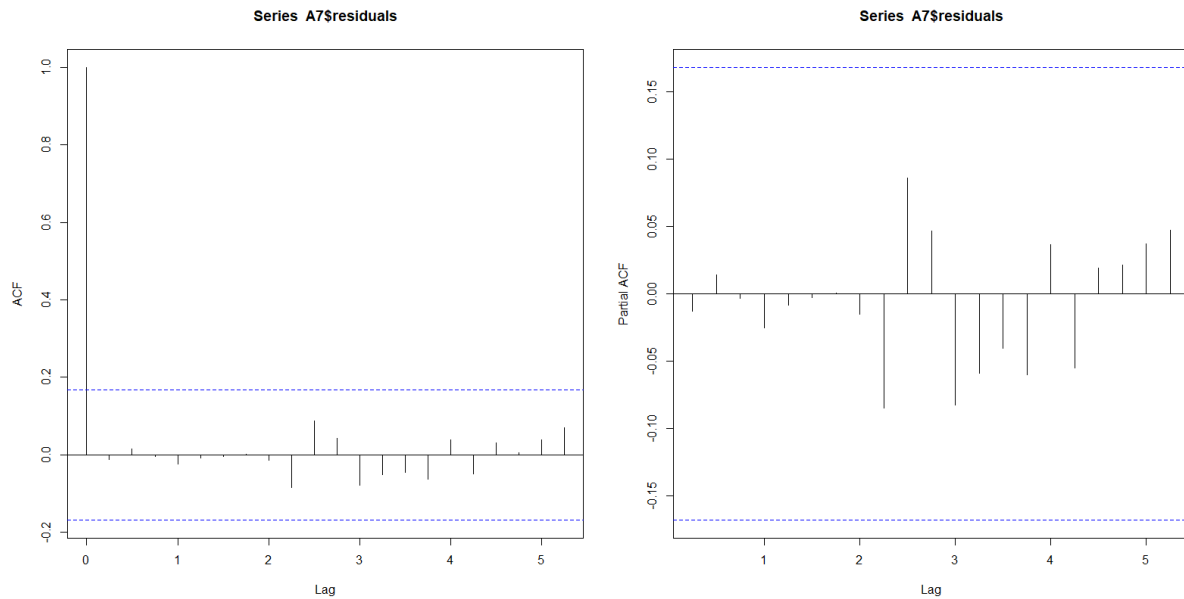
About the four Box Pierce's tests, the p-value is far beyond 0.05 for (2,1), (4,1) and (8,1) which is very good because we want the residuals to be white noise, meaning that the residuals have a null mean. These three models have a mean value around zero.

For the Shapiro test, we need to have a high p-value (>0.05) to conclude that the residual is following a normal distribution, which is the case for every model. However, (4,1) as the highest p-value : 0,4782.
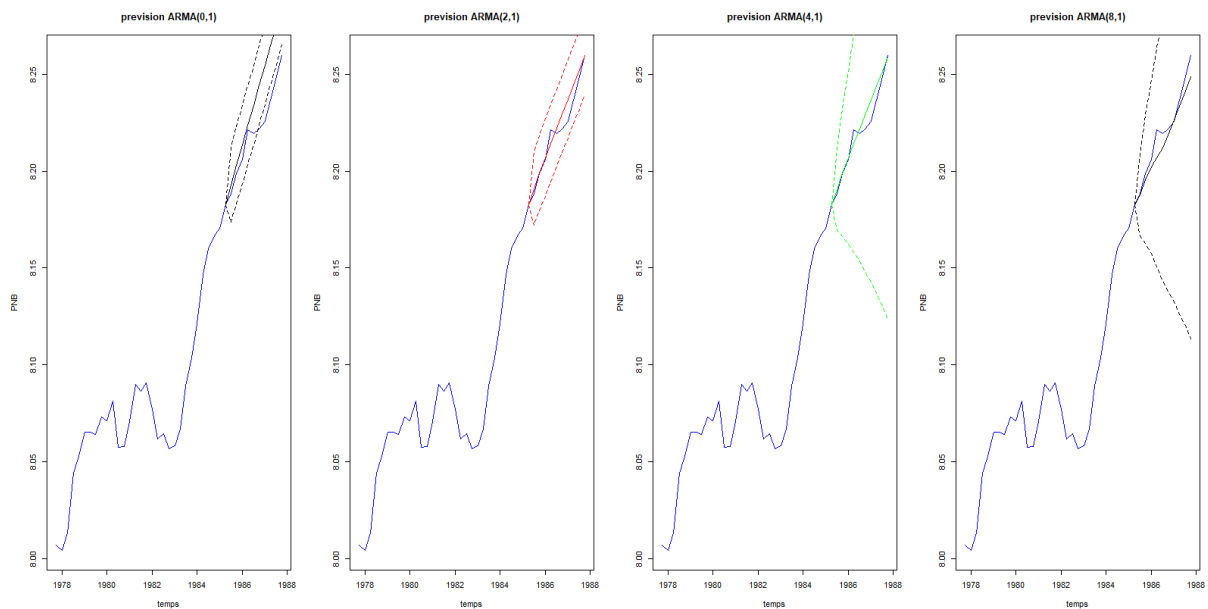
Looking at the acf, we can see that (2,1), (4,1) and (8,1) seems to be uncorrelated because "the model goes into the box and doesn't go outside after".

Finally, we can keep the model (2,1) which seems the best (especially due to Q5). But keeping the models (4,1) and (8,1) is far from absurd and can be a correct answer too so far.
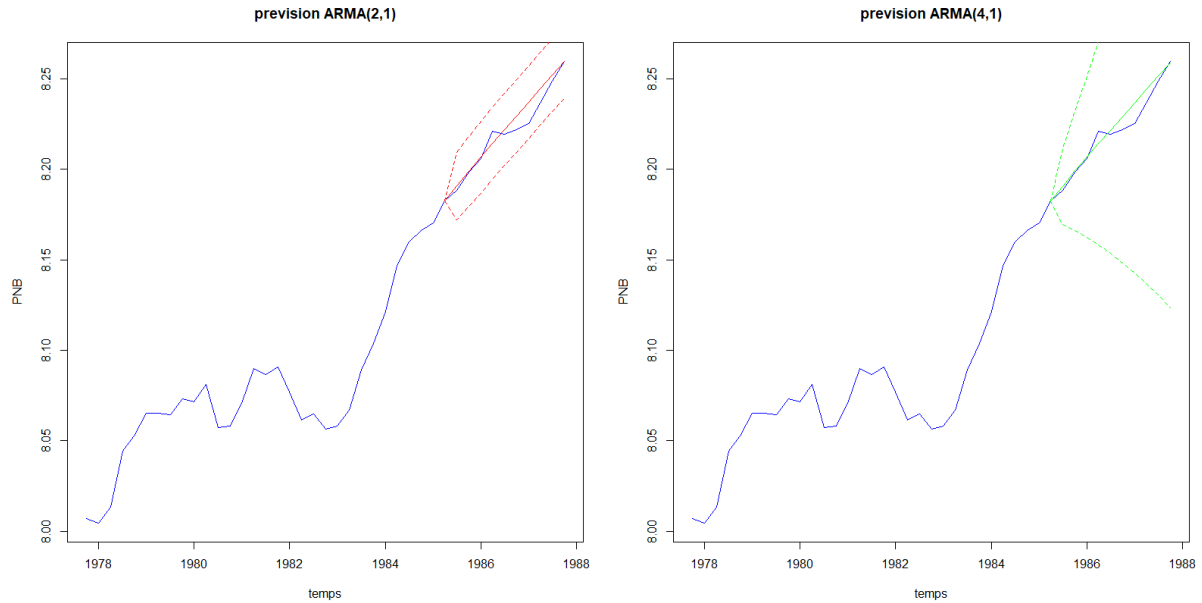
We are now going the plot the predictions that these models can produce in order to assess which one is the closest to the real values.

```
n <- 10
T=length(annee)
index <- 1:(T - n - 1)
res01 <- predict(arima(PNB[index], c(0, 2, 1)), n)
res21 <- predict(arima(PNB[index], c(2, 2, 1)), n)
res41 <- predict(arima(PNB[index], c(4, 2, 1)), n)
res81 <- predict(arima(PNB[index], c(8, 2, 1)), n)


par(mfrow = c(1, 4))
plot(annee[(T - 4 * n):T], PNB[(T - 4 * n):T - 1], main = "prevision ARMA(0,1)", t = "l", col = "blue", xlab = "temps", ylab = "PNB")
lines(annee[(T - n):T], c(PNB[T - n - 1], res01$pred))
lines(annee[(T - n):T], c(PNB[T - n - 1], res01$pred) + c(0,res02$se) * 1.96, lty = 2)
lines(annee[(T - n):T], c(PNB[T - n - 1], res01$pred) - c(0,res02$se) * 1.96, lty = 2)
plot(annee[(T - 4 * n):T], PNB[(T - 4 * n):T - 1], main = "prevision ARMA(2,1)",t = "l", col = "blue", xlab = "temps", ylab = "PNB")
lines(annee[(T - n):T], c(PNB[T - n - 1], res21$pred), col = "red")
lines(annee[(T - n):T], c(PNB[T - n - 1], res21$pred) + c(0,res82$se) * 1.96, lty = 2, col = "red")
lines(annee[(T - n):T], c(PNB[T - n - 1], res21$pred) - c(0,res82$se) * 1.96, lty = 2, col = "red")
plot(annee[(T - 4 * n):T], PNB[(T - 4 * n):T - 1], main = "prevision ARMA(4,1)",t = "l", col = "blue", xlab = "temps", ylab = "PNB")
lines(annee[(T - n):T], c(PNB[T - n - 1], res41$pred), col = "green")
lines(annee[(T - n):T], c(PNB[T - n - 1], res41$pred) + c(0,res01$se) * 1.96, lty = 2, col = "green")
lines(annee[(T - n):T], c(PNB[T - n - 1], res41$pred) - c(0,res01$se) * 1.96, lty = 2, col = "green")
plot(annee[(T - 4 * n):T], PNB[(T - 4 * n):T - 1], main = "prevision ARMA(8,1)",t = "l", col = "blue", xlab = "temps", ylab = "PNB")
lines(annee[(T - n):T], c(PNB[T - n - 1], res81$pred), col = "black")
lines(annee[(T - n):T], c(PNB[T - n - 1], res81$pred) + c(0,res01$se) * 1.96, lty = 2, col = "black")
lines(annee[(T - n):T], c(PNB[T - n - 1], res81$pred) - c(0,res01$se) * 1.96, lty = 2, col = "black")
```



We can see from these results that the models (2,1) and (4,1) seem to be more accurate than the (8,1) model. We compute only the (2,1) and (4,1) models in order to get a clearer view of the predictions.

prevision ARMA(2,1)          prevision ARMA(4,1)

We observe that the confidence interval of the model (2,1) model is way smaller than the confidence interval of the model (4,1). Both predictions are very similar. Therefore the prediction is way more precise with the model (2,1).

We can conclude that the prediction from the model ARIMA(2,2,1) seems to be the best out of these four models.

Considering all of the previous results that we have obtained on the log likelihood, the aic coefficient and the results on the residuals, we conclude that the best model of this time series is the model ARIMA(2,2,1).