

# PROJECTE ANÀLISI DE DADES: Entrada turística a Espanya del 2000 a l'actualitat

*David Anglada Rotger i Andreu Huguet Segarra*

*17/5/2019*

## Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introducció</b>  | <b>1</b>  |
| <b>2</b> | <b>Identificació del model</b>  | <b>2</b>  |
| 2.1      | Representació gràfica de les dades . . . . .                                      | 2         |
| 2.2      | Transformació de les dades . . . . .  | 3         |
| 2.3      | ACF/PACF de les dades i proposta de models . . . . .                              | 9         |
| 2.4      | Models proposats . . . . .  | 11        |
| <b>3</b> | <b>Estimació dels models</b>  | <b>11</b> |
| <b>4</b> | <b>Validació dels Models</b>  | <b>12</b> |
| 4.1      | Estudi dels residus dels models . . . . .   | 13        |
| 4.2      | Estabilitat dels Models . . . . .   | 23        |
| 4.3      | Capacitat de predicció . . . . .  | 24        |
| 4.4      | Elecció de model . . . . .  | 25        |
| <b>5</b> | <b>Predicció a llarg termini</b>  | <b>25</b> |
| <b>6</b> | <b>Tractament de outliers</b>   | <b>26</b> |
| 6.1      | Identificació i estimació del model per la sèrie linealitzada . . . . .           | 28        |
| 6.2      | Validació del model per la sèrie linealitzada . . . . .                           | 29        |
| 6.3      | Estabilitat del model proposat per la sèrie linealitzada . . . . .                | 32        |
| 6.4      | Capacitat de predicció del model proposat per la sèrie linealitzada . . . . .     | 33        |
| 6.5      | Previsions a llarg termini pel model proposat per la sèrie linealitzada . . . . . | 34        |
| <b>7</b> | <b>Comparació dels dos models</b>   | <b>34</b> |
| <b>8</b> | <b>Comentaris finals</b>  | <b>36</b> |

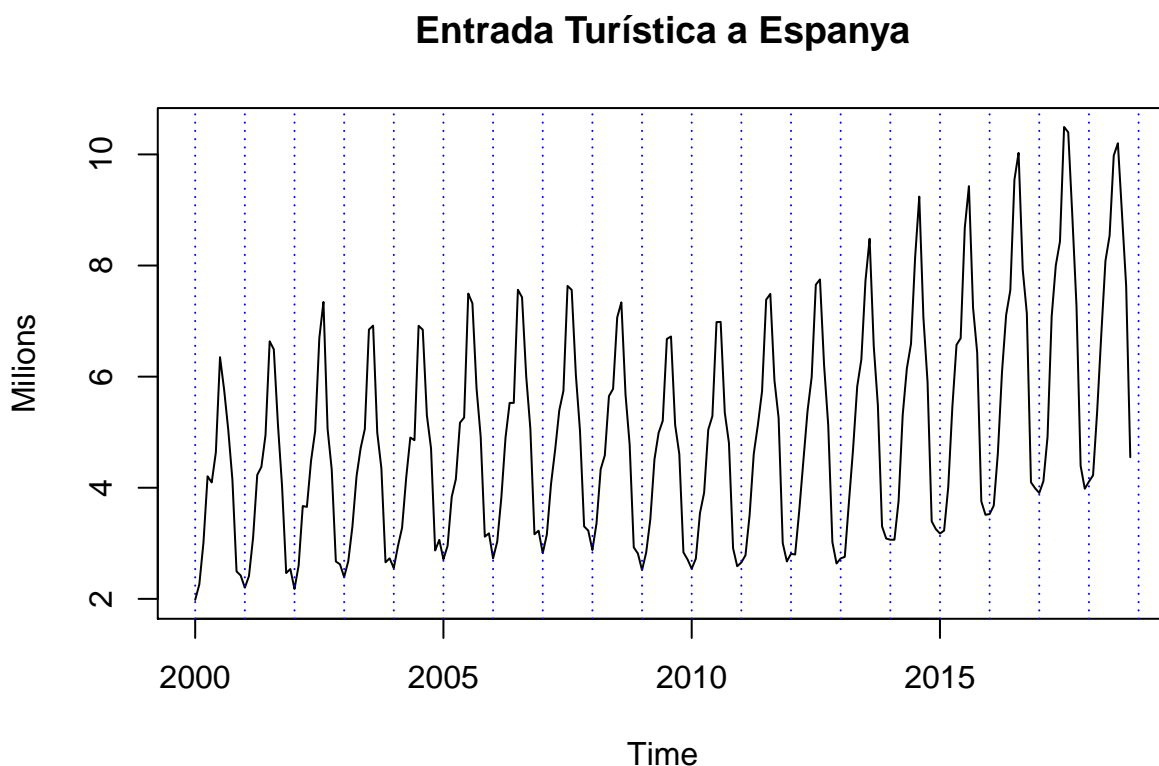
## 1 Introducció

En aquest estudi, es realitzarà l'anàlisi de l'entrada turística a Espanya des de l'any 2000 fins a l'actualitat. S'analitzaran les tendències tant generals com estacionals, així com la presència de valors atípics que puguin dificultar-ne l'examinació. L'objectiu de tot plegat és arribar a proposar un model vàlid per poder predir amb la màxima precisió l'entrada turística del proper any.

**FONT DE DADES:** Ministeri d'Indústria, Comerç i Turisme. <https://sedeaplicaciones.minetur.gob.es/Badase/BadasiUI/lstSeriesInformesPostBack.aspx> U64A2: FRONTUR ENTRADA DE TURISTAS

## 2 Identificació del model

### 2.1 Representació gràfica de les dades



Un cop feta la representació de les dades, s'observa una clara tendència variant, generalment creixent al llarg dels anys. És menys pronunciada entre els anys 2000 i 2010, fins i tot amb una petita baixada entre els anys 2007 i 2010, i sembla que es pronuncia a partir de l'any 2011.

Pel que fa a la variància, s'observa que va augmentant a mesura que es pronuncia la tendència creixent. És a dir, en els anys 2000-2010, la variància és menor que en els anys 2011-2019, on el creixement augmenta.

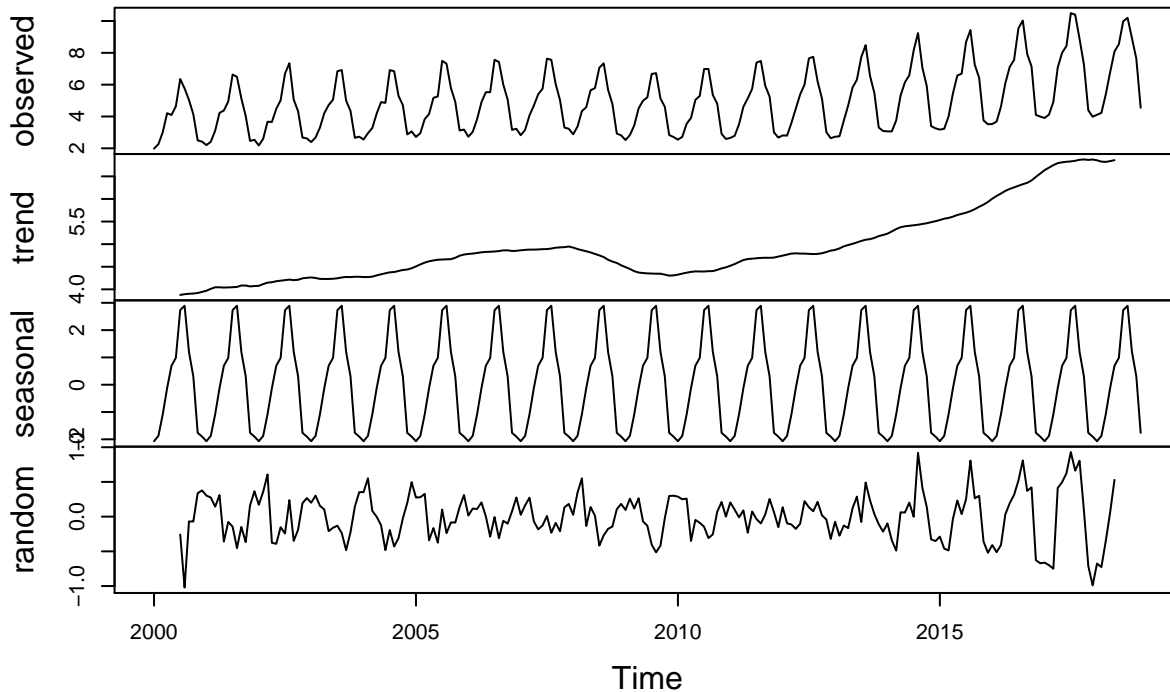
#### 2.1.1 Descomposició en components bàsiques

Per poder analitzar millor les dades, es realitza la seva descomposició en les seves components bàsiques, és a dir, el model aditiu de la sèrie:

$$X_t = T_t + S_t + C_t + \omega_t$$

on: \*  $T_t$  és la **tendència** de la sèrie a llarg termini. \*  $S_t$  és el **seasonal** de la sèrie (patró repetit periòdicament amb període constant). \*  $C_t$  és el **cicle** de la sèrie (patró repetit periòdicament amb període no constant). Aquesta part no surt representada en la descomposició. \*  $\omega_t$  és el soroll aleatori.

## Decomposition of additive time series



S'observa, tal i com s'havia comentat anteriorment, la clara tendència creixent de la sèrie, amb un creixement menys pronunciat a l'inici, una petita baixada entre els anys 2007 i 2010 i una pujada més pronunciada més cap a l'actualitat.

Pel que fa al patró estacional, observem que durant els mesos d'estiu, el número de turistes a Espanya augmenta molt considerablement. Aquest fet que no crida l'atenció, ja que és durant els mesos d'estiu quan més vacances s'agafa la gent i més aprofiten per venir a les costes espanyoles. Durant els mesos de tardor-hivern, observem que el número de turistes cau en picat.

## 2.2 Transformació de les dades

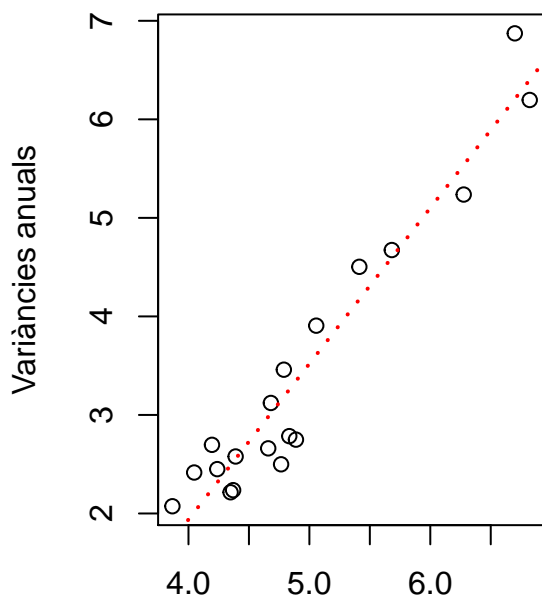
La clara tendència creixent de la mitjana i la variància a mesura que  $t$  augmenta ens fa rebutjar la hipòtesi de l'estacionarietat en les dades.

A continuació s'analitzarà la necessitat de realitzar una sèrie de transformacions amb l'objectiu d'aconseguir estacionarietat en la nostra sèrie temporal i poder aplicar la teoria dels models ARIMA per a sèries estacionàries.

### 2.2.1 Variància constant

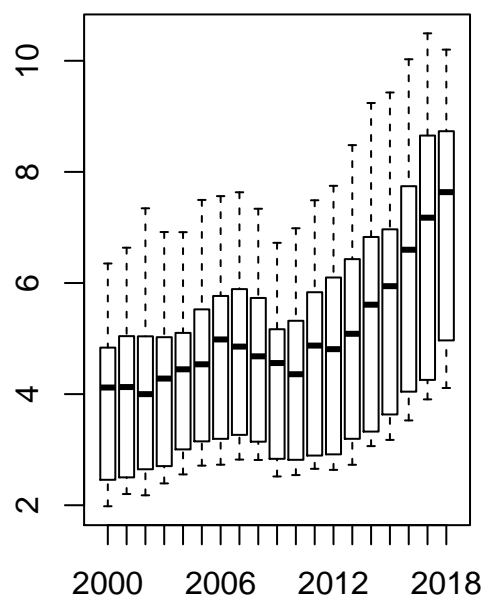
En primer lloc, s'estudiarà si es pot considerar que la variància de les dades sigui constant en el temps. Es comprova amb un plot de la variància front la mitjana i un *boxplot* de les dades cada 12 mesos (que és la freqüència de les nostres dades) que la variància, efectivament, no és constant.

### Mean–Variance Plot



Mitjanes anuals

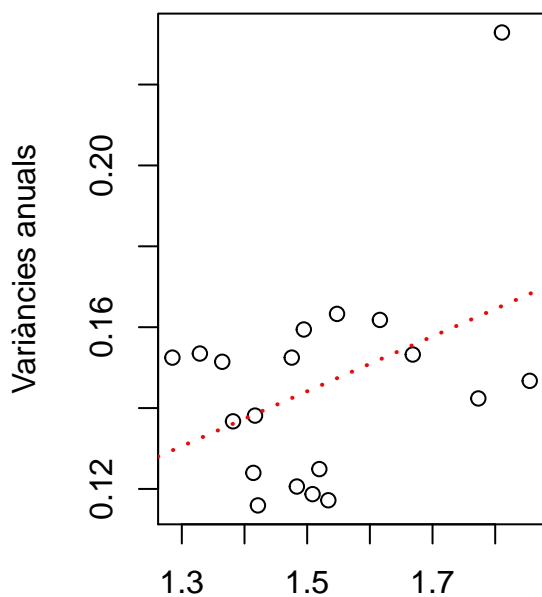
### Monthly Boxplot



Tal i com s'havia observat a simple vista, la variància augmenta a mesura que augmenta la mitja. Per tant, no podem assumir variància constant. Amb el *boxplot* es confirma aquesta hipòtesis.

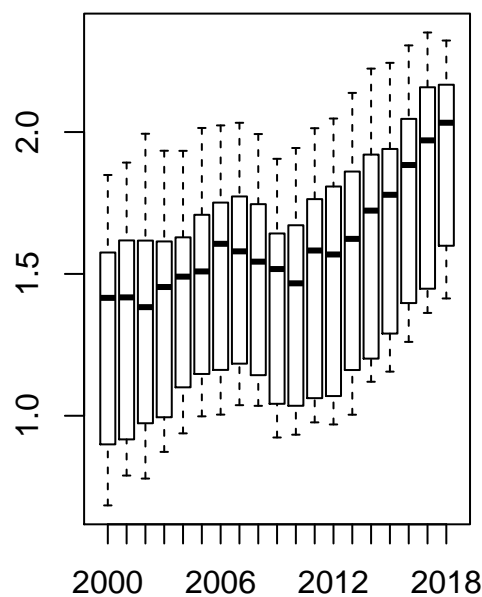
Així doncs, es procedeix a realitzar una transformació logarítmica (*Box-Cox* amb  $\lambda = 0$ ) de la sèrie per homogeneïtzar la variància. Els resultats obtinguts són els següents:

### (Log) Mean–Variance Plot



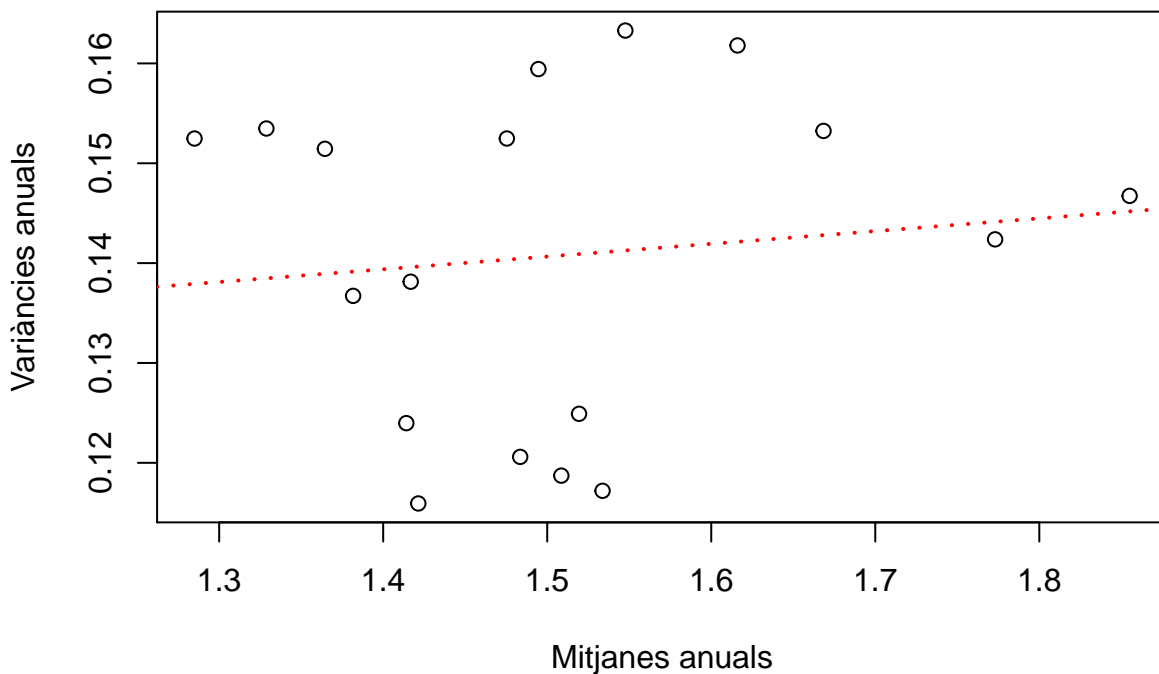
Mitjanes anuals

### (Log) Monthly Boxplot

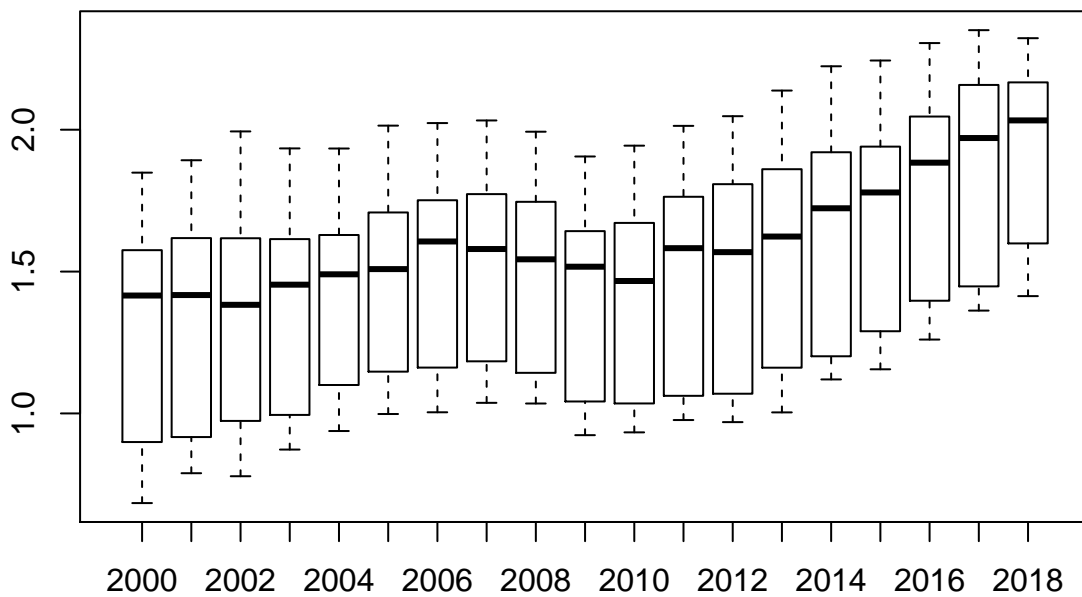


Veiem que encara hi ha una tendència creixent causada per un valor anormalment gran, si l'eliminem i tornem a fer el plot veiem que, en efecte, s'homogenitza la variància tal i com volíem.

### (Log) Mean-Variance Plot



### (Log) Monthly Boxplot



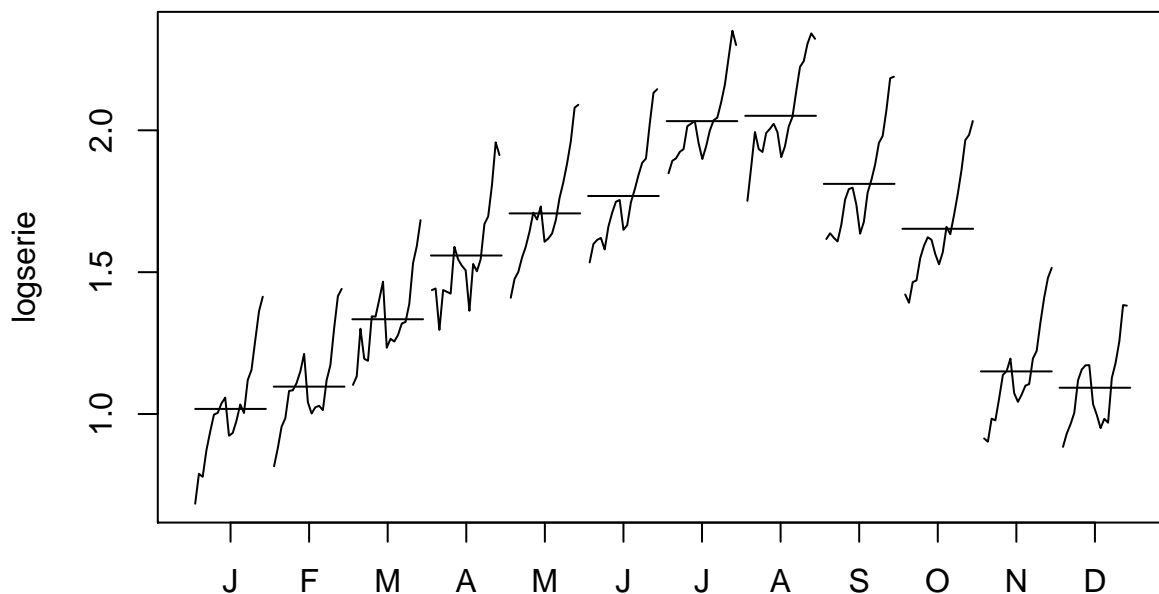
Ara la variància ja es pot considerar constant.

### 2.2.2 Patró estacional

En segon lloc, s'estudiarà l'existència d'un patró estacional en les nostres dades. En cas que hi sigui present, es realitzarà una diferenciació d'ordre 12 per eliminar aquest patró, és a dir,

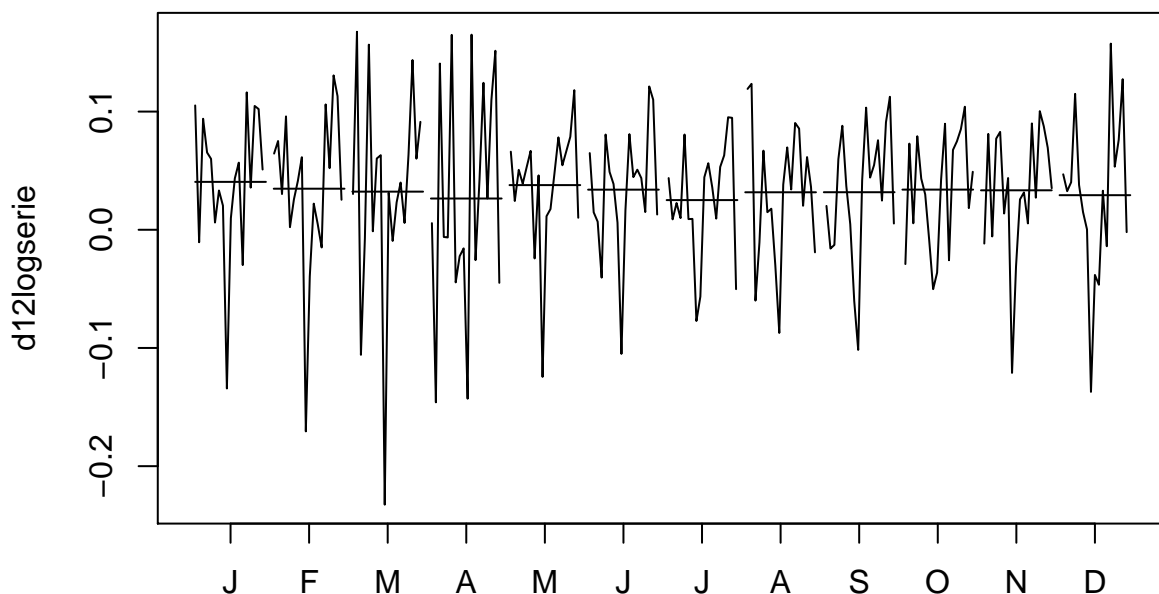
$$W_t := X_t - X_{t-12} = (1 - B^{12})X_t$$

on  $B$  és el *backshift operator*. Es realitza un *monthplot* per comprovar-ne l'existència.



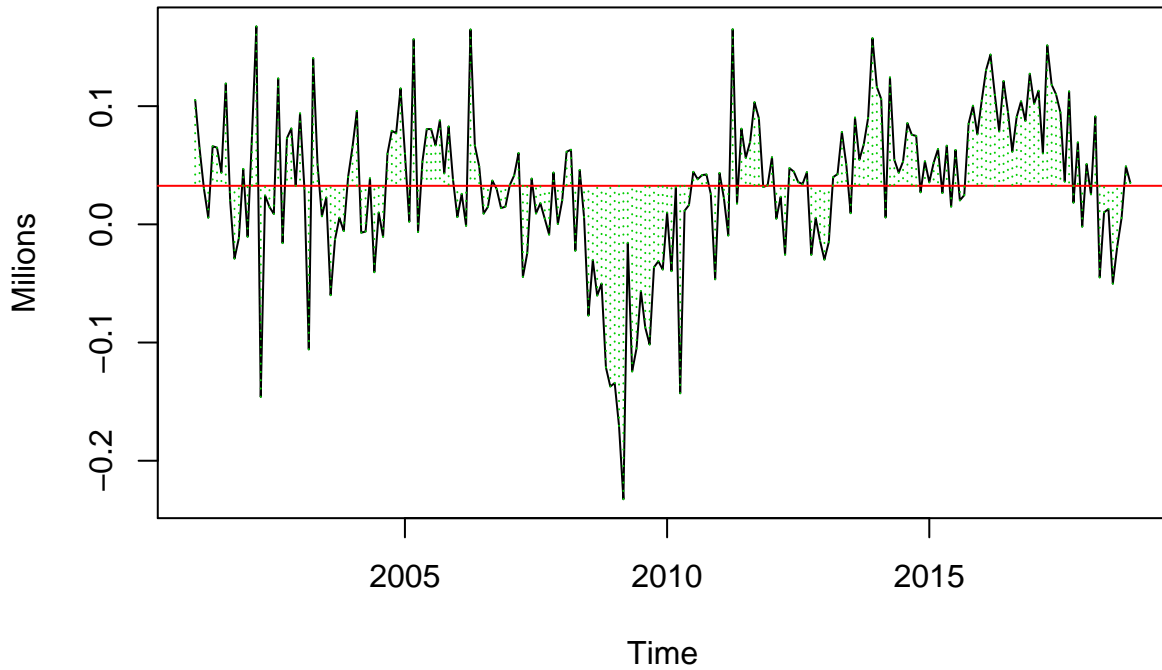
Tal i com s'havia comentat, s'observa una clara pujada de la presència de turistes durant els mesos d'estiu i una baixada en picat en l'entrada de l'hivern/tardor. Així doncs, és necessària una diferenciació d'ordre 12 per eliminar aquest patró.

### Monthplot of logserie



S'observa que amb una diferenciació d'ordre 12 s'ha eliminat el patró estacional.

## Distance from value to mean of d12logserie



Ara bé, la mitjana de la sèrie encara no és constant.

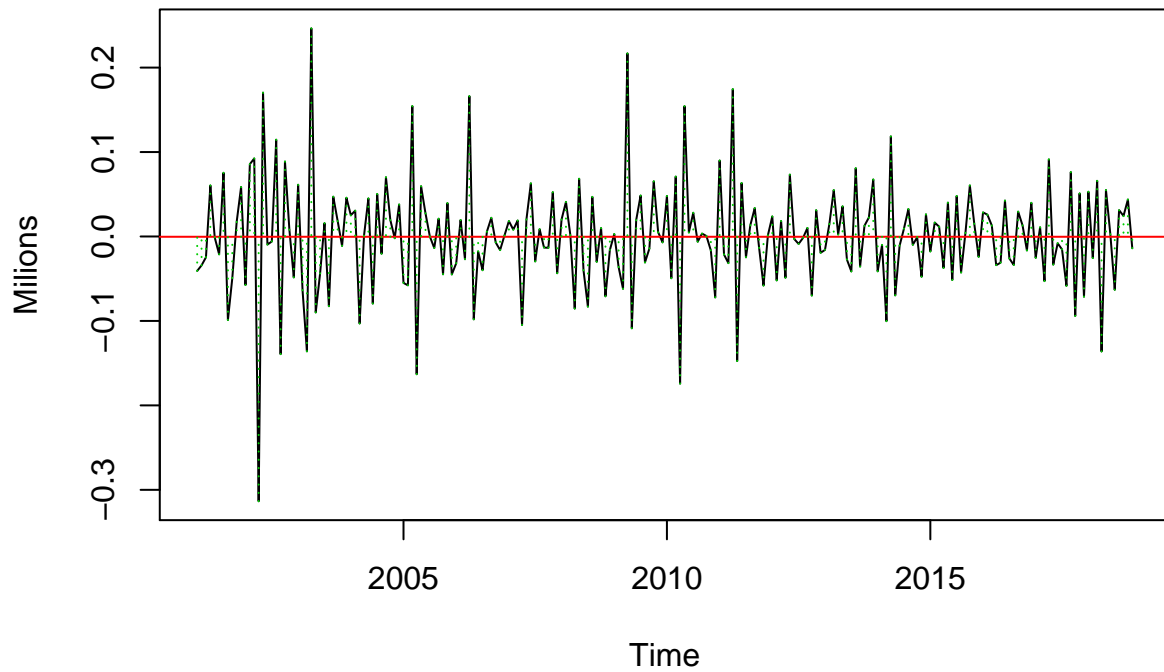
### 2.2.3 Mitjana constant

Per últim, es vol aconseguir que la sèrie tingui mitjana constant igual (i si és possible igual a 0) per a poder considerar definitivament la sèrie com un procés estacionari. Per aconseguir-ho, es realitzaran diferenciacions regulars de la sèrie fins que s'obtingui el resultat desitjat

$$W_t = X_t - X_{t-1} = (1 - B)X_t$$

Es realitza la primera diferenciació. Els valors de mitjana i variància aconseguits són els següents:

## Distance from value to mean of d1d12logserie



```
## [1] "Mean of d1d12logserie: -0.0003272784840803"
```

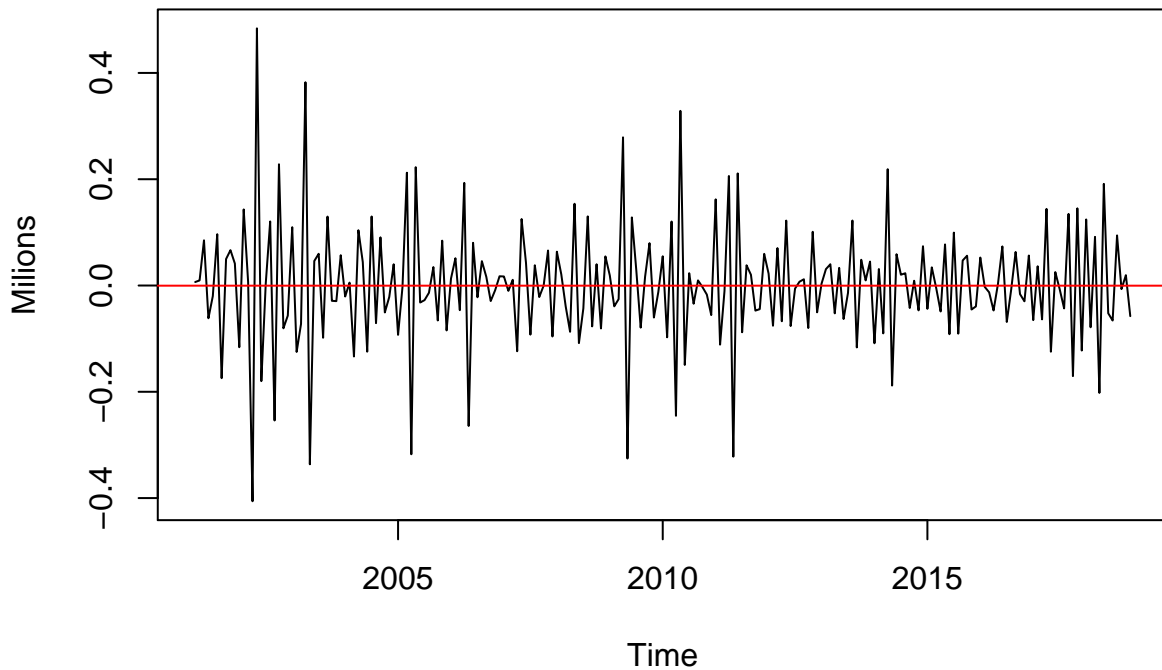
```
## [1] "Variance of d1d12logserie: 0.00426696529209441"
```

Com es pot observar, la mitjana del procés diferenciat regularment un cop es pot arribar a considerar constant i nula.

Es mira de diferenciar un segon cop per veure si la variància augmenta i, per tant, hi ha *overdifferentiation*.



## Plot of d1d12logserie

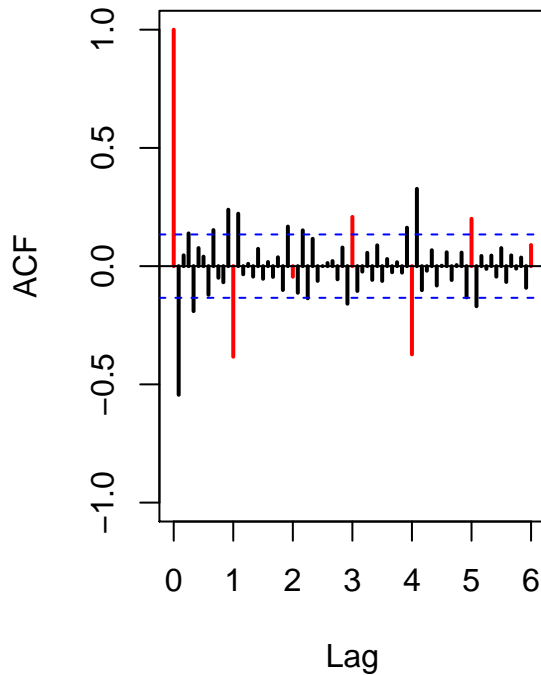
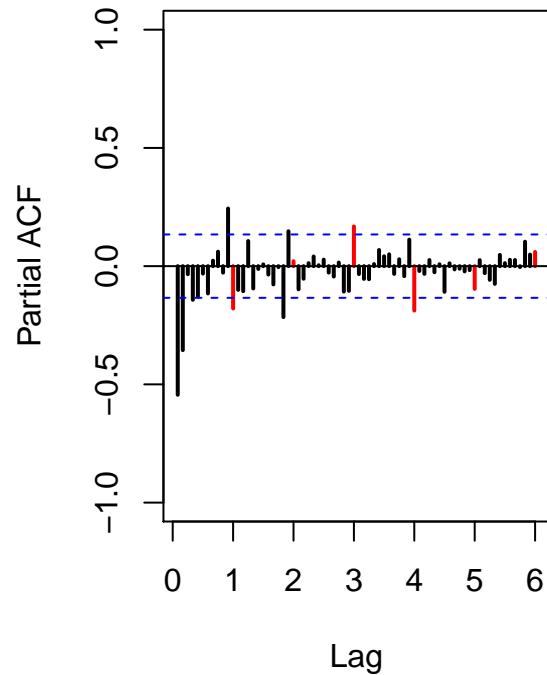


S'observa que la variància augmenta. Així que ens quedarem amb la primera sèrie.

En definitiva, la sèrie transformada pel logaritme, diferenciada un cop i amb una diferenciació d'ordre 12 per eliminar el patró estacional (`d1d12logserie`) és un procés estacionari de mitjana 0.

### 2.3 ACF/PACF de les dades i proposta de models

Tot seguit, es realitza un anàlisi de les funcions d'*AutoCorrelació* i de *Correlació Parcial* de la sèrie transformada, és a dir, de la sèrie estacionària.

**ACF de d1d12logserie****PACF de d1d12logserie**

### 2.3.1 Models proposats per la part regular (p,d,q)

En relació a la part regular de la sèrie, en la funció d'AutoCorrelació (ACF) s'observa un decreixement exponencial alternat en tots els valors. S'observen també valors fora de la banda de confiança en retards llunyans, però poden ser assignats a la aleatorietat del cas 5%. Per tant, en aquest cas, es proposaria  $q = 0$  pel decreixement exponencial.

En tot cas, si es volgués mirar d'incloure el primer valor que sobresurt més que la resta, es podria considerar també  $q = 1$ .

Pel que fa a la funció de Correlació Parcial (PACF) s'observa que els dos primers valors sobresurten més significativament que la resta. La resta de valors es poden considerar nuls, ja que o bé estan dintre de l'interval de confiança, o bé es poden assignar al cas d'aleatorietat del 5%. Per tant, per la part regular, es proposaria  $p = 2$ .

Donat que s'ha realitzat diferenciació un cop, es té que  $d = 1$ . Per tant, els models proposats per la part regular serien  $AR(2)$  o, en tot cas,  $ARMA(2,1)$  sobre la sèrie transformada regular.

### 2.3.2 Models proposats per la part estacional (P,D,Q)

En relació a la part estacional de la sèrie, en la funció d'AutoCorrelació (ACF) s'observa que el primer valor es força significatiu, però també ho són el tercer, el quart i el cinquè, sobretot el quart. Donat que volem intentar proposar un model simplificat, es proposa  $Q = 0$ .

Pel que fa a la funció de Correlació Parcial (PACF) s'observa que sobresurt el primer valor una mica i també sobresurten el tercer i el quart valor lleugerament. Ara bé, no sobresurten de manera tant significativa com en el cas dels valors del ACF i, per tant, podem assignar-ho al cas d'aleatorietat del 5%.

Per tant, en aquest cas, es proposaria  $P = 1$ . En tot cas, es podria proposar  $P = 4$  per mirar d'inclure aquests valors que sobresurten de la banda de confiança.

Donat que s'ha realitzat una diferenciació d'ordre 12 per eliminar el patró estacional, es té que  $D = 1$ . Per tant, el model proposat per la part regular seria un  $AR(1)$

## 2.4 Models proposats

En conclusió, es proposen per la sèrie diferenciada els models estacionals:

- $ARMA(2, 0)(1, 0)_s$
- $ARMA(2, 0)(4, 0)_s$
- $ARMA(2, 1)(1, 0)_s$
- $ARMA(2, 1)(4, 0)_s$

I per la sèrie original, tenint en compte les diferenciacions, es proposen:

- $ARIMA(2, 1, 0)(1, 1, 0)_{12}$
- $ARIMA(2, 1, 0)(4, 1, 0)_{12}$
- $ARIMA(2, 1, 1)(1, 1, 0)_{12}$
- $ARIMA(2, 1, 1)(4, 1, 0)_{12}$

## 3 Estimació dels models

A continuació, s'estimen els coeficients dels dos models proposats i es mira que tots siguin significatius. Per mirar-ho, es realitza el test següent (suposant que estem davant d'un model MA:

$$H_0 : \theta_i = 0$$

$$H_1 : \theta_i \neq 0$$

amb l'estadístic

$$\hat{t} = \frac{\hat{\theta}_i}{\text{se}(\hat{\theta}_i)} \sim t - \text{student}_{T-k}$$

on  $k$  és el nombre total de paràmetres i  $T = 12$  és el període. Ara bé, a la pràctica es diu que un coeficient és significant si  $|\hat{t}| > 2$ .

En primer lloc, s'estimen els coeficients dels models proposats, amb intercept i sense.

```
## ##### ARIMA(2,1,0)(0,1,0) #####
##
## Significance of the coefficients:
##      ar1      ar2      sar1 intercept
##      TRUE      TRUE      TRUE      FALSE
##
## Significance of the coefficients:
##  ar1  ar2  sar1
## TRUE TRUE TRUE

## ##### ARIMA(2,1,0)(4,1,0) #####
##
## Significance of the coefficients:
##      ar1      ar2      sar1      sar2      sar3      sar4 intercept
##      TRUE      TRUE      TRUE      TRUE      FALSE      TRUE      FALSE
```

```

##
## Significance of the coefficients:
##   ar1   ar2   sar1   sar2   sar3   sar4
##  TRUE  TRUE  TRUE   TRUE  FALSE  TRUE
## ##### ARIMA(2,1,1)(0,1,0) #####
##
## Significance of the coefficients:
##       ar1       ar2       ma1       sar1 intercept
##      FALSE      FALSE      TRUE      TRUE      FALSE
##
## Significance of the coefficients:
##   ar1   ar2   ma1   sar1
## FALSE FALSE  TRUE  TRUE
## ##### ARIMA(2,1,1)(4,1,0) #####
##
## Significance of the coefficients:
##       ar1       ar2       ma1       sar1       sar2       sar3       sar4
##      TRUE      TRUE      FALSE      TRUE      TRUE      FALSE      TRUE
## intercept
##      FALSE
##
## Significance of the coefficients:
##   ar1   ar2   ma1   sar1   sar2   sar3   sar4
##  TRUE  TRUE  FALSE  TRUE   TRUE  FALSE  TRUE

```

S'observa que, en cap dels casos, l'intercept no és significatiu i, per tant, es descarten els models amb aquest paràmetre.

1. En el cas del model  $ARIMA(2, 1, 0)(1, 1, 0)_{12}$ , els tres coeficients són significatius.
2. Al model  $ARIMA(2, 1, 0)(4, 1, 0)_{12}$ , s'observa que el **sar3** no és significatiu, però la resta de coeficients sí que ho són.
3. El model  $ARIMA(2, 1, 1)(1, 1, 0)_{12}$  té els coeficients **ar1**, **ar2** no significatius i els altres dos significatius.
4. El model  $ARIMA(2, 1, 1)(4, 1, 0)_{12}$  té els coeficients **ma1**, **sar3** no significatius i els altres significatius.

En termes de *loglikelihood* i de *AIC*, el model que sembla el millor és el segon model  $ARIMA(2, 1, 0)(4, 1, 0)_{12}$ , que és el millor tant en *AIC* com en *loglikelihood*.

Dels altres dos models, donat que el  $ARIMA(2, 1, 0)(1, 1, 0)_{12}$  té tots els coeficients significatius, es descarta el model  $ARIMA(2, 1, 1)(1, 1, 0)_{12}$ . Tot i tenir una mica millor la *loglikelihood*.

Així doncs, *a-priori*, s'escolliria el primer model  $ARIMA(2, 1, 0)(4, 1, 0)_{12}$  com a millor model. Tot i així, es realitzarà la validació i la predicció dels dos models escollits en aquest pas.

## 4 Validació dels Models

Tot seguit, es realitzarà la validació dels dos models proposat. En el procés de validació es realitzarà un anàlisi dels residus ( $Z_t$ ) dels models, es comprovarà que aquests siguin estacionaris i invertibles, es verificarà la seva estabilitat i s'evaluarà la seva capacitat de previsió.

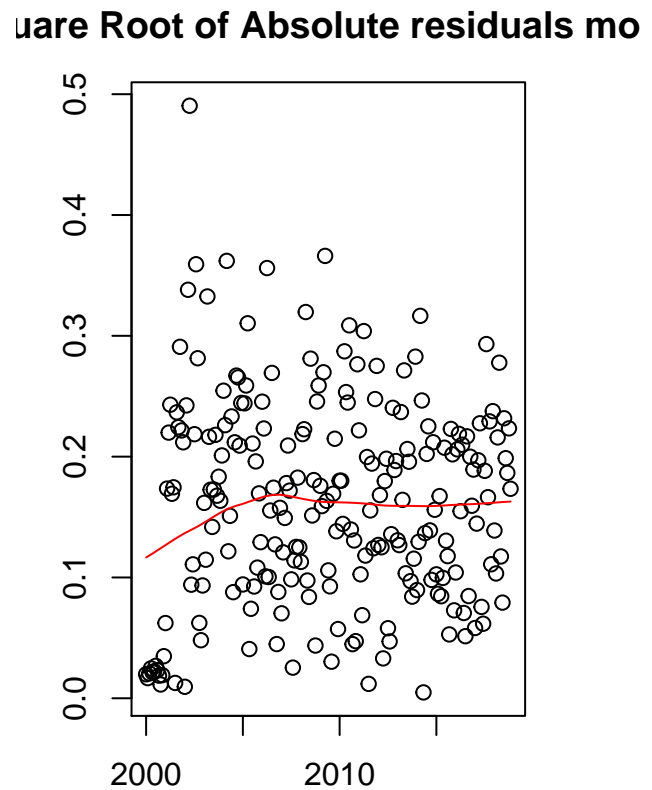
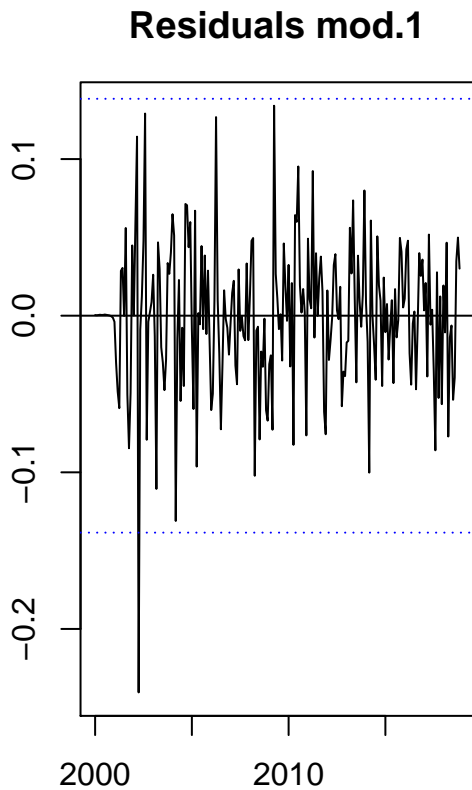
## 4.1 Estudi dels residus dels models

Així doncs, en primer lloc, s'estudiaràn els residus del model i es comprovaran els següents aspectes:

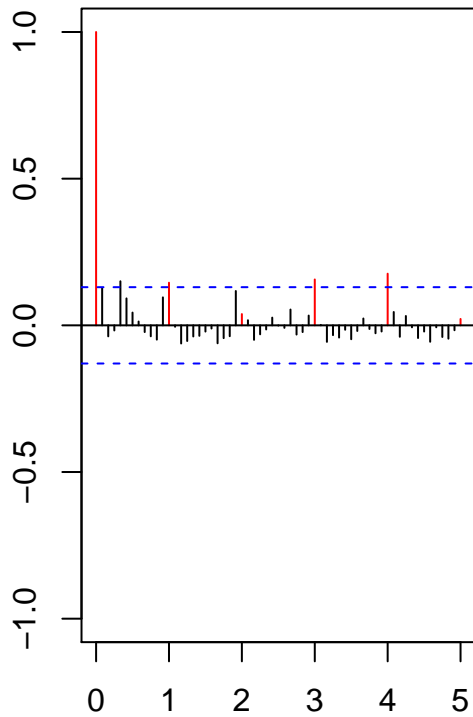
- Homogeneïtat de la variància residual ( $\sigma_Z^2$  constant).
- Normalitat ( $Z_T \sim \mathcal{N}(0, \sigma_Z^2)$ ).
- Independència ( $\rho(k) = 0, \forall k > 0$ ).

### 4.1.1 Homogeneïtat de la variància

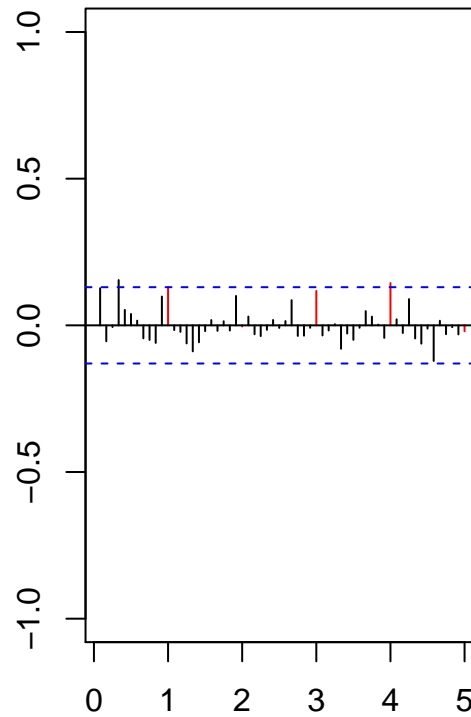
Per comprovar l'homogeneïtat de la variància dels residus, s'analitzen el plot dels mateixos residus, el plot de l'arrel quadrada del seu valor absolut i les funcions ACF i PACF del seu quadrat.



**Series resid^2 mod.1 ACF**

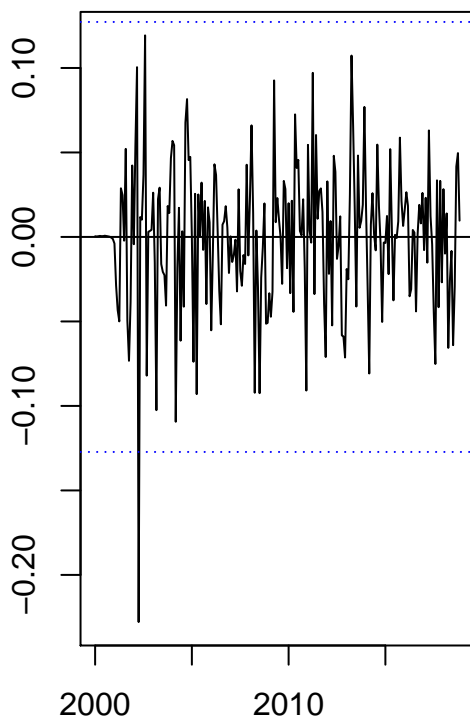


**Series resid^2 mod.1 PACF**

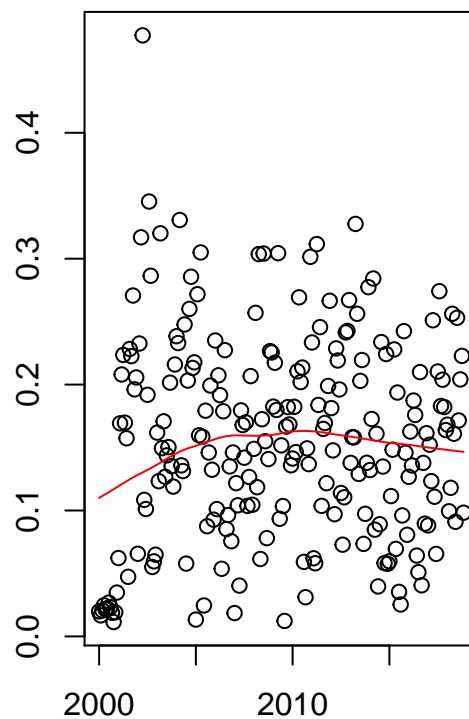


En el cas del primer model (mod.1) no s'observa cap tipus de patró (ni creixent ni decreixent) en el plot dels residus o en el plot de l'arrel quadrada del seu valor absolut. A més, en l'ACF i el PACF del quadrat dels residus tots els valors estan dintre de la banda de confiança i, per tant, els podem considerar nuls.

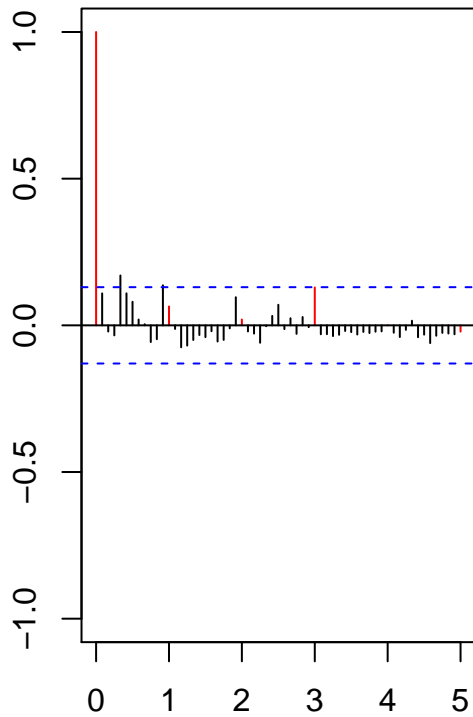
**Residuals mod.2**



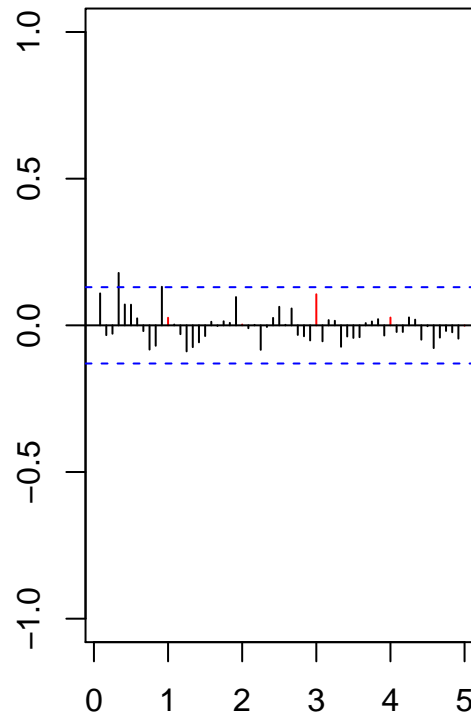
**quare Root of Absolute residuals mo**



**Series resid^2 mod.2 ACF**



**Series resid^2 mod.2 PACF**

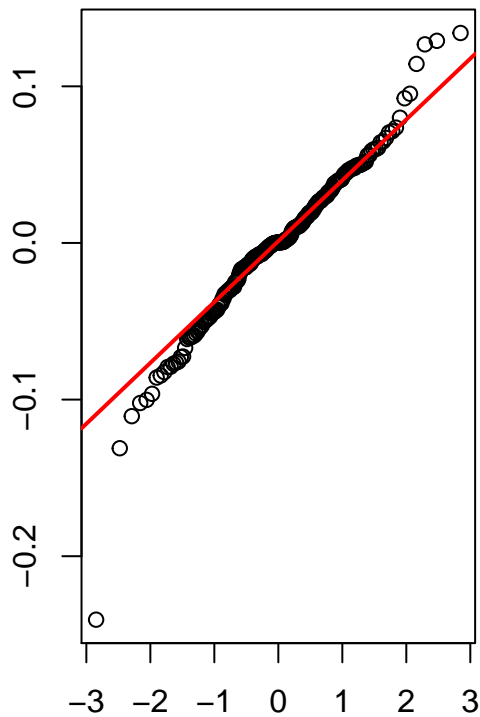


En el cas del segon model (`mod.2`), es poden extreure les mateixes conclusions que en el primer model i, per tant, també es pot assumir homogeneïtat de variància residual.

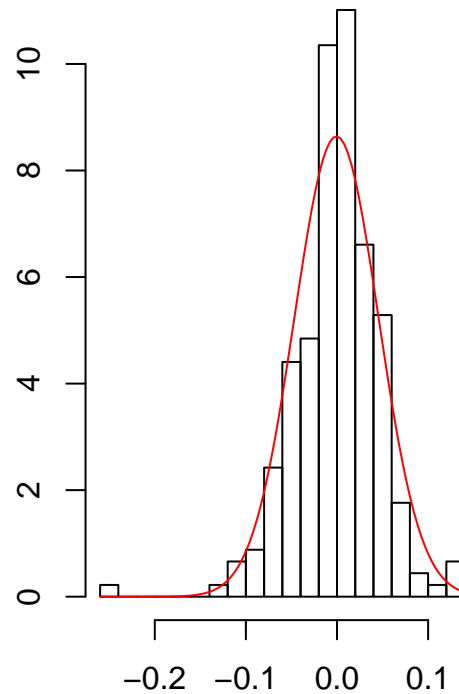
#### 4.1.2 Normalitat

Per comprovar la normalitat dels residus dels models proposats s'estudiarà el Q-Q plot, l'histograma dels residus amb la normal que s'hauria de seguir sobreposada i es realitzarà el test de Sharipo-Wilks.

### Normal Q-Q Plot



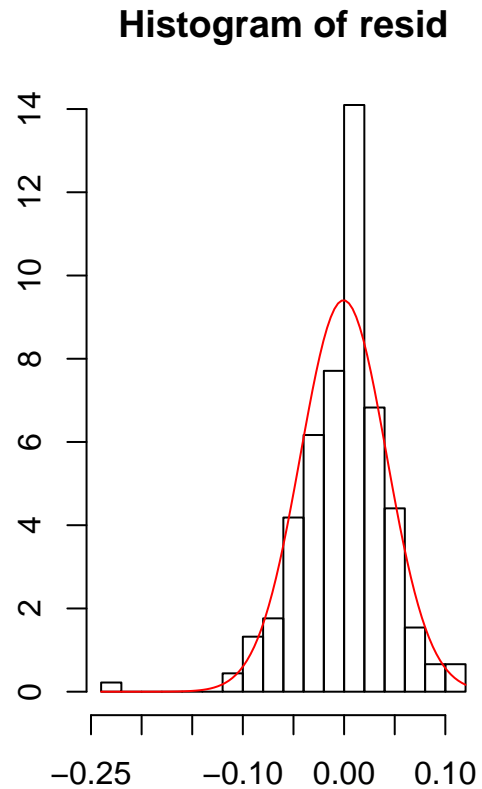
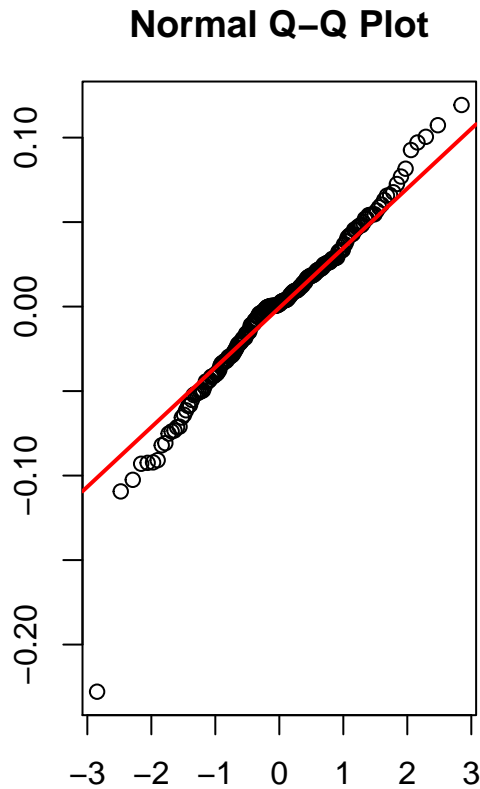
### Histogram of resid



```
##  
## Shapiro-Wilk normality test  
##  
## data: resid(model)  
## W = 0.96408, p-value = 1.713e-05
```

En el cas del model `mod.1`, s'observa en el Q-Q plot que els quantils es situen sobre la línia dels quantils teòrics i que l'histograma s'ajusta a la distribució normal a la que s'hauria d'ajustar (tot i tenir les dues barres més grans una mica per fora de la corba normal). A més, el *p-value* del test de Sharipo-Wilks és  $1.713 \times 10^{-05}$ , menor que 0.05 i, per tant, es pot assumir la hipòtesi de normalitat en els residus.



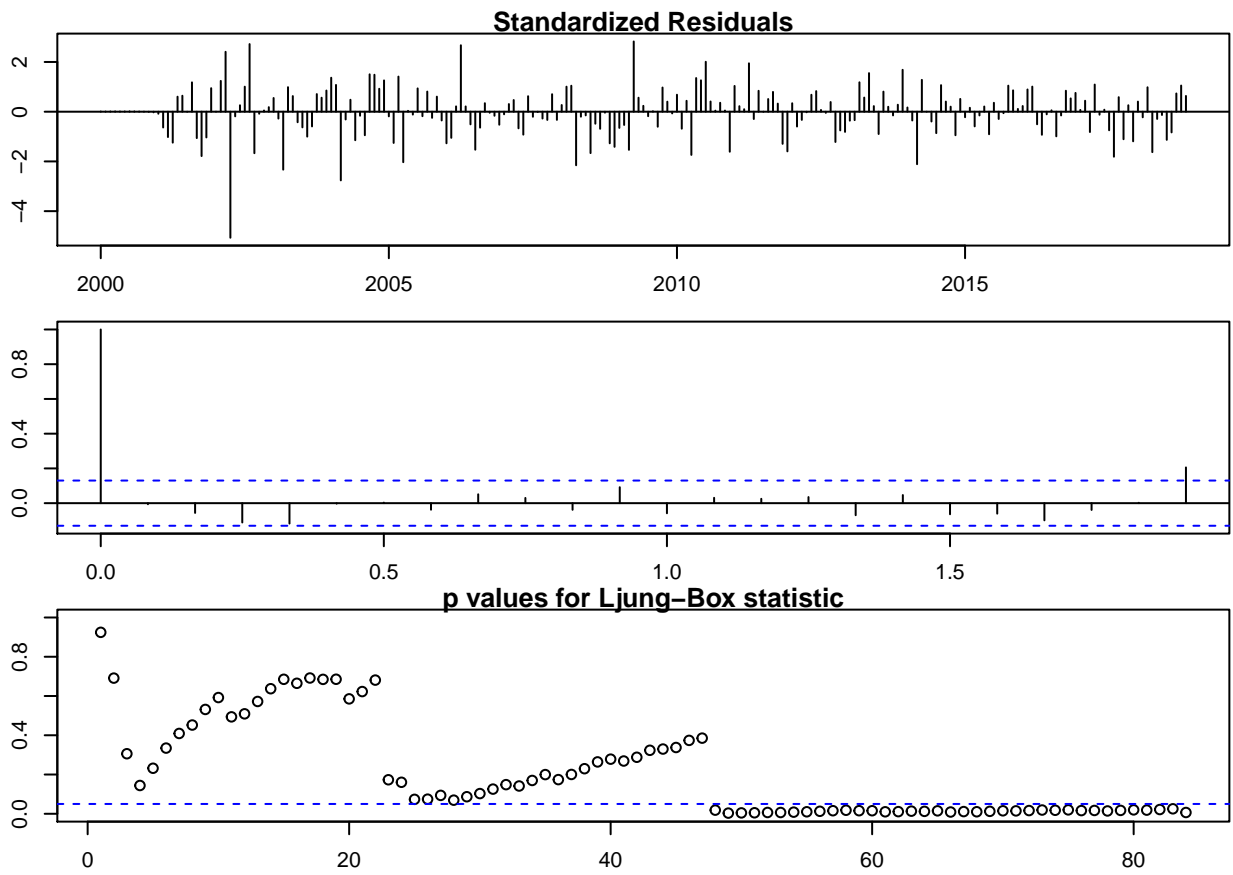
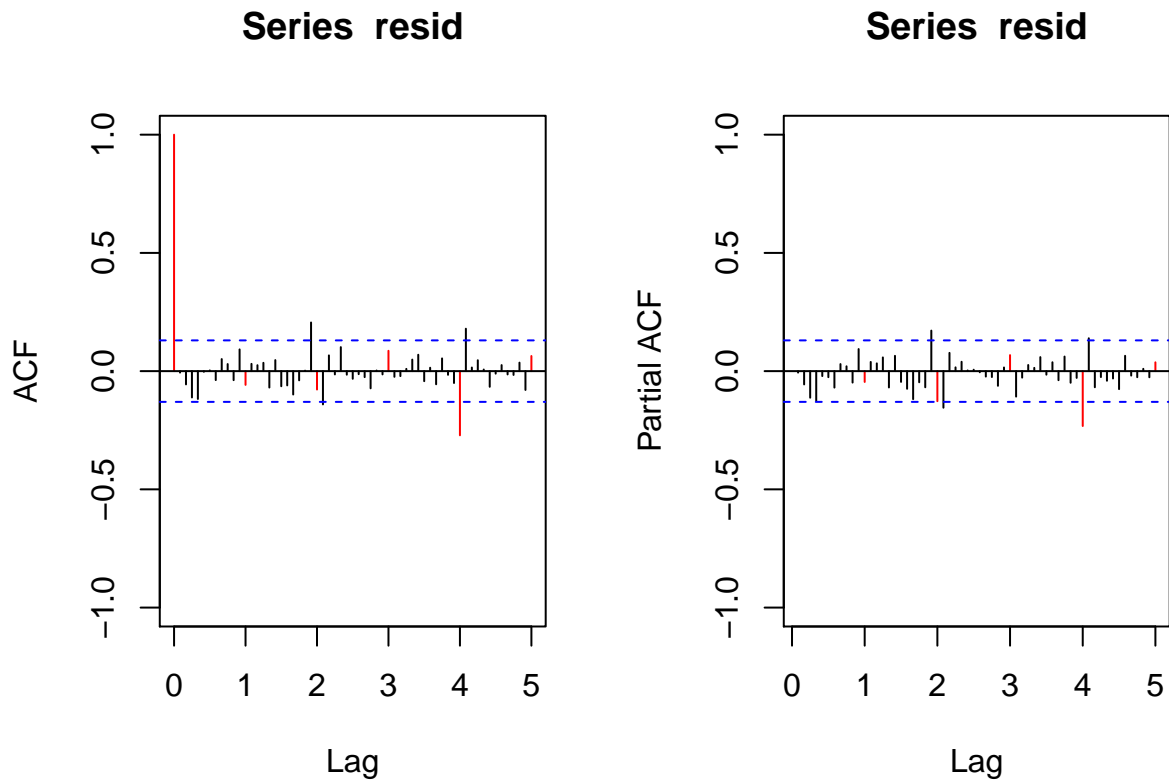


```
##
## Shapiro-Wilk normality test
##
## data:  resid(model)
## W = 0.96197, p-value = 9.536e-06
```

En el cas del model `mod.2`, les conclusions que s'extreuen són les mateixes. En aquest cas, el *p-value* és de  $9.53 \times 10^{-06}$ . Per tant, també assumim normalitat en aquest cas. En l'histograma, en aquest cas, només hi ha una barra que sobresurt de la corba normal.

#### 4.1.3 Independència

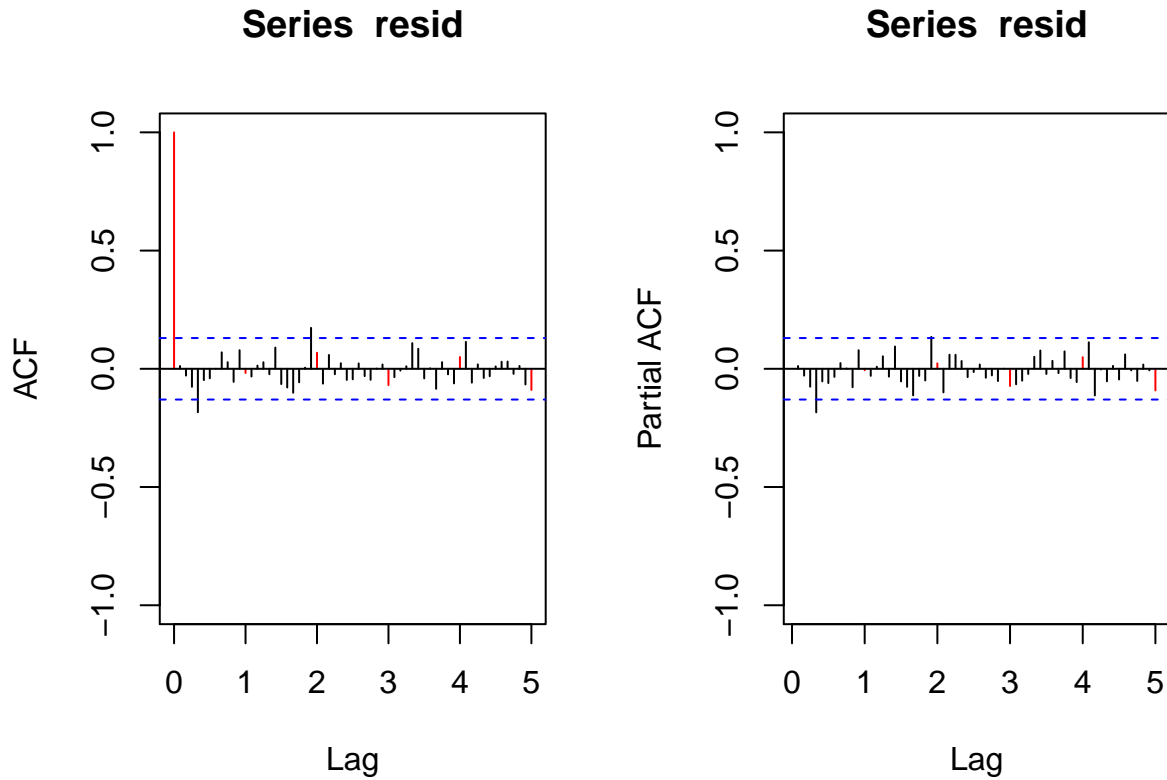
Per comprovar la independència en els residus buscarem que els coeficients de l'autocovariància  $\rho(k) = 0$ ,  $\forall k > 0$ . S'estudiarà el ACF i el PACF dels residus i es realitzarà el test de Ljung-Box.

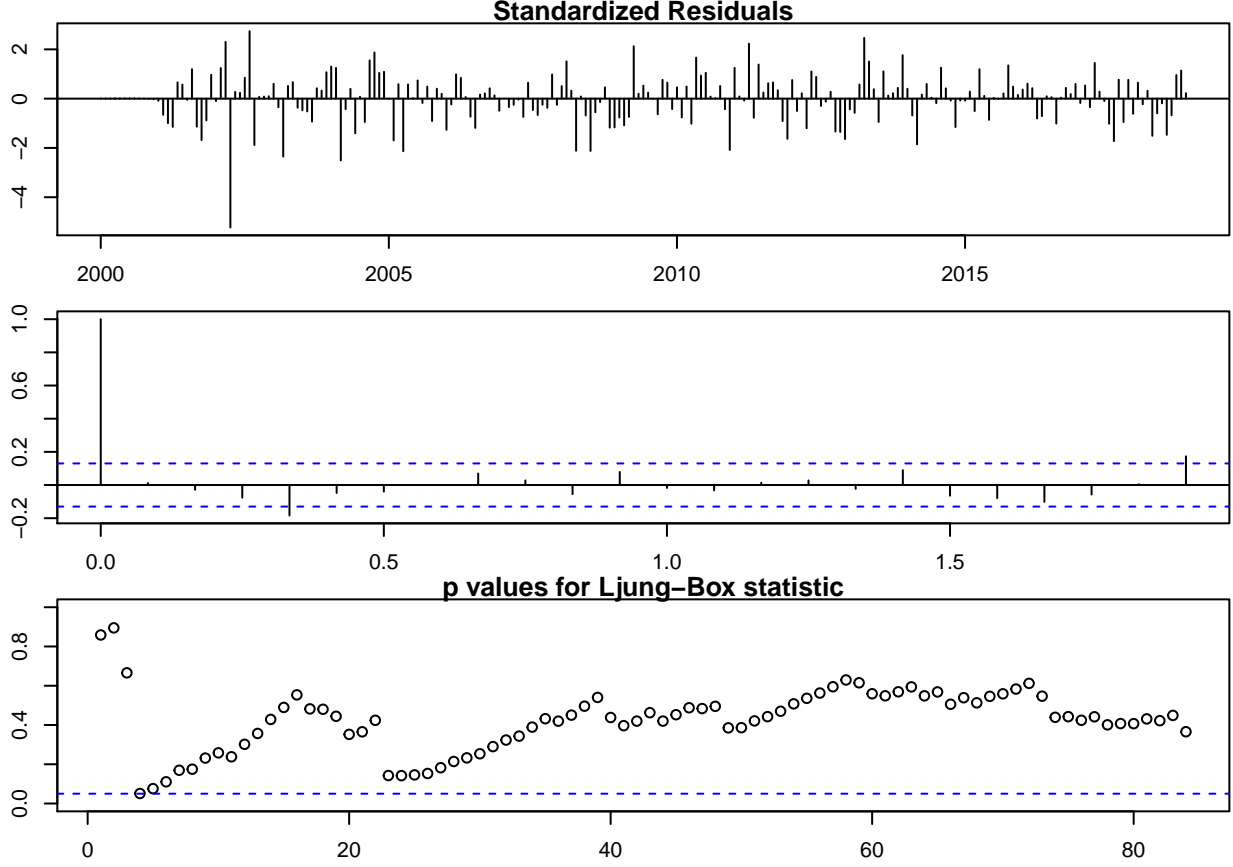


Pel que fa al primer model, en primer lloc observem que les funcions ACF i PACF prenen valors pràcticament

iguals (menys alguns retards llunyans que en un són positius i en l'altre negatius, però estan dins la banda de confiança en els dos casos), cosa que ja fa intuïr que es complirà la independència.

Els residus estandaritzats prenen valors dintre de la franja de  $(-2,2)$ , la gran majoria, que és el comportament esperat. Ara bé, els *p-values* del test Ljung-Box tenen valors superiors a 0.05 en els primers retards però en els retards llunyans no es pot assumir la independència. Tot i així, donat que en els primers retards sí que es té independència, s'assumeix aquesta hipòtesi pel model.





En el segon model, en canvi, observem més diferències entre les gràfiques del ACF i el PACF. Ara bé, en aquest cas els p-values del test de Ljung-Box estan **tots** per sobre de 0.05 i, per tant, podem assumir independència en tots els residus.

#### 4.1.4 Estacionaritat i invertibilitat dels models

Per analitzar l'estacionaritat i la invertibilitat dels models proposats, s'expressaran els models com a models  $AR(\infty)$  i  $MA(\infty)$ :

$$(1 - \phi_1 B - \dots - \phi_p B^p) X_t = (1 + \theta_1 B + \dots + \theta_q B^q) Z_t$$

$$AR(\infty) : \frac{1 - \phi_1 B - \dots - \phi_p B^p}{1 + \theta_1 B + \dots + \theta_q B^q} X_t = (1 - \pi_1 B - \pi_2 B - \dots) X_t = Z_t$$

$$MA(\infty) : \frac{1 + \theta_1 B + \dots + \theta_q B^q}{1 - \phi_1 B - \dots - \phi_p B^p} Z_t = (1 + \psi_1 B + \psi_2 B + \dots) Z_t = X_t$$

A partir d'aquí el models seran *invertibles* si el mòdul de totes les arrels del polinomi característic  $\theta_q(B) = 1 + \theta_1 B + \dots + \theta_q B^q$  és major que 1, és a dir, si

$$\sum_{i \geq 0} \pi_i^2 < \infty$$

Per per altra banda, seran *estacionaris* si el mòdul de totes les arrels del polinomi característic  $\phi_q(B) = 1 - \phi_1 B - \dots - \phi_q B^q$  és major que 1, és a dir, si

$$\sum_{i \geq 0} \psi_i^2 < \infty$$

En el cas del primer model, s'observa que es compleixen totes les condicions i, per tant, el mod.1 és estacionari i invertible.

```
##
## Modul of AR Characteristic polynomial Roots:
## [1] 1.086772 1.086772 1.086772 1.086772 1.086772 1.086772 1.086772 1.086772
## [8] 1.086772 1.086772 1.086772 1.086772 1.086772 1.675221 1.675221
##
## Modul of MA Characteristic polynomial Roots:
## numeric(0)
##
## Psi-weights (MA(inf))
##
## -----
##      psi 1      psi 2      psi 3      psi 4      psi 5
## -0.712679068 0.151578740 0.145923872 -0.158009352 0.060612509
##      psi 6      psi 7      psi 8      psi 9      psi 10
## 0.013106635 -0.030939044 0.017379286 -0.001361260 -0.005222667
##      psi 11     psi 12     psi 13     psi 14     psi 15
## 0.004207147 -0.369553394 0.261873824 -0.054948029 -0.054153900
##      psi 16     psi 17     psi 18     psi 19     psi 20
## 0.058174132 -0.022162680 -0.004934468 0.011413980 -0.006376192
##
## Pi-weights (AR(inf))
##
## -----
##      pi 1      pi 2      pi 3      pi 4      pi 5      pi 6
## -0.7126791 -0.3563327 0.0000000 0.0000000 0.0000000 0.0000000
##      pi 7      pi 8      pi 9      pi 10     pi 11     pi 12
## 0.0000000 0.0000000 0.0000000 0.0000000 0.0000000 -0.3684161
##      pi 13     pi 14     pi 15     pi 16     pi 17     pi 18
## -0.2625624 -0.1312787 0.0000000 0.0000000 0.0000000 0.0000000
##      pi 19     pi 20
## 0.0000000 0.0000000
```

En el cas del segon model, també es compleix tot i, per tant, també és invertible i estacionari.

```
##
## Modul of AR Characteristic polynomial Roots:
## [1] 1.014940 1.014940 1.031881 1.031881 1.031881 1.014940 1.014940
## [8] 1.031881 1.014940 1.014940 1.014940 1.014940 1.031881 1.031881
## [15] 1.014940 1.014940 1.014940 1.014940 1.014940 1.014940 1.031881
## [22] 1.031881 1.014940 1.031881 1.031881 1.014940 1.014940 1.031881
## [29] 1.031881 1.031881 1.031881 1.014940 1.014940 1.031881 1.014940
## [36] 1.031881 1.031881 1.014940 1.031881 1.031881 1.031881 1.031881
## [43] 1.031881 1.014940 1.014940 1.031881 1.031881 1.784666 1.784666
## [50] 1.014940
##
## Modul of MA Characteristic polynomial Roots:
```

```

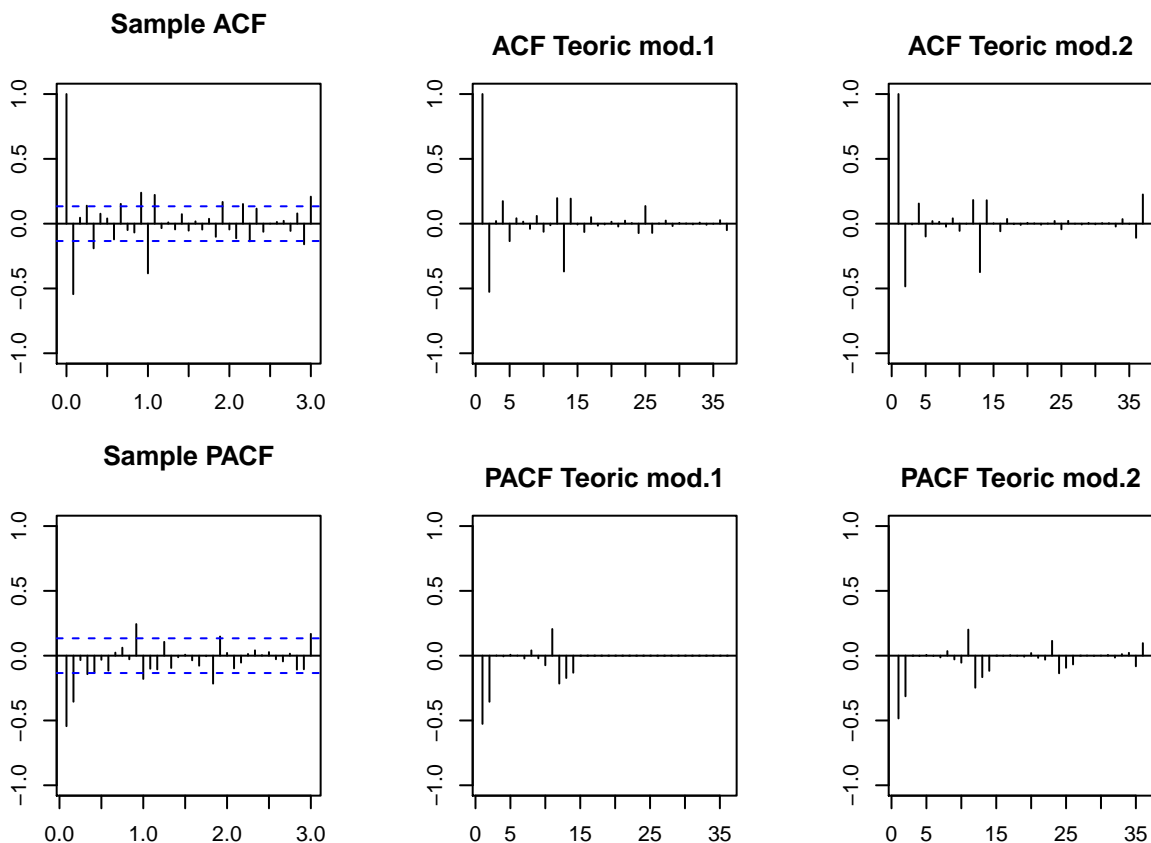
## numeric(0)

##
## Psi-weights (MA(inf))
##
## -----
##      psi 1      psi 2      psi 3      psi 4      psi 5
## -0.635389324  0.089751056  0.142465394 -0.118699998  0.030691060
##      psi 6      psi 7      psi 8      psi 9      psi 10
##  0.017767293 -0.020925175  0.007717262  0.001666381 -0.003481778
##      psi 11     psi 12     psi 13     psi 14     psi 15
##  0.001689093 -0.370760402  0.235046879 -0.032939177 -0.052868124
##      psi 16     psi 17     psi 18     psi 19     psi 20
##  0.043933706 -0.011316081 -0.006603685  0.007748804 -0.002850158
##
## Pi-weights (AR(inf))
##
## -----
##      pi 1      pi 2      pi 3      pi 4      pi 5      pi 6
## -0.6353893 -0.3139685  0.0000000  0.0000000  0.0000000  0.0000000
##      pi 7      pi 8      pi 9      pi 10     pi 11     pi 12
##  0.0000000  0.0000000  0.0000000  0.0000000  0.0000000 -0.3707803
##      pi 13     pi 14     pi 15     pi 16     pi 17     pi 18
## -0.2355899 -0.1164134  0.0000000  0.0000000  0.0000000  0.0000000
##      pi 19     pi 20
##  0.0000000  0.0000000

```

#### 4.1.5 Comparació entre els ACF/PACF mostrals i els ACF/PACF teòrics

Per últim, comparem els valors del ACF i el PACF de les dades amb els valors teòric. S'observa que, en el cas del model mod.2, els valors teòrics s'aproximen gairebé perfectament als valors mostrals. En el cas del segon model també es podria dir el mateix. Per tant, ambdós models s'aproximen als valors de ACF/PACF de les mostres, potser una mica millor el mod.2 (ja que té més coeficients per calcular els valors teòrics).



## 4.2 Estabilitat dels Models

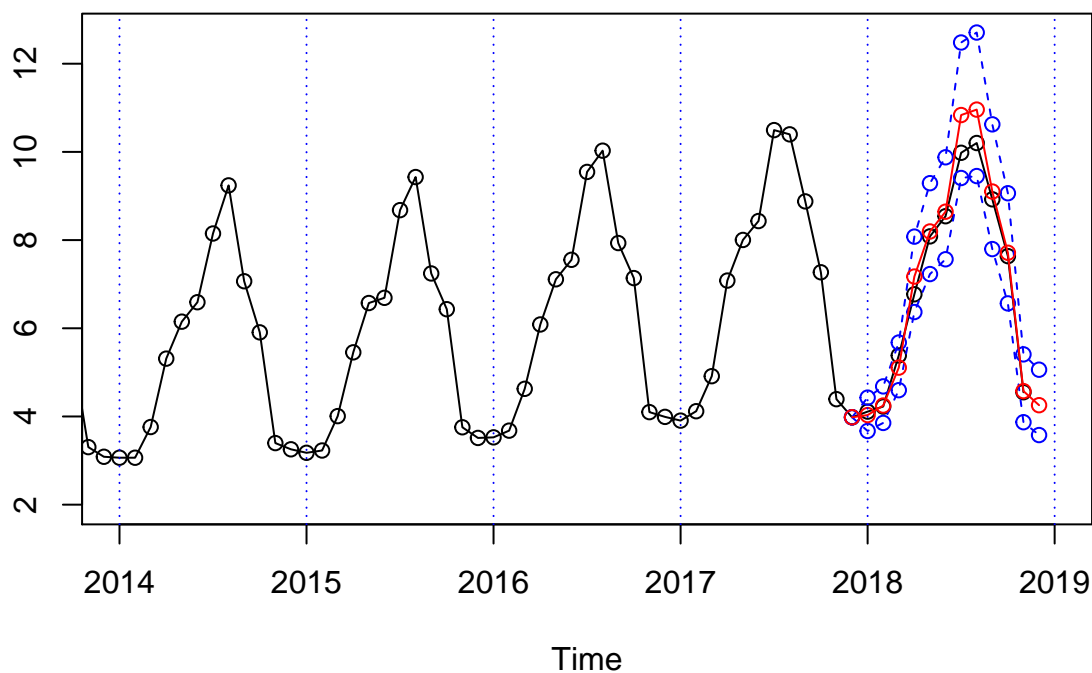
Per comprovar l'estabilitat dels models proposats, calculem els models de la serie ocultant les 12 últimes observacions, és a dir, l'últim període d'observacions. Així doncs, s'observa que el valor dels coeficients varia molt poc, de l'ordre de menys de 0.02 en gairebé tots els casos. Per tant, podem confirmar que els models són estables.

```
## ##### Model ARIMA(2,1,0)(1,1,0)12 amb i sense les 12 últimes observacions #####
##
##          ar1          ar2          sar1
## mod.1.full -0.712679068 -0.356332714 -0.368416056
## mod.1.cut  -0.717347180 -0.361884118 -0.359686924
## abs.diff    0.004668112  0.005551404  0.008729132
## significance 1.000000000  1.000000000  1.000000000
## ##### Model ARIMA(2,1,0)(4,1,0)12 amb i sense les 12 últimes observacions #####
##
##          ar1          ar2          sar1          sar2          sar3
## mod.2.full -0.63538932 -0.313968537 -0.370780339 -0.1911898  0.014182477
## mod.2.cut  -0.64580060 -0.319460181 -0.374558744 -0.2094321  0.002321698
## abs.diff    0.01041128  0.005491644  0.003778405  0.0182423  0.011860779
## significance 1.000000000  1.000000000  1.000000000  1.0000000  0.000000000
##
##          sar4
## mod.2.full -0.329852718
## mod.2.cut  -0.323641262
## abs.diff    0.006211457
## significance 1.000000000
```

### 4.3 Capacitat de predicció

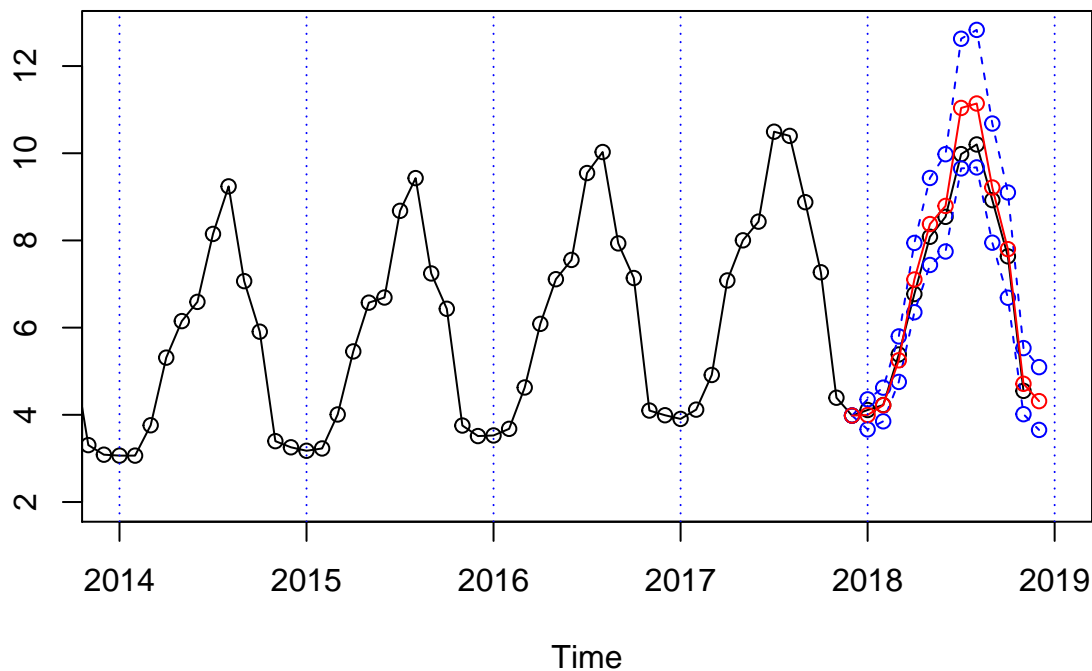
A continuació s'avaluarà la capacitat de predicció dels dos models proposats fent-los predir el valor de les 12 últimes observacions utilitzant la resta d'observacions conegudes.

**Model mod.1 ARIMA(2,1,0)(1,1,0)12**



## ##### Errors de predicció del model mod.1 #####

**Model mod.2 ARIMA(2,1,1)(4,1,0)12**





## ##### Errors de predicció del model mod.2 #####

Com es pot veure, les prediccions (en vermell a la gràfica) de les últimes 12 observacions són semblants i prou bones, ja que en ambdós casos s'apropen força a la realitat. A més, el valor real de les observacions (en negre a la gràfica) queda dins l'interval de confiança (en blau a la gràfica) dels valors predits. Per tant, es pot concloure que els models tenen bona capacitat de predicció. A més, els errors de predicció (l'**Error Quadràtic Mitjà** i l'**Error Absolut Mitjà**) dels dos models són semblants i molt petits ( $> 0.05$  en ambdós casos).

#### 4.4 Elecció de model

En definitiva, els dos presenten un comportament similar en la predicció de les últimes 12 observacions, però que el test d'independència ha sortit molt millor en el segon model, s'escull el segon model, el model  $ARIMA(2, 1, 1)(4, 1, 0)_{12}$ . A més, com ja s'havia comentat abans, aquest model era el millor en  $AIC$  i  $loglikelihood$ .

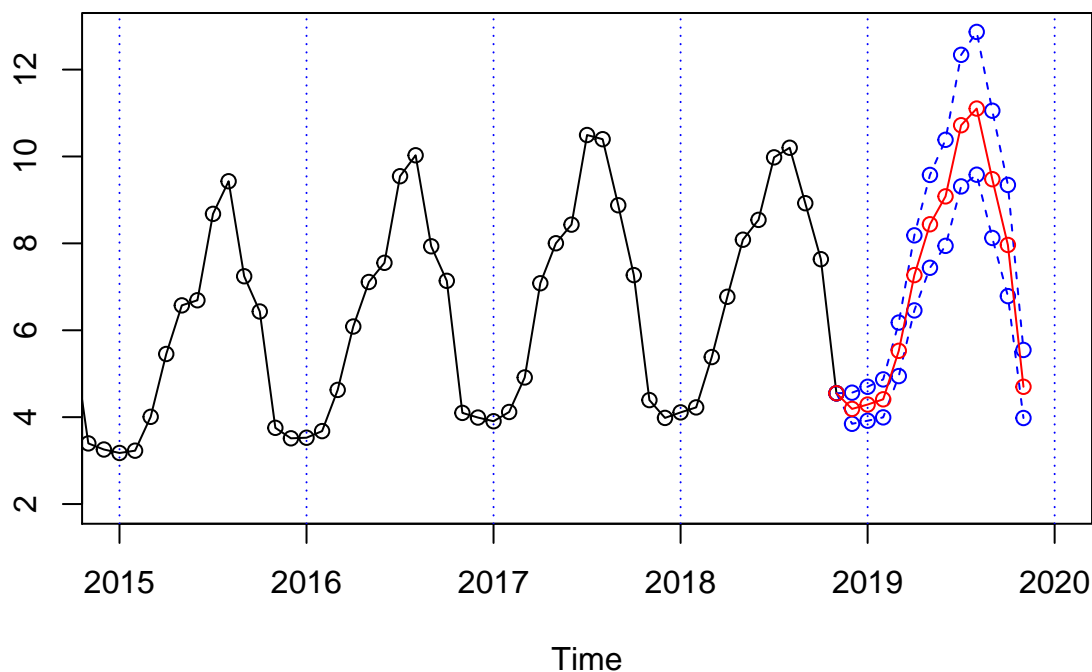
### 5 Predicció a llarg termini

A continuació, s'utilitza el model escollit  $ARIMA(0, 1, 1)(1, 1, 0)_{12}$  per predir el valor de la sèrie els 12 mesos posteriors a l'última dada que es té. Com es pot observar, el valor predit (en vermell a la gràfica) sembla prou raonable per 2 motius:

- Segueix amb la tendència general de la sèrie de creixement lleu.
- Segueix amb el patró estacional vist al llarg de tota la sèrie: pujada molt pronunciada del número de turistes durant els mesos de primavera i estiu i baixada en picat els mesos de tardor i hivern.

A més, els intervals de confiança també segueixen aquestes tendències estacionals.

#### Model $ARIMA(0,1,1)(1,1,0)_{12}$



## t11 pr1 tu1

```
## Nov 2018 4.549899 4.549899 4.549899
## Dec 2018 3.846837 4.189916 4.563591
## Jan 2019 3.917671 4.290606 4.699042
## Feb 2019 3.995362 4.410566 4.868919
## Mar 2019 4.945079 5.526941 6.177267
## Apr 2019 6.457154 7.269534 8.184119
## May 2019 7.438742 8.439524 9.574947
## Jun 2019 7.942888 9.081850 10.384132
## Jul 2019 9.311565 10.719334 12.339936
## Aug 2019 9.580515 11.102785 12.866932
## Sep 2019 8.125041 9.476734 11.053296
## Oct 2019 6.786354 7.963653 9.345191
## Nov 2019 3.980613 4.698813 5.546594
```

## 6 Tractament de *outliers*

Per acabar l'anàlisi d'aquesta serie, es centrarà l'atenció en la detecció i la correcció de possibles *outliers* en la serie. Aquests valors atípics poden ser de tres tipus diferents:

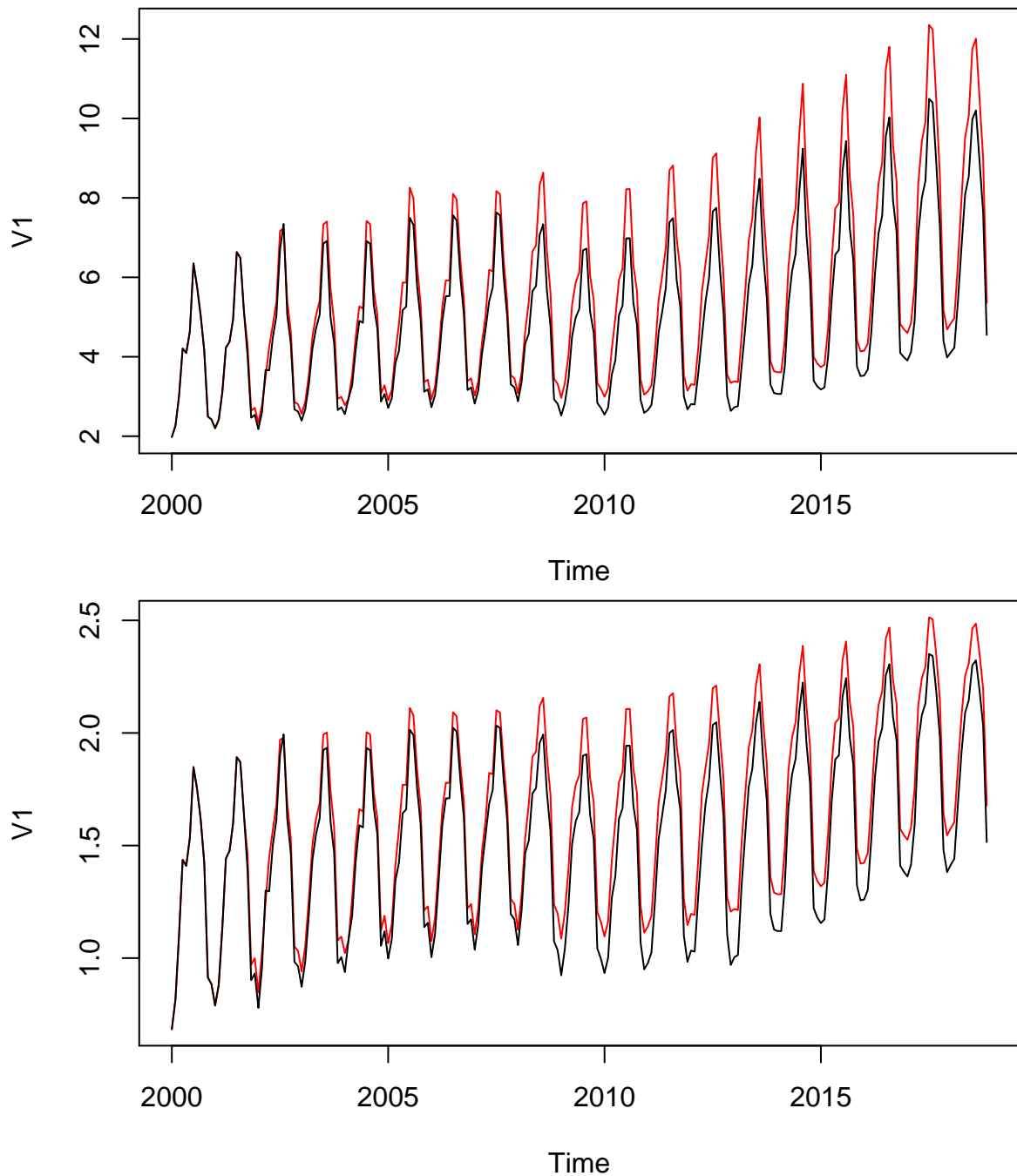
- **Outlier Aitíu (AO):** Afecta només a un període ( $X_t = 1_{t=TO}(t)$ ).
- **Canvi Transitori (TC):** Afecta a un període i el seu efecte decreix exponencialment en els següents períodes ( $X_t = \delta^{(t-TO)} 1_{t \geq TO}(t)$ ).
- **Level Shift (LS):** Afecta a un període i el seu efecte es manté durant els següents períodes ( $X_t = 1_{t \geq TO}(t)$ ).

A continuació es mostren els *outliers* detectats, així com la seva influència sobre les dades i la data en què tenen aquest efecte. S'observa que tenim 10 valors atípics en total, dels quals la majoria són puntuals (6 són AO), 2 són TC i tenim 2 LS. L'*outlier* que més efecte ha tingut és del març del 2002, que ha tingut a més un efecte positiu, és a dir, ha fet créixer la sèrie, tot i només afectar a aquell període, ja que era de tipus AO. En segon lloc, tenim un del tipus LS a l'abril del 2008. Aquest últim es podria associar perfectament a l'entrada en crisi econòmica del país.

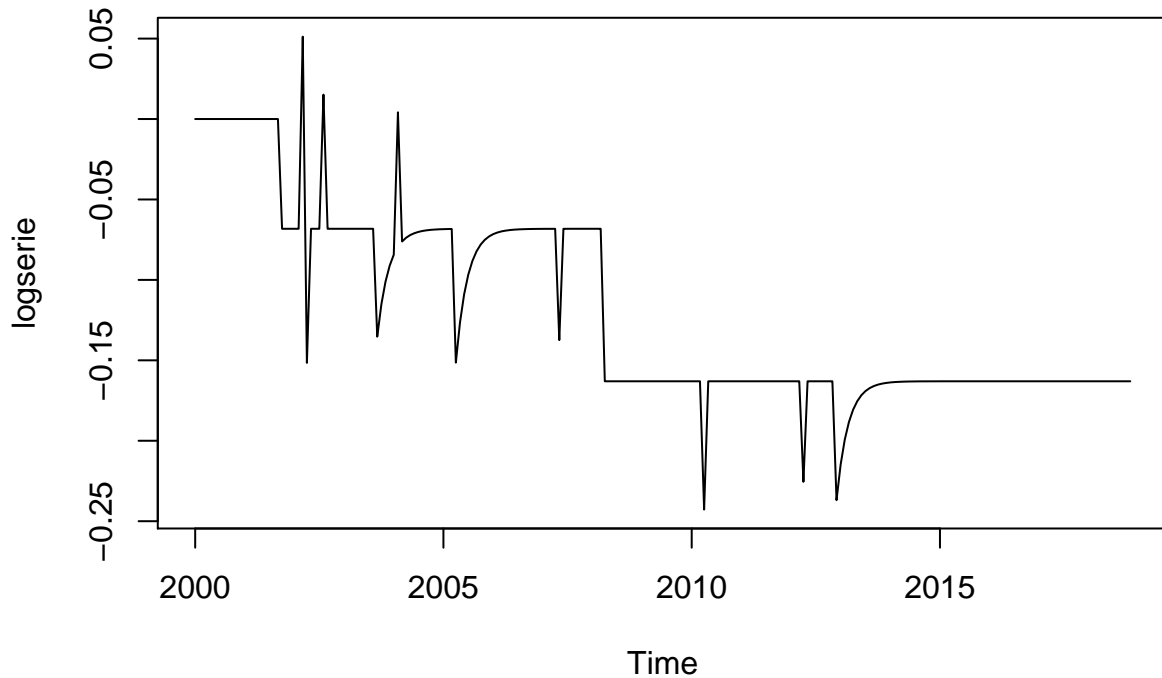
```
## Sigma^2 of the model with outliers:
```

```
## [1] 0.000959009
```

En el següent gràfic, es poden observar la gràfica amb els *outliers* (en negre) i la gràfica linealitzada, és a dir, sense *outliers* (en vermell). S'observa com hi ha outliers que han provocat que la sèrie tingui valors més baixos i outliers que han provocat que la sèrie tingui valors més alts. El fet que crida més l'atenció és que, a causa dels outliers, la sèrie des del 2005 (sobretot a partir del 2008) té valors més baixos del que hauria de tenir (s'observen els pics negres per sota dels vermells).



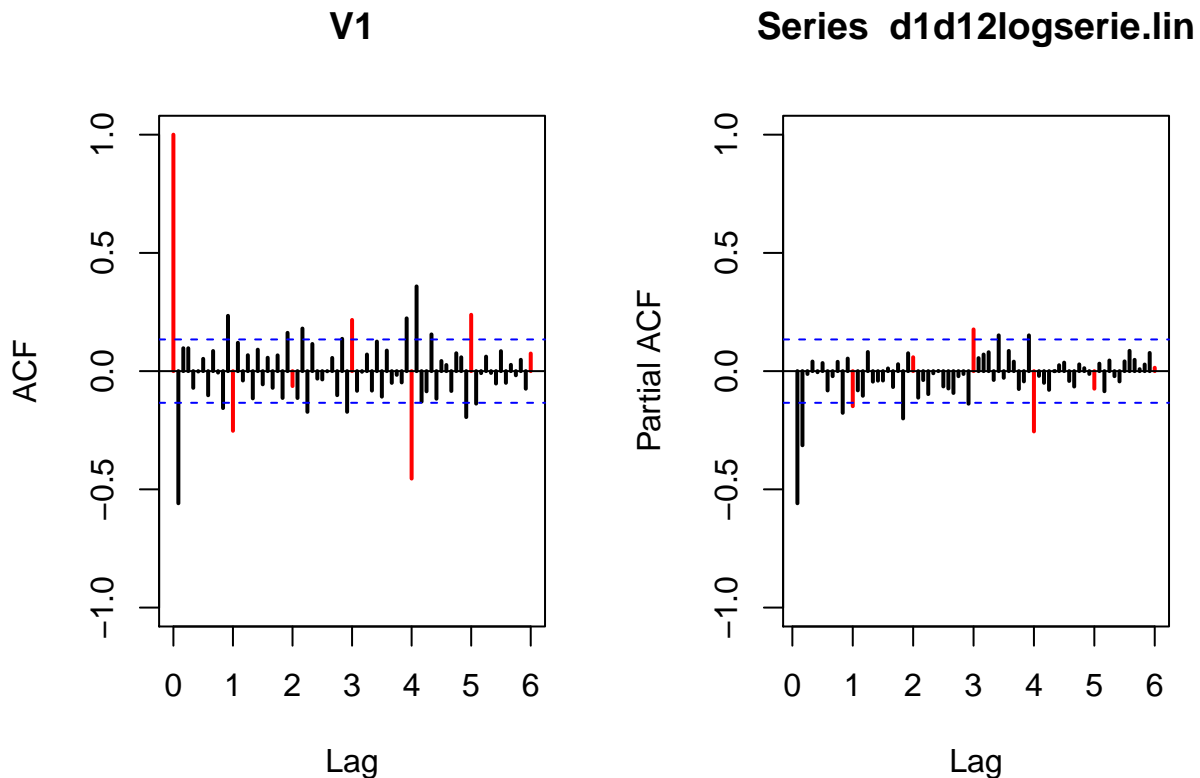
Per veure més clar l'efecte dels outliers, s'exposa la gràfica dels valors de la sèrie menys els valors de la sèrie linealitzada. Com es pot veure, l'abril del 2008 hi ha un *outlier* de tipus LS que fa que la sèrie agafi valors més petits des d'aquesta data en endavant. Crida l'atenció també els pics amunt i avall del principi de la sèrie, a les dates de març del 2002 i abril de 2002, que en només un més de diferència es va tenir una gran entrada turística el març i una baixa entrada turística l'abril. Per últim, cal remarcar un *outlier* de tipus LS a l'octubre de 2001, que també fa que des de llavors la sèrie prengui valors inferiors.



## 6.1 Identificació i estimació del model per la sèrie linealitzada

Un cop eliminats els *outliers* de la sèrie, es calculen un altre cop les funcions ACF i PACF de la sèrie linealitzada. S'observa clarament que, pel que fa a la part regular, al ACF es té decreixement exponencial alternat durant tots els valors, tinguent valors infinits fora de la banda de confiança als valors inicials. Pel que fa al PACF, observem que els dos primers valors són encara més significants que en el cas del PACF de la sèrie sense linealitzar. També hi ha altres valors fora de la banda que poden associar-se al cas del 5%. Per tant, pel que fa a la part regular, igual que en el cas de la sèrie amb valors atípics, es confirma la hipòtesi que el model adequat és un AR(2).

Pel que fa a la part estacional, l'anàlisi és molt semblant al de la sèrie sense linealitzar: un ACF amb força valors fora de la banda de confiança, sobretot el quart valor (molt més que el primer) i un PACF on es podrien considerar significatius el primer, el tercer i el quart valor. Per tant, igual que en el cas de la sèrie sense linealitzar, proposem un AR(4) per la part estacional.



Ara bé, a diferència del que s'ha vist anteriorment, en aquest cas surten no significatius tant el coeficient de `sar2` com el de `sar3`

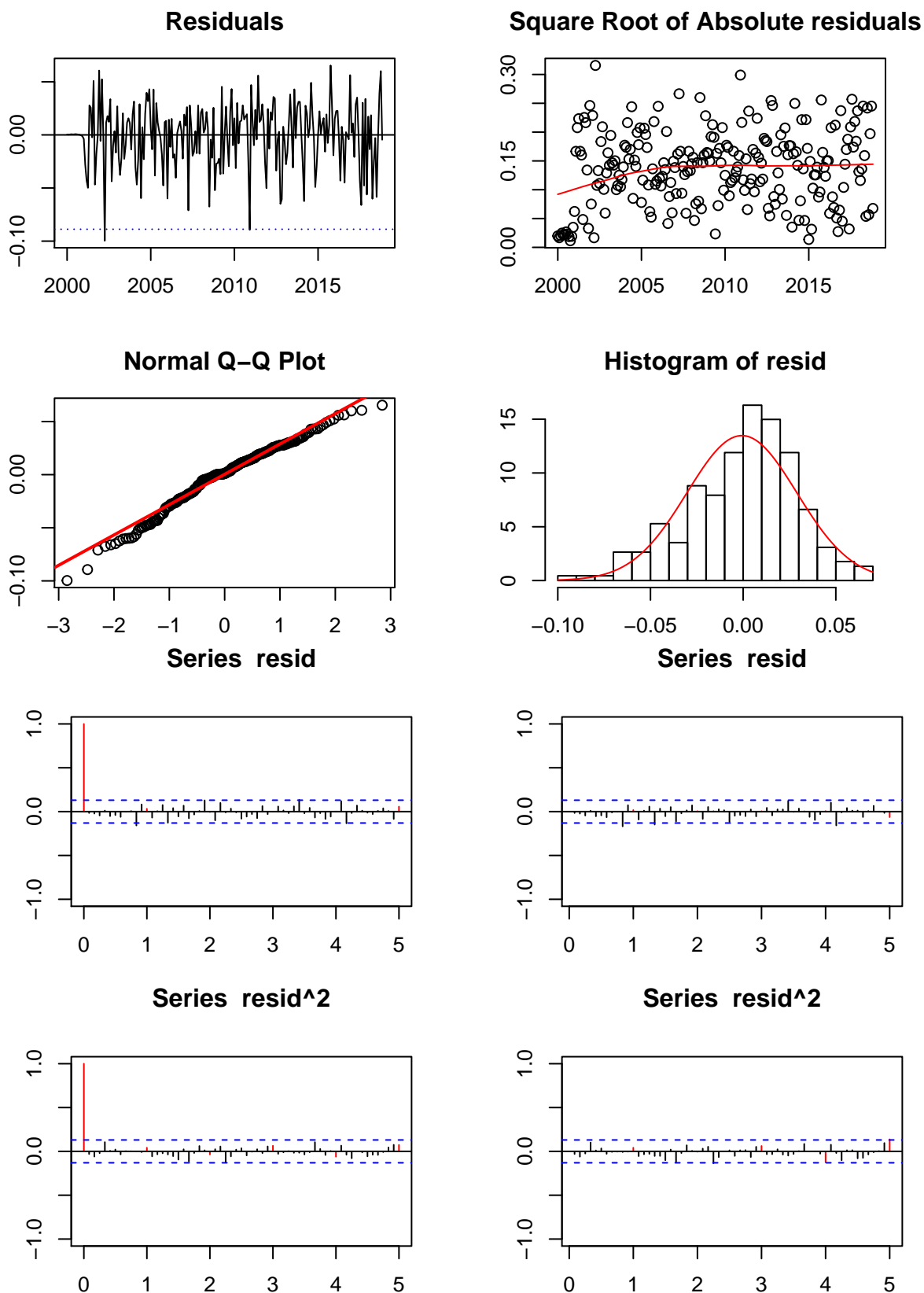
```
##
## Call:
## arima(x = logserie.lin, order = pdq.1, seasonal = list(order = PDQ.2, period = 12))
##
## Coefficients:
##          ar1      ar2      sar1      sar2      sar3      sar4
##      -0.6611 -0.3138 -0.1435 -0.1237  0.1143 -0.4575
## s.e.   0.0663  0.0662  0.0652  0.0726  0.0703  0.0684
##
## sigma^2 estimated as 0.0009259:  log likelihood = 436.23,  aic = -858.46
##   ar1  ar2  sar1  sar2  sar3  sar4
## TRUE TRUE TRUE FALSE FALSE TRUE
```

## 6.2 Validació del model per la sèrie linealitzada

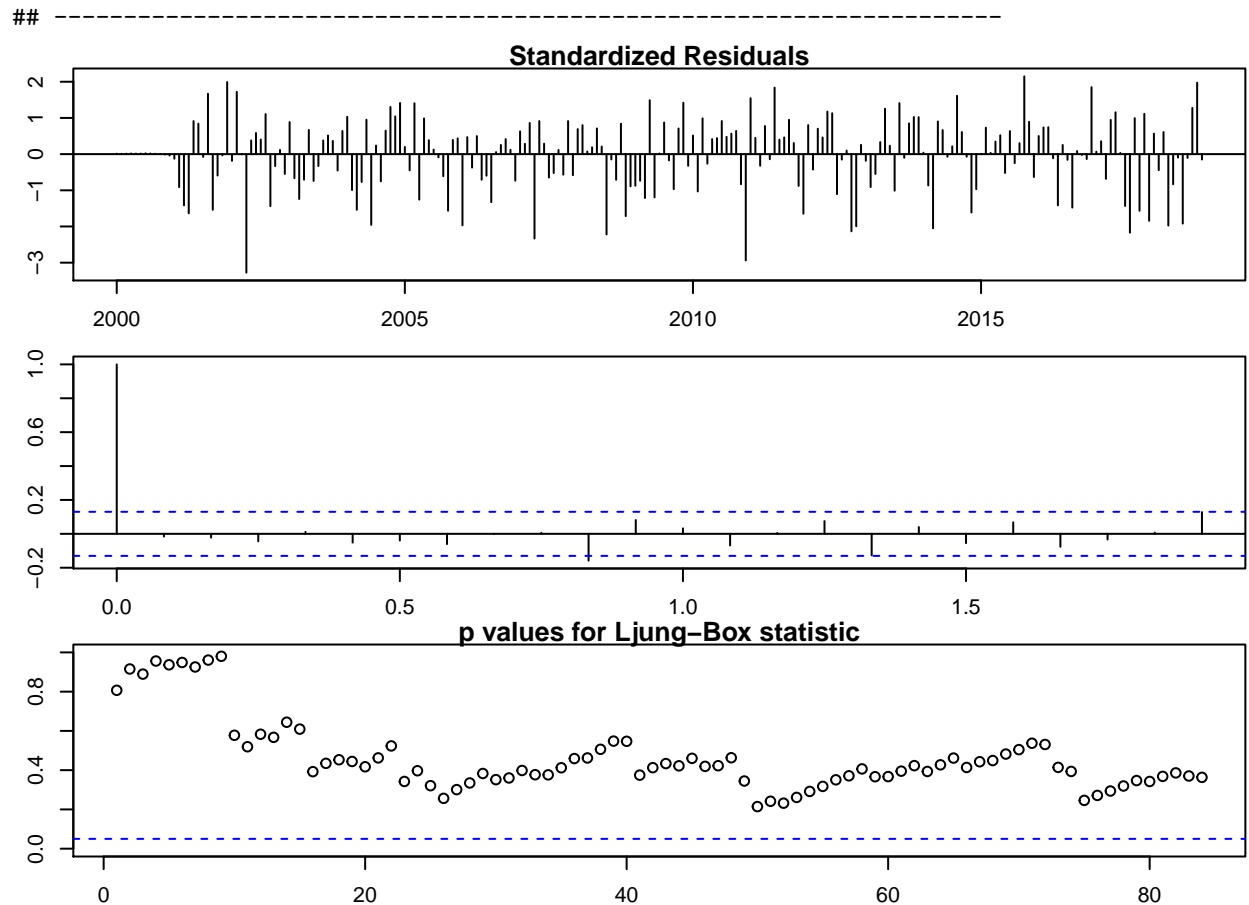
Pel que fa a la validació del model per la sèrie linealitzada, s'observa en les gràfiques els mateixos anàlisis realitzats anteriorment, és a dir:

- Es conclueix que es poden assumir les hipòtesis d'homogeneïtat en la variància dels residus (no patrons en les gràfiques de la seva variància i ACF i PACF dels residus al quadrat nuls), de normalitat dels residus (Q-Q plot amb relació lineal, histograma s'ajusta a la corba normal) i d'independència dels residus (*p-values* de Ljung-Box per sobre de 0.05 i ACF i PACF dels residus molt iguals).
- Es pot dir que el model és causal i invertible, ja que totes les arrels dels polinomis característics tenen mòdul major que un.
- El ACF i el PACF teòrics són molt semblants al ACF i PACF mostrals.

Per tant, es conclou que el model per la sèrie linealitzada és un model vàlid.



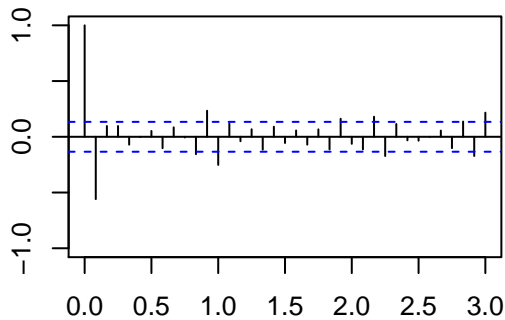
##



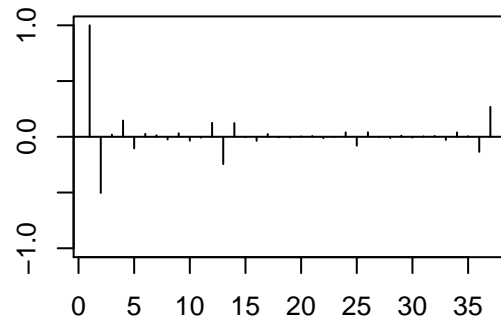
```
##
## -----
##
## Call:
## arima(x = logserie.lin, order = pdq.1, seasonal = list(order = PDQ.2, period = 12))
##
## Coefficients:
##      ar1      ar2      sar1      sar2      sar3      sar4
## -0.6611 -0.3138 -0.1435 -0.1237  0.1143 -0.4575
## s.e.   0.0663  0.0662  0.0652  0.0726  0.0703  0.0684
##
## sigma^2 estimated as 0.0009259:  log likelihood = 436.23,  aic = -858.46
##
## Modul of AR Characteristic polynomial Roots:
##
## Modul of MA Characteristic polynomial Roots:
##
## Psi-weights (MA(inf))
## -----
##      psi 1      psi 2      psi 3      psi 4      psi 5
## -0.6611482178  0.1232716501  0.1259974395 -0.1219912125  0.0411105666
##      psi 6      psi 7      psi 8      psi 9      psi 10
##  0.0111061928 -0.0202451983  0.0098994502 -0.0001911432 -0.0029805221
##      psi 11      psi 12      psi 13      psi 14      psi 15
```

```
## 0.0020305563 -0.1439246207 0.0945182259 -0.0173204886 -0.0182126923
##      psi 16      psi 17      psi 18      psi 19      psi 20
## 0.0174772433 -0.0058390801 -0.0016246536 0.0029067047 -0.0014118727
##
## Pi-weights (AR(inf))
##
## -----
##      pi 1      pi 2      pi 3      pi 4      pi 5      pi 6
## -0.66114822 -0.31384532 0.00000000 0.00000000 0.00000000 0.00000000
##      pi 7      pi 8      pi 9      pi 10     pi 11     pi 12
## 0.00000000 0.00000000 0.00000000 0.00000000 0.00000000 -0.14351754
##      pi 13     pi 14     pi 15     pi 16     pi 17     pi 18
## -0.09488637 -0.04504231 0.00000000 0.00000000 0.00000000 0.00000000
##      pi 19     pi 20
## 0.00000000 0.00000000
##
## Shapiro-Wilk normality test
##
## data: resid(model)
## W = 0.97956, p-value = 0.002271
```

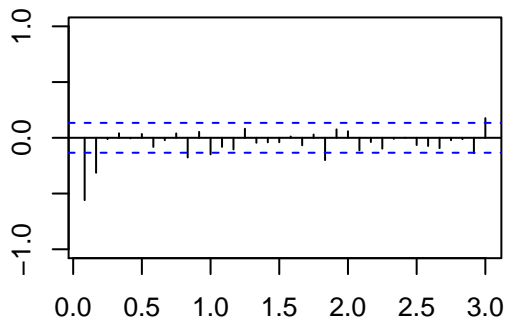
**Sample ACF**



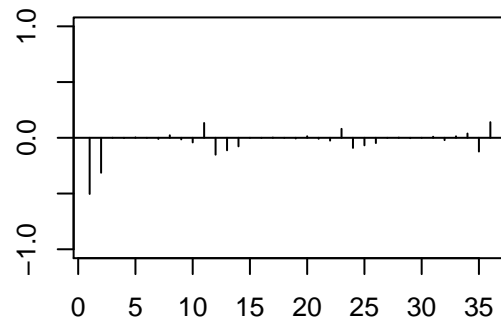
**ACF Teoric**



**Sample PACF**



**PACF Teoric**



### 6.3 Estabilitat del model proposat per la sèrie linealitzada

En relació a l'estabilitat del model per la sèrie linealitzada, s'observa que el valor dels coeficients varia molt poc, de l'ordre de menys de 0.02 en gairebé tots els casos. Per tant, podem confirmar que és estable.

```
##
```

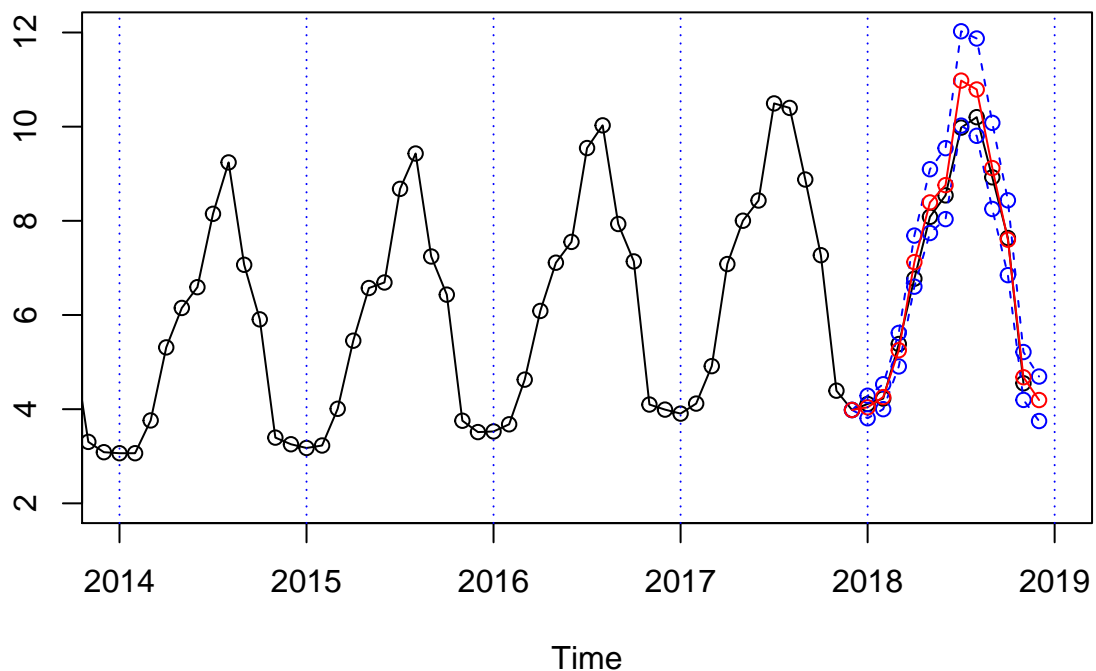


```
## Call:
## arima(x = logserie1.lin, order = pdq.1, seasonal = list(order = PDQ.2, period = 12))
##
## Coefficients:
##          ar1          ar2          sar1          sar2          sar3          sar4
##      -0.6611  -0.3138  -0.1435  -0.1237   0.1143  -0.4575
## s.e.    0.0663   0.0662   0.0652   0.0726   0.0703   0.0684
##
## sigma^2 estimated as 0.0009259:  log likelihood = 436.23,  aic = -858.46
##
## Call:
## arima(x = logserie2.lin, order = pdq.1, seasonal = list(order = PDQ.2, period = 12))
##
## Coefficients:
##          ar1          ar2          sar1          sar2          sar3          sar4
##      -0.6778  -0.3276  -0.1295  -0.1432   0.1111  -0.4494
## s.e.    0.0676   0.0682   0.0671   0.0742   0.0713   0.0705
##
## sigma^2 estimated as 0.000911:  log likelihood = 415.46,  aic = -816.91
```

## 6.4 Capacitat de predicció del model proposat per la sèrie linealitzada

Pel que fa a la capacitat de predicció del model per la sèrie linealitzada, es pot observar que és millor que el model per la sèrie sense linealitzar, ja que s'ajusta molt més (de fet, tenim tant el EQM com el EAM més baix en aquest cas). De fet, en les zones de pujada i baixada els intervals de confiança estan gairebé a sobre del valor real de la sèrie. En els valors més alts és on es té més error, fet que no extranya, ja que és en els mesos de l'any on més ha anat variant el valor de la sèrie al llarg dels anys

### Model ARIMA(2,1,0)(4,1,0)12

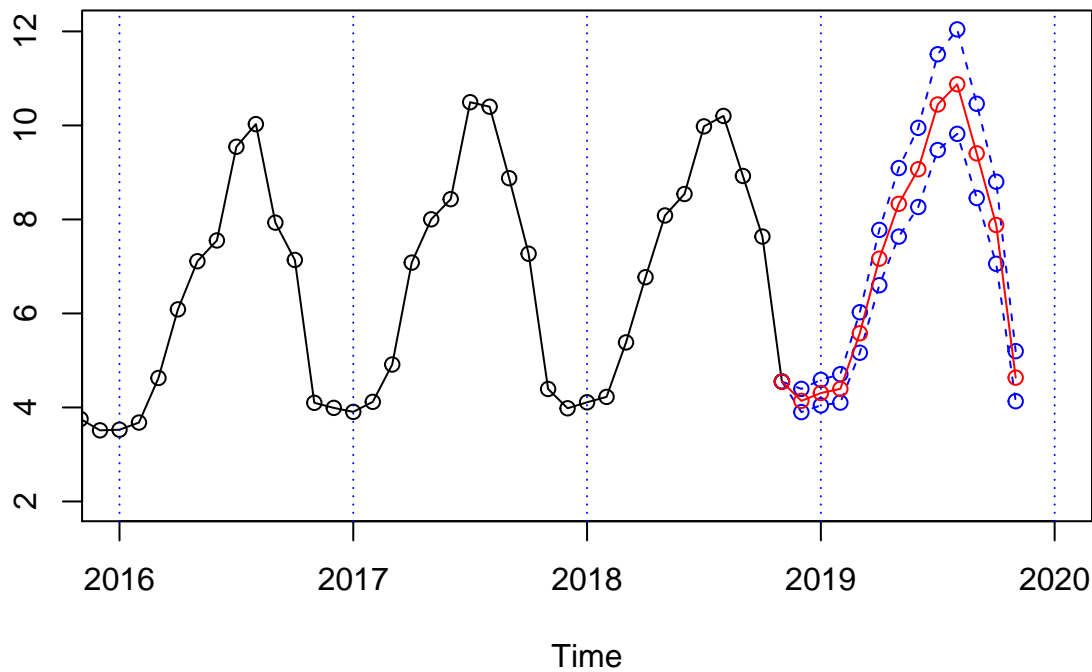


```
##          t1          pr          tu          serie          error
```

```
## Dec 2017 3.982530 3.982530 3.982530 3.982530 0.000
## Jan 2018 3.807304 4.039328 4.285492 4.110137 0.071
## Feb 2018 3.999645 4.256121 4.529043 4.224826 -0.031
## Mar 2018 4.907500 5.251276 5.619135 5.383687 0.132
## Apr 2018 6.601429 7.123301 7.686429 6.770845 -0.352
## May 2018 7.740691 8.391079 9.096113 8.084173 -0.307
## Jun 2018 8.036484 8.757653 9.543537 8.541181 -0.216
## Jul 2018 10.022139 10.978290 12.025662 9.979779 -0.999
## Aug 2018 9.805936 10.789681 11.872117 10.201456 -0.588
## Sep 2018 8.251762 9.120408 10.080496 8.924326 -0.196
## Oct 2018 6.846393 7.599599 8.435669 7.635569 0.036
## Nov 2018 4.196912 4.677580 5.213299 4.549899 -0.128
## Dec 2018 3.748080 4.193886 4.692716 NA NA
## [1] 0.04120862
## [1] 0.03142105
```

## 6.5 Previsions a llarg termini pel model proposat per la sèrie linealitzada

### Model ARIMA(2,1,0)(4,1,0)12



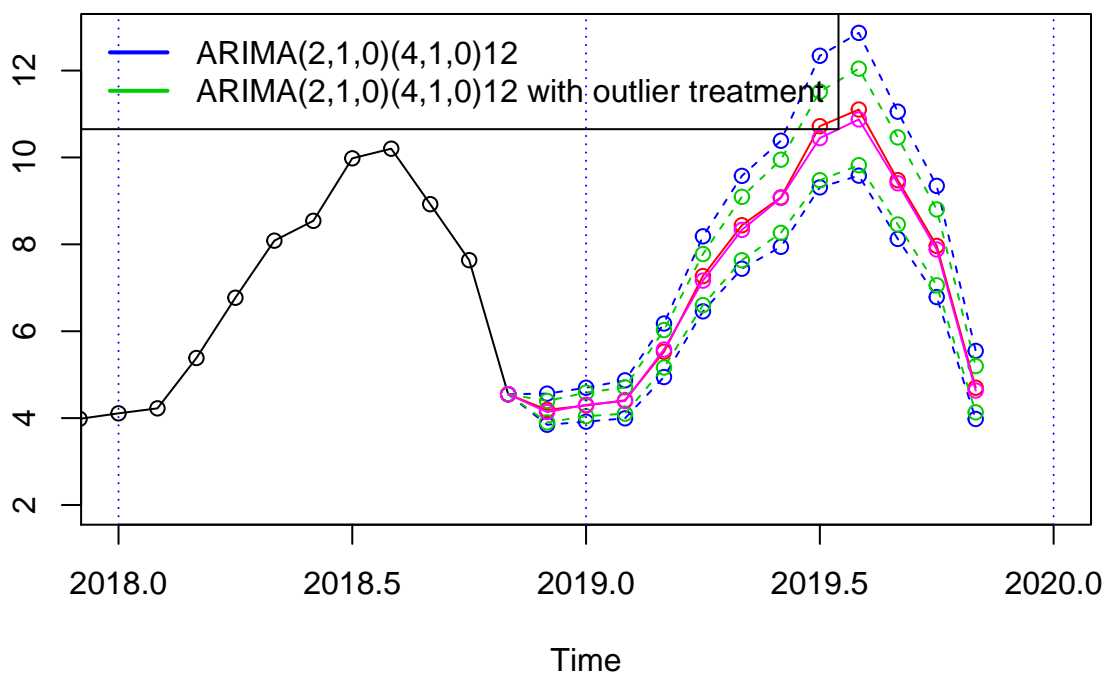
## 7 Comparació dels dos models

Com a última part d'aquest anàlisi, es realitza una comparació entre els valors predits pel model per la sèrie sense linealitzar i els valors predits per la sèrie linealitzada (és a dir, sense *outliers*). Tot i que els valors són força semblants, el que crida més l'atenció és que l'interval de confiança d'aquests valors és més estret en el cas del model per la sèrie linealitzada. Aquest fet es deu a que la variància dels valors sense *outliers* és molt més petita i, per tant, hi ha més marge d'acotació a l'hora de calcular tant els valors com els intervals. Per tant, podem dir que el model sense *outliers* és millor

```
##          t11          pr1          tu1
## Nov 2018 4.549899  4.549899  4.549899
## Dec 2018 3.846837  4.189916  4.563591
## Jan 2019 3.917671  4.290606  4.699042
## Feb 2019 3.995362  4.410566  4.868919
## Mar 2019 4.945079  5.526941  6.177267
## Apr 2019 6.457154  7.269534  8.184119
## May 2019 7.438742  8.439524  9.574947
## Jun 2019 7.942888  9.081850 10.384132
## Jul 2019 9.311565 10.719334 12.339936
## Aug 2019 9.580515 11.102785 12.866932
## Sep 2019 8.125041  9.476734 11.053296
## Oct 2019 6.786354  7.963653  9.345191
## Nov 2019 3.980613  4.698813  5.546594
```

```
##          t12          pr2          tu2
## Nov 2018 4.549899  4.549899  4.549899
## Dec 2018 3.900358  4.140052  4.394477
## Jan 2019 4.044126  4.306979  4.586916
## Feb 2019 4.102884  4.394829  4.707548
## Mar 2019 5.164192  5.578496  6.026037
## Apr 2019 6.600659  7.164904  7.777382
## May 2019 7.634515  8.332016  9.093241
## Jun 2019 8.264364  9.067381  9.948425
## Jul 2019 9.476601 10.445786 11.514091
## Aug 2019 9.821204 10.875595 12.043183
## Sep 2019 8.456203  9.405362 10.461059
## Oct 2019 7.056392  7.881295  8.802630
## Nov 2019 4.132904  4.634817  5.197685
```

## Entrada Turística a Espanya



Per acabar de confirmar-ho, es mostren tota una sèrie de mesures de bondat d'ajust dels models. Sobretot

ens crida l'atenció el  $AIC$  i el  $BIC$ , on el model sense *outliers* és clarament millor. Aquest model també té millor RMSPE i MAPE.

## 8 Comentaris finals

Així doncs, un cop realitzat aquest anàlisi, es conclou que la presència de valors atípics en una sèrie pot influir (i molt) en les previsions que pugui fer un model basant-se en ella. Per tant, es confirma la importància de la seva detecció i correcció. Pel que fa a les previsions de la sèrie, sembla que la tendència general de creixement es mantindrà, així com també el patró estacional (és a dir, que Espanya tenint un munt de turistes a l'estiu).