

# A Time Series Analysis: Total Exports from Spain

Joel Solé and Carlos Arbonés

GCED, UPC.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Data description</b>	<b>3</b>
<b>3</b>	<b>Box-Jenkins methodology</b>	<b>3</b>
<b>4</b>	<b>Analysis of the time series</b>	<b>4</b>
4.1	Identification . . . . .	4
4.1.1	Analysis of Variance . . . . .	4
4.1.2	Seasonality . . . . .	4
4.1.3	Constant mean . . . . .	5
4.1.4	ACF and PACF . . . . .	6
4.2	Estimation . . . . .	7
4.3	Validation . . . . .	8
4.3.1	Residual Analysis . . . . .	8
4.3.2	Model properties . . . . .	9
4.3.3	Goodness-of-fit measures . . . . .	11
4.3.4	Stability . . . . .	11
4.3.5	Forecasting Capacity and model selection . . . . .	12
4.4	Forecasting . . . . .	13
4.5	Calendar effects . . . . .	13
4.5.1	Applying Box-Jenkins to $X_t^*$ . . . . .	14
4.6	Outlier treatment . . . . .	14
<b>5</b>	<b>Appendix</b>	<b>17</b>

## 1 Introduction

The purpose of this project is to conduct a comprehensive analysis of the monthly total exportations in Spain and develop a reliable forecasting model that can predict its future behavior with a high degree of accuracy. Time series data is a sequential collection of observations that are recorded over time and can be used to identify patterns, trends, and other underlying characteristics.

This project will utilize data analysis techniques, statistical models, and time series analysis methodologies to gain insights into the behavior of the time series and develop a forecasting model that can provide reliable predictions of future values. We will conduct a thorough review, develop hypotheses based on our initial observations, and test these hypotheses using statistical tests and other relevant analytical methods.

In particular, we will employ the Box-Jenkins methodology, a widely used approach to time series analysis that involves a series of steps for model identification, parameter estimation, and model validation. This methodology provides a structured and systematic approach to analyzing time series data and has proven to be effective in producing accurate forecasts.

The success of this project will be measured by the quality of the forecasting model that we develop, and its ability to accurately predict future values of the time series.

## 2 Data description

The time series of total exportations in Spain is a monthly series that shows the total value of goods exported by Spain in billions of euros at current prices and without seasonal adjustment.

The data source is the *Base de Datos de Análisis Sectorial (BADASE)* [1], which is a database maintained by the Ministry of Industry, Commerce and Tourism of Spain that provides statistical information on various economic sectors.

The sample size is the total number of transactions reported by customs authorities, which is 240 total observations. The series covers the period from January 1999 to December 2018.

## 3 Box-Jenkins methodology

The Box-Jenkins method, is a systematic approach to time series analysis and forecasting. The method consists of four key steps: identification, estimation, diagnostic checking, and forecasting as it can be seen in figure 5.1.

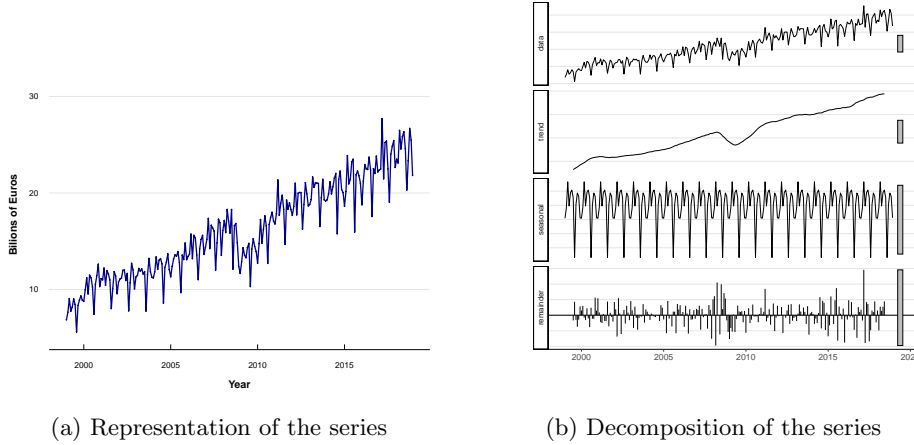
The first step, identification, involves analyzing the autocorrelation and partial autocorrelation functions to determine the order of the autoregressive, differencing, and moving average terms that are appropriate for the time series data. Once the order of the terms has been determined, the next step is to estimate the model parameters using maximum likelihood estimation.

After estimating the model parameters, the third step is to perform diagnostic checking by analyzing the residuals to ensure that they are independent, normally distributed with zero mean and constant variance. If the residuals do not meet these criteria, the model may need to be re-specified.

Finally, the model can be used to forecast future values of the time series. Additionally, other considerations such as seasonality, outliers, and non-stationary must be taken into account to build an accurate and reliable model.

## 4 Analysis of the time series

In this chapter we will apply the methods before mentioned and give a detail explanation of the reasoning behind the steps in the analysis.



**Fig. 4.1:** Visualization of the time series

We can observe that from 1999 (start of the series) until 2008, there is a clearly increasing linear trend. In 2008, there was a decrease in Spanish exports due to the crisis. Following this event, there was another similar linear growth to what existed prior to the crisis. The existence of negative peaks during the summer months related to a decrease in exports due to vacations is clear.

### 4.1 Identification

Following the Box-Jenkins methodology, the first step is to convert our series  $X_t$  into a stationary<sup>1</sup> series  $W_t$ .

#### 4.1.1 Analysis of Variance

The variance of the series appears to increase as the values increase as it can be seen in figure 4.2a. In the box plot (figure 5.2a) it appears that the higher values in the series have greater variance, although we can also observe smaller values with higher variance. The sizes of the boxes (*IQR*) vary through time, which indicates different levels of dispersion. Is clear that the variance is not constant.

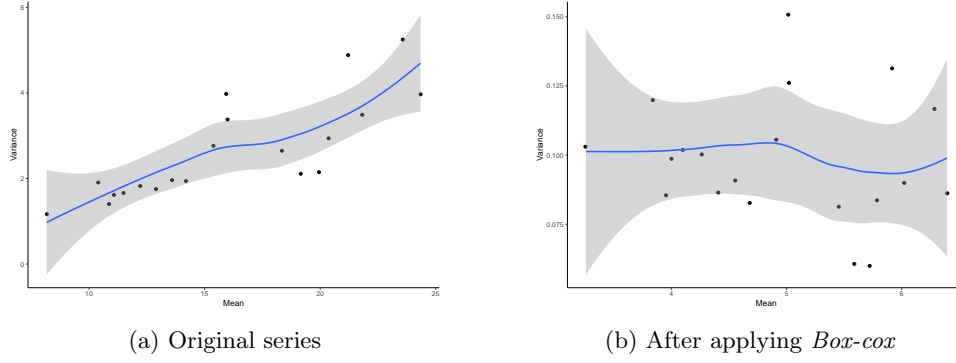
In order to achieve homoscedasticity, we will apply a [Box-Cox](#) transformation to the values of the series. We notice an improvement, as we see in figure 4.2b the is no a pattern between the values of the series and the variance, what leads to think the constant variance is achieved.

#### 4.1.2 Seasonality

Next step is to remove the seasonality of the series. As it has been mentioned before there is a clear seasonality pattern in the series as in the summer months there is a remarkable decrease in the exportations.

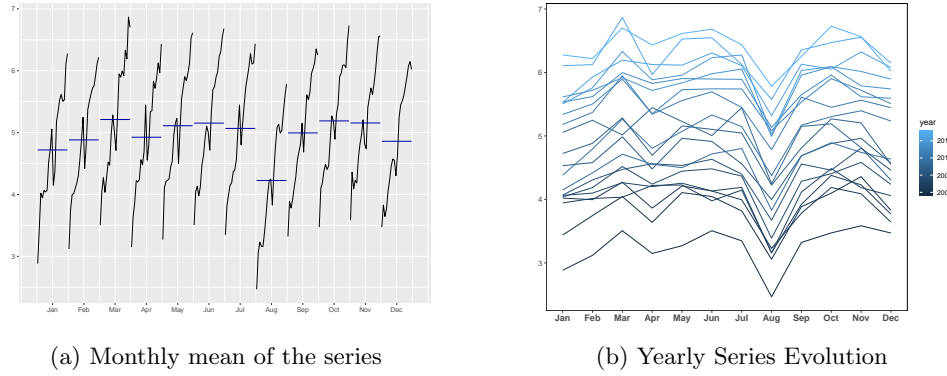
---

<sup>1</sup>An stationary series is a series with constant variance, constant mean and auto-correlation structure independent from the origin.



**Fig. 4.2:** Relation Mean vs Variance

It can be seen more clearly in figure 4.3a, where the monthly mean of the series is not constant meaning there is seasonality in the series as the values are different for each month. In figure 4.3 we can see the yearly evolution and notice again the peak in August and a clearly pattern.



**Fig. 4.3:** Seasonality of the original series

To remove the seasonality it is applied a seasonal differentiation<sup>2</sup>. As it can be seen in figure 5.3a, the seasonal pattern has been removed, the monthly mean is constant and there is no clear pattern as we can verify in figure 5.3b.

#### 4.1.3 Constant mean

The last step to convert the series in stationary is to ensure it has a constant mean. After applying a *Box-cox* transformation and a seasonal differentiation to the series we have

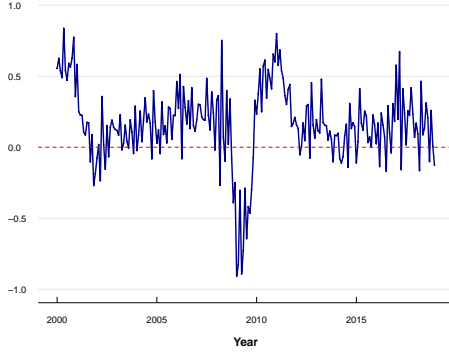
$$W'_t = (1 - B^{12})\left(\frac{X_t^\lambda - 1}{\lambda} - \mu\right)$$

---

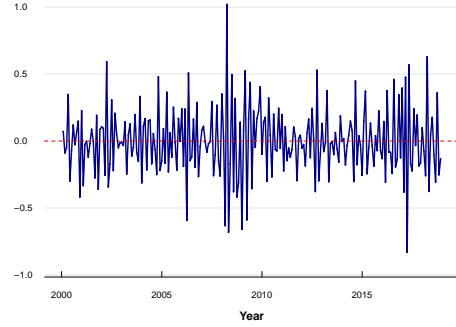
<sup>2</sup>Seasonal differentiation consist in subtracting to each observation the observation from the previous year.

The current series has no constant mean as it can be seen in figure 4.4, thus it is necessary to apply regular differentiations<sup>3</sup> to the series until a constant mean is achieved. After applying a regular differentiation to  $W'_t$  we obtain:

$$W_t = (1 - B)(1 - B^{12})\left(\frac{X_t^\lambda - 1}{\lambda} - \mu\right)$$



**Fig. 4.4:** Representation of  $W'_t$



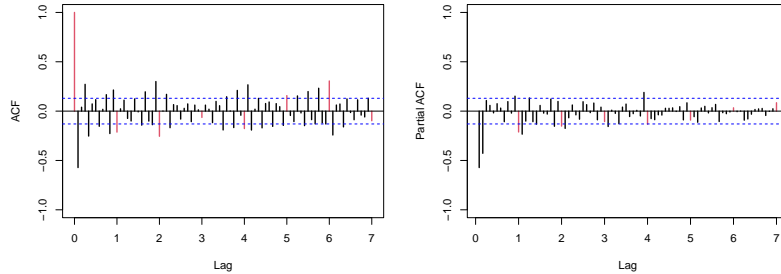
**Fig. 4.5:** Representation of  $W_t$

From Figure 4.5, it is evident that the series displays a stationary characteristic with a constant mean, as indicated by the clustering of values around zero. Nevertheless, we decided to attempt a different level of differentiation on the series. In case this leads to an increase in variance, we will revert to the previous level of differentiation, as over-differentiating the series may not be appropriate.

As seen in table 7, applying an additional differentiation to the series results in a higher variance. Hence, we chose to stick with the previous order of differentiation, which is the series denoted by  $W_t$ .

#### 4.1.4 ACF and PACF

We will examine the Auto-Correlation Function (ACF) and Partial Auto-Correlation Function (PACF) of the time series denoted by  $W_t$ , with the aim of proposing appropriate models. We are utilizing the  $ARIMA(p, d, q)(P, D, Q)_{12}$  modeling, where  $d$  and  $D$  denote the number of regular and seasonal differentiations, respectively.



**Fig. 4.6:** ACF and PACF of  $W_t$

<sup>3</sup>Regular differentiations consist in subtracting to each observation its previous one.

### Seasonal Lags

On the one hand, it is clear that in the PACF, the seasonal lags follow an exponential decay pattern, typical of an  $MA(q)$  model. In this case, we identify an  $MA(6)$ , as we observe a significant 6<sup>th</sup> seasonal lag that we want to capture, and therefore assume that after this, the rest of the lags are null.

On the other hand, a model with fewer parameters where the infinite lags are in the ACF and the finite lags are in the PACF seems more evident. In this case, we identify an  $AR(4)$ , considering that from the 4<sup>th</sup> lag onwards, all are null, as they fall within the confidence bands.

Finally, we identify patterns of decay in both the ACF and the PACF, thus an  $ARMA(1,1)$  model can also be identified.

### Regular Lags

We can identify patterns of decay in the ACF where many significant lags are observed far from the origin. Therefore, it will be an  $AR(2)$  process since we clearly observe 2 significant lags and consider the rest to be null since they are within the confidence bands. Additionally, there were significant lags near the seasonal ones, but these were deemed to be satellite lags and were not taken into account.

Based on the previous analysis the models<sup>4</sup> identified for  $X_t$  are:

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

$$ARIMA(2, 1, 0)(0, 1, 6)_{12}$$

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

## 4.2 Estimation

In this section, we will [estimate](#) the parameters of the first two models. To perform the estimation correctly, we will follow the following methodology for both models. First, we will fit the model on the stationary series  $W_t$  to estimate the mean of the series. In case the mean is not [significant](#), we will proceed with the fitting on the model with only the original transformation<sup>5</sup>, so that it will be easier later on to predict and obtain the original values of the series. Next, we will estimate the parameters of the model to check if they are significant. If any parameter is not significant, we will eliminate it, refit the model and keep the model with less parameters in case the AIC improves. We will do this until all the remaining parameters of the model are significant.

### $ARIMA(2, 1, 0)(4, 1, 0)$

We obtained that the mean is not significant, and therefore we can consider  $\mu = 0$ . After estimating the parameters, we can see that all of them are significant, and the model we obtain is:

$$(1 + 0.076B + 0.38B^2)(1 + 0.59B^{12} + 0.64B^{24} + 0.45B^{36} + 0.4251B^{48})W_t = Z_t$$

---

<sup>4</sup>All the ARIMA models we will identify through this report will have a period of 12 observations, but will not be indicated for notation simplicity

<sup>5</sup>Original transformation refers to the series with constant variance, that is, with *Box-cox* transformation applied.

### ***ARIMA(2, 1, 0)(0, 1, 6)***

Based on our hypothesis test, we found that the mean is not significant. After re-estimating the parameters, we observed that the third, fourth, and fifth seasonal parameters have  $|\hat{t}| < 2$ , indicating they are not significant. Therefore, we decided to repeat the process without the fifth parameter, as it had the lowest t-test value. After re-estimating the model, we found that the remaining parameters are significant, and the model we obtain is:

$$(1 + 0.80B + 0.44B^2)W_t = (1 - 0.90B^{12} - 0.27B^{24} + 0.22B^{36} + 0.17B^{48} - 0.24B^{72})Z_t$$

## **4.3 Validation**

In this section, we will perform the validation of the two previously estimated models. Specifically, we will review whether the residuals of the models meet the corresponding assumptions, whether the models meet the properties of invertibility and causality, and whether the models are stable and allow for the generalization of the obtained predictions measurements.

### **4.3.1 Residual Analysis**

For each of the models, we will conduct a complete residual analysis and verify that they meet the assumptions of homoscedasticity (constant variance), normality, and, most critically, independence. In the event that the residuals do not meet the independence assumption, we will need to re-identify the model to better capture the autocorrelation structure present in the data.

#### ***Homocedasticity***

To check for constant variance, we will first create a plot of the residuals with probability marks referring to  $3 \cdot \sigma$ . In this plot, we need to verify that the variability of the residuals over time is constant and check for any outliers (residuals that fall outside the probability marks). We will also create a smooth estimate of the square root of the residuals' mean to verify that the variance is constant. We should observe that the mean follows an horizontal line.

#### ***Normality***

We want the residuals to follow a normal distribution. To check this, we will use a QQ-plot where the theoretical quantiles of the normal distribution are plotted against the residuals' quantiles. The residuals should lie along the line. We will also plot a histogram of the residuals with the normal density function plotted on top to see if it fits. Finally, we will perform the [Shapiro-Wilk](#) test to verify normality.

#### ***Independence***

To verify that the residuals are independent, that is, that there is no autocorrelation structure in the residuals and therefore the model captures all the information present in the data, we will perform the ACF and PACF of the residuals. We should see that there are no significant lags (they should be within the confidence bands), and if there is any, it should be far from the origin. We will also calculate the p-values of the [Ljung-Box test](#) since checking the significance of individual autocorrelation may ignore that the configuration of all lags together could be significant.



### ***ARIMA(2, 1, 0)(4, 1, 0)***

The variance of the residuals appears to be constant as we can see in figure 5.4a, where the variability of the values is constant over time. We can also detect an outlier in the year 2008. In figure 5.4b, we can see that the line is almost horizontal and therefore, we can conclude that the variance is constant.

Regarding normality, we can observe in the QQ-plot (figure 5.4c) how the residuals fit quite well on the line, which gives us evidence that the residuals are normally distributed. In the histogram (figure 5.4d), seem that the residuals are normally distributed, since they fit the normal density function, although and we can also observe large values in the tails. Finally, the *p-value* for the *Shapiro-Wilk* test gives us 0.0471, suggesting that there is not enough evidence to claim that the residuals are normally distributed. The graphical results are sufficiently convincing to validate the hypothesis of normality of the residuals.

In the ACF and PACF plots of the residuals (5.4e), we can observe that all lags are not significant as they are within the confidence intervals, and the lags closer to the origin are close to zero. This indicates that there is no autocorrelation structure in the residuals and they are independent. By performing the LjungBox test (5.4f), we observe that the *p-values* are above 0.05, which indicates that there is not enough evidence to reject the null hypothesis of randomness of the residuals, and therefore the independence hypothesis is validated.

### ***ARIMA(2, 1, 0)(0, 1, 6)***

Based on the plot in figure 5.5a, we can observe that the residuals' variance remains stable across time, indicating constant variability. Additionally, there appears to be an outlier in the year 2008. Moreover, figure 5.5b illustrates that the blue line remains almost horizontal. Therefore, we can conclude that the variance is constant.

The QQplot (5.5c) shows that the residuals are located on the line, although there are some points at the beginning that are not. The histogram (5.5d) does not appear to follow a normal distribution, as it does not fit completely the density function. Finally, the *Shapiro-Wilk* test yields a very low *p-value* of 0.006, therefore the residuals are not normally distributed.

In the ACF and PACF plots (5.5e), the nearest lags to the origin are not significant. There is some significant lag, but it is very far from the origin, indicating that the residuals are independent. Nevertheless, we will confirm this with the *Ljung-Box Test* (5.5f). As we can see, the *p-values* are above 0.05, except in some specific cases but quite far from the origin. Therefore, the hypothesis of independence is validated.

## **4.3.2 Model properties**

In order to determine whether the models are invertible and/or causal, we need to check these properties as they allow us to convert the models into *AR*( $\infty$ ) and *MA*( $\infty$ ) models. If the properties are not satisfied, the weights  $\pi$  and  $\psi$  will not tend to zero, and we won't be able to truncate the expressions, making it difficult to predict the future as the distant observations will be given too much importance.

For the model to be invertible, we must satisfy the condition that  $\sum_{i=0}^{\infty} \pi^2 < \infty$ . In other words, the modulus of all roots of the polynomial  $\theta(B)$  must be greater than 1.

For the model to be causal, we must satisfy the condition that  $\sum_{i=0}^{\infty} \psi^2 < \infty$ . This is equivalent to saying that the modulus of all roots of the polynomial  $\phi(B)$  must be greater than 1.

### ***ARIMA(2,1,0)(4,1,0)***

To calculate the modulus of the roots of the polynomial:

$$\phi(B) = (1 + 0.076B + 0.38B^2)(1 + 0.5954B^{12} + 0.6411B^{24} + 0.4560B^{36} + 0.4251B^{48})$$

We obtain that all module are greater than 1 and therefore the model is causal. Since it is an AR model, it will always be invertible. In figure 5.6a we can see the inverse of the modulus of the roots, and we can verify that all of them are inside the unit circle. This means that we can express the model as an  $AR(\infty)$  and as an  $MA(\infty)$  since the weights tend to 0. In figures 5.6b and 5.6c, we can observe how the weights of the  $\pi$  and  $\psi$  parameters tend towards 0, respectively, which is in line with what we have seen previously.

The expression for the  $AR(\infty)$  is:

$$(1 - 0.7602B - 0.3825B^2 - 0.5953B^{12} - 0.4526B^{13} - 0.2277B^{14} - \dots)X_t = Z_t$$

The expression for the  $MA(\infty)$  is:

$$X_t = (1 - 0.7602B + 0.1954B^2 + 0.1422B^3 - 0.1828B^4 + 0.0846B^5 - \dots)Z_t$$

### ***ARIMA(2,1,0)(0,1,6)***

We have the polynommials:

$$\phi(B) = (1 + 0.8024B + 0.4485B^2)$$

$$\theta(B) = (1 - 0.9014B^{12} - 0.2779B^{24} + 0.2242B^{36} + 0.1711B^{48} - 0.2405B^{72})$$

From  $\phi(B)$ , we obtain that all modules are greater than 1 and therefore the model is causal. When calculating the module of the roots of the second polynomial  $\theta(B)$ , we see that the smallest module is equal to  $0.9986 < 1$  (figure 5.7a), which means that the model is not invertible.

In figure 5.7b, we can observe how the weights  $\pi$  do not tend towards zero, as the model is not invertible, and therefore we cannot truncate the expression. In figure 5.7c, we can see how the weights  $\psi$  do tend towards zero due to the causality of the model.

The expression for the  $AR(\infty)$  is:

$$(1 - 0.8023B - 0.4484B^2 - 0.9014B^{12} - 0.7232B^{13} - 0.4042B^{14} - \dots)X_t = Z_t$$

The expression for the  $MA(\infty)$  is:

$$X_t = (1 - 0.8023B + 0.1953B^2 + 0.2031B^3 - 0.2505B^4 + 0.1099B^5 - \dots)Z_t$$

### 4.3.3 Goodness-of-fit measures

	AIC	BIC	$\sigma_z^2$
<b>ARIMA(2,1,0)(4,1,0)</b>	-191.74	-167.77	0.021
<b>ARIMA(2,1,0)(0,1,6)</b>	-191.92	-164.52	0.019

**Table 1:** Comparison of adequacy metrics

The AIC is very similar for both models, although it is slightly better for the second model. It seems that both models fit the data similarly well. It was expected that the first model would have a better BIC than the second model, since the BIC selects more parsimonious models by penalizing the number of parameters more heavily, and in this case the second model has one more parameter. However, the BIC is also similar for both models.

### 4.3.4 Stability

In this section, we will verify if the models are stable. A model is stable if it maintains its predictive ability across different datasets, that is, if it is practically the same with or without the last  $h$  observations. It is important to verify stability to ensure that the metrics used to evaluate the forecasting ability (RMSPE and MAPE) according to the second model<sup>6</sup> are extrapolatable to the first model, which is the one we will use to calculate future predictions.

We will estimate both models without the last 12 observations and compare them with their respective original models. We will conclude that the model is stable if the parameters have the same sign, magnitude, and significance.

#### *ARIMA(2,1,0)(4,1,0)*

$m$ coefficients	-0.760	-0.383	-0.595	-0.641	-0.456	-0.425
$m_{12}$ coefficients <sup>1</sup>	-0.749	-0.371	-0.591	-0.649	-0.449	-0.420
Difference ( $m - m_{12}$ )	0.011	0.012	0.004	0.008	0.007	0.005
T-ratios ( $m$ )	-11.57	-5.70	-9.44	-9.01	-6.60	-6.44
T-ratios ( $m_{12}$ )	-11.15	-5.36	-8.75	-8.87	-6.29	-6.21

<sup>1</sup>  $m_{12}$  refers to the model estimated without using the last 12 observations.

**Table 2:** Comparison of coefficients

We can see that the magnitude of the coefficients is very similar, if we look at the difference between the coefficients of both models, it is very small, and they all have the same sign. We can also see that they have the same significance since the *T-ratios* are also similar. Therefore, we conclude that the model is stable

---

<sup>6</sup>Second model refers to the model estimated without the last  $h = 12$  observations.

### ARIMA(2,1,0)(0,1,6)

$m$ coefficients	-0.802	-0.448	-0.901	-0.278	0.224	0.171	0	-0.241
$m_{12}$ coefficients	-0.789	-0.441	-0.923	-0.283	0.237	0.183	0	-0.256
Difference ( $m - m_{12}$ )	0.014	0.008	0.021	0.005	0.013	0.012	0	0.016
T-ratios ( $m$ )	-13.07	-7.30	-7.78	-2.94	2.38	2.18	0.00	-3.92
T-ratios ( $m_{12}$ )	-12.61	-6.97	-6.09	-2.82	2.45	2.14	0.00	-4.09

**Table 3:** Comparison of coefficients

We can observe that the magnitude of the coefficients is very similar, and they all have the same sign. Additionally, all of the coefficients have the same significance, so we can conclude that the model is stable.

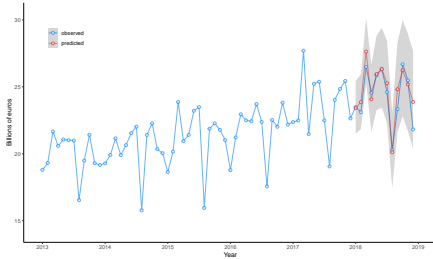
#### 4.3.5 Forecasting Capacity and model selection

To compute the confidence intervals of the predictions, we will take a significance level of  $\alpha = 0.05$ , with a 95% confidence level. Therefore, we have  $Z_{1-\alpha} = 1.96$ . We can write:

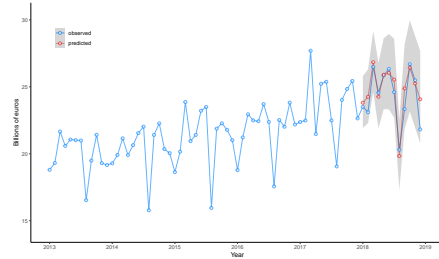
$$IC_{95\%}(X_{t+h}) = \tilde{X}_{t+h|t} \pm 1.96 \cdot \sqrt{Var(\tilde{X}_{t+h|t})}$$

To obtain the original values of the series after the predictions, we need to perform the inverse operation of the *Box-cox* transformation:

$$X_t = (\lambda \cdot X + 1)^{\frac{1}{\lambda}}$$



(a) ARIMA(2,1,0)(4,1,0)



(b) ARIMA(2,1,0)(0,1,6)

**Fig. 4.7:** Predictions for the 2 models

We can observe that the predictions of both models are very good, with the pointwise predictions fitting the actual observations very well. In both models, we see that the December prediction is the least accurate.

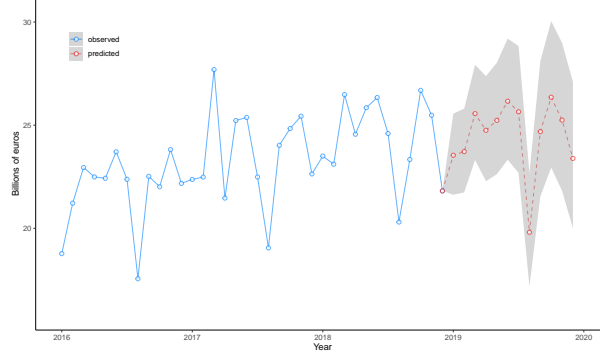
Model	RMSE	MAE	RMSPE	MAPE	CI
ARIMA(2,1,0)(4,1,0)	0.8812	0.6477	0.0381	0.0273	5.8040
ARIMA(2,1,0)(0,1,6)	0.9309	0.6798	0.0411	0.0292	5.5249

**Table 4:** Performance metrics comparison

The first model outperforms the second in all measures of forecasting accuracy, as we can see in the metrics where the errors are slightly smaller. Despite having larger confidence intervals, we will choose the first model for making predictions.

## 4.4 Forecasting

From the model we selected, we can obtain the forecast of the exportations in Spain for the next 12 months, as well as the 95% confidence intervals of the forecasts. The resulting forecast is presented in 4.8.



**Fig. 4.8:** Forecast for the 2019 exportations

## 4.5 Calendar effects

Next, we will analyze the calendar effects of the Holy Week and Working Days on the series and see if they are significant.

Depending on which month Holy Week falls on, it can affect the predictions. If it falls in March one year, exports during that month will likely decrease and be noticeable in the series. If the following year it falls in April, the prediction for March will be lower than it should be. In order to control for these effects, we will distribute the effect of the Holy Week between the months of March and April.

Regarding the effect of working days, we will maintain a proportion of working days to holidays of 5/2. This way, if for example a month does not meet this ratio, we will add/remove working days depending on whether it has fewer/more than the required number.

The resulting series with the calendar effects removed will be as follows:

$$X_t^* = X_t - w_{TD}TD_t - w_{Ea}Ea_t$$

To obtain an estimation for both coefficients, we can fit all possible models with and without each of the calendar effects, to check if they are significant or not. To do so we will use the model selected earlier.

After adjusting all the possible models, we observe that the model with the best AIC is the one accounting for both the Holy Week and trading days effects as we can see in table 9. Furthermore, both the absolute t-ratios of  $w_{Ea}$  and  $w_{TD}$  are above 2, indicating that both effects are statistically significant. We observe that  $w_{Ea} = -0.2075$ , which means that having Holy Week fall entirely within one month results in a reduction in exports. On the other hand,  $w_{TD} = 0.0285$  is positive, indicating that increasing the number of working days in a month has a positive effect on exports. Finally, since we see that the AIC improves when taking into account both calendar effects, we will use  $X_t^*$  to fit our models.

#### 4.5.1 Applying Box-Jenkins to $X_t^*$

After obtaining the series  $X_t^*$ , we need to fit a model using the Box-Jenkins methodology that we have used previously. Therefore, we will skip the details of the process.

After transforming the new series into stationary, we obtain the ACF and PACF, which can be seen in figure 5.8.

For the regular part, the ACF follows a decreasing pattern that indicates infinite lags, while in the PACF only 2 lags are significant, suggesting an  $AR(2)$  model.

For the seasonal part, we observe that the seasonal lags in the PACF follow a clear exponential decay, indicating that the infinite lags are found in the PACF. We see that there is only one significant lag in the ACF, since the rest are considered null due to being within the confidence bands. Therefore, we identify an  $MA(1)$  model. We will adjust an  $ARIMA(2, 1, 0)(0, 1, 1)$ .

After estimating our model, we obtain that all the parameters are significant, and an AIC of  $-310.7$  which is better than all the models adjusted before.

Regarding the residuals, we observe that our model satisfies the [homoscedasticity](#), [normality](#), and [independence](#) hypothesis.

Furthermore, as we can see in figure 5.12, our model is causal and invertible. In terms of forecasting capacity our model is [stable](#), and we observe that it gives good [predictions](#) for the last 12 observations of the series.

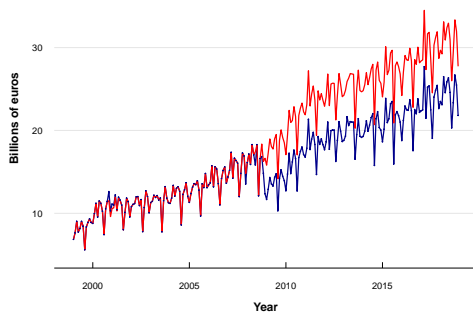
Model	RMSE	MAE	RMSPE	MAPE	CI
ARIMA(2,1,0)(4,1,0) (with $X_t$ )	0.8812	0.6478	0.0382	0.0274	5.8041
ARIMA(2,1,0)(0,1,1) (with $X_t^*$ )	0.9070	0.7268	0.0393	0.0307	5.4066

**Table 5:** Performance metrics comparison

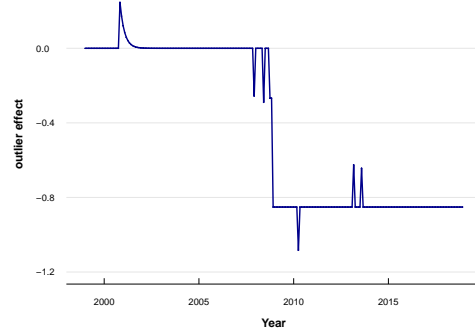
We also observe in table 5 that the model with calendar effects presents slightly worse performance metrics when compared to the previous selected model, but presents less variance in the predictions since the average length of the confidence interval is smaller.

#### 4.6 Outlier treatment

After observing the plots of the residuals of the previously fitted models, we notice the presence of significant outliers, being the most notorious ones the ones corresponding to the 2008 financial crisis, as can be seen in figures 4.9 and 4.10. Therefore, identifying and removing the outliers in our series may be a good idea when fitting our model.



**Fig. 4.9:** Series without the outlier effects (red) vs actual series (blue)



**Fig. 4.10:** Outlier effects

### Outlier interpretation

In our case, we have identified 8 significant outliers present between the years 2000 and 2013 (see table 10). The most evident explanations we found for some of the detected outliers are the following: The two AO and two LS (all with a negative impact) found between the end of 2007 and the end of 2008 are caused due to the crisis of the real estate bubble in Spain[2]. This economic crisis had a negative impact on Spain's export sector, leading to a decline in export volumes and foreign investment. Note that this outliers are the most significant for the series, as there are 2 level shifts that result in a permanent effect throughout the rest of the series.

On the other hand, in April 2010 we observe a negative AO that can be due to the closure of European airspace due to the volcanic ash cloud from Iceland's Eyjafjallajökull volcano [3].

### Model with outlier treatment

After removing the outliers, we proceed to fit new models. After transforming the series to a stationary one, we obtain the ACF and PACF plots shown in figure 5.9. We decided to fit both an ARIMA(2,1,0)(0,1,1) and an ARIMA(0,1,1)(0,1,1) models, and both have the same AIC of -407.

After checking that both models fulfill the assumptions of homoscedasticity, normality, and independence on the residuals in figures 5.14 and 5.15, we can proceed to make predictions using the models. Note that both models are also causal and invertible.

When making predictions on the original series, we will have to add the effect of the outliers that have an impact on the predictions, which in our case are the 2 level shifts of the 2008 crisis since our predictions are made on the series without the effect of exogenous variables.

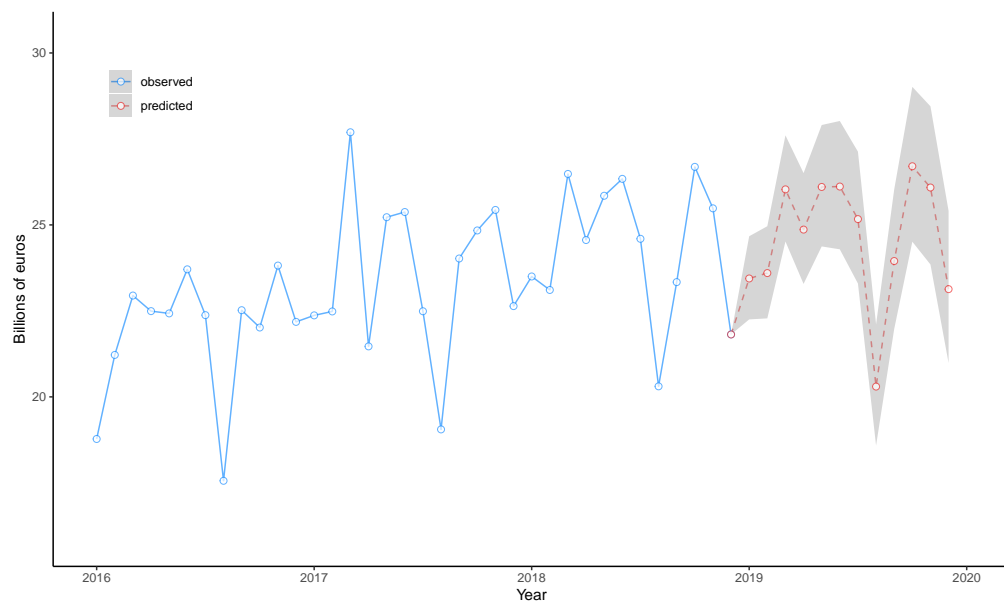
Table 6 shows the metrics of the predictions made on the last 12 samples of the series. We observe that, in general, the second model obtains better criterion values and has a smaller confidence interval. Therefore, we will opt for this second model when making predictions.

Model	RMSE	MAE	RMSPE	MAPE	CI	AIC
ARIMA(2,1,0)(4,1,0) (with $X_t$ )	0.8812	0.6478	0.0382	0.0274	5.8041	-191.7
ARIMA(2,1,0)(0,1,1) (with $X_t^*$ )	0.9070	0.7268	0.0393	0.0307	5.4066	-305.2
ARIMA(2,1,0)(0,1,1) (with OT) <sup>1</sup>	0.9613	0.7716	0.0415	0.0325	3.9286	-407.7
ARIMA(0,1,1)(0,1,1) (with OT)	0.9232	0.7310	0.0399	0.0308	3.6626	-407.8

<sup>1</sup> Outlier Treatment (OT)

**Table 6:** Performance metrics comparison

We observe that the predictions of the second model (figure 4.11) are similar to those obtained previously in section 4.4 but with a smaller confidence interval. In addition, although the metrics related to predictive ability are slightly worse than those of the models without OT, both models with outlier treatment have a considerably better AIC compared to the other two models, and present a much narrower confidence interval. Furthermore, the ARIMA(0,1,1)(0,1,1) model has fewer parameters than the other models, which may be advantageous for interpreting the results.



**Fig. 4.11:** Forecast for the 2019 exportations



## 5 Appendix

The theoretical content presented in this appendix is largely based on the work of [4].

### *Box-Cox Transformation*

$$X_t^{(\lambda)} = \begin{cases} \frac{X_t^\lambda - 1}{\lambda}, & \text{if } \lambda \neq 0 \\ \log(X_t), & \text{if } \lambda = 0 \end{cases}$$

### *Parameter Estimation*

The estimation of parameters is done using Maximum Likelihood (ML) estimation. ML estimators are consistent, asymptotically efficient, and Gaussian.

$$\hat{\Lambda}_{ML} = (\hat{\phi}, \hat{\theta}, \hat{\sigma}_z^2) = \mathbf{argmax}(L(\phi, \theta, \sigma_z^2; X))$$

$$\hat{\Lambda}_{ML} \sim N(\Lambda, I_\Lambda^{-1})$$

### *Parameter Signification*

In order to test whether a parameter is statistically significant, we will perform the following hypothesis test:

$$\begin{cases} H_0 : \phi = 0 \\ H_1 : \phi \neq 0 \end{cases}$$

where

$$\hat{\phi}_i \sim N(\phi_i, \sigma_{\phi_i}) \implies \hat{t} = \frac{\hat{\phi}}{se(\hat{\phi}_i)} \sim t - Student_{T-k}$$

We will conclude that a parameter is significant if  $|\hat{t}| > 2$ . If this value is very close to 2, we may still consider it significant if we observe that removing the parameter results in a decrease in the AIC of the model and the model appears to worsen.

### *Shapiro-Wilk Normality Test*

The Shapiro-Wilk test is a hypothesis test. The null hypothesis is that the sample comes from a normal distribution, and the alternative hypothesis is that it does not. To ascertain normal distribution of residuals, the p-value of the test must exceed 0.05 to accept the null hypothesis that the residuals are normally distributed. If the p-value is less than 0.05, the null hypothesis should be rejected, and it can be concluded that the residuals do not follow a normal distribution.

### *Ljung-Box Test*

For a lag  $k$  test the joint hypothesis that the first  $k$  autocorrelation of the residuals are jointly zero:

$$H_0 : \rho_z(1) = \rho_z(2) = \dots = \rho_z(k)$$

To verify that the residuals are independent, we should see that all calculated p-values are above 0.05.

### ***Infinite coefficients models***

Under certain conditions of stationarity, the ARMA(p, q) models can be expressed as an  $AR(\infty)$  or  $MA(\infty)$  model.

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

Expression as  $AR(\infty)$ :

$$\frac{(1 - \phi_1 B - \dots - \phi_p B^p)}{(1 + \theta_1 B + \dots + \theta_q B^q)}X_t = (1 - \pi_1 B - \pi_2 B^2 - \dots)X_t = Z_t$$

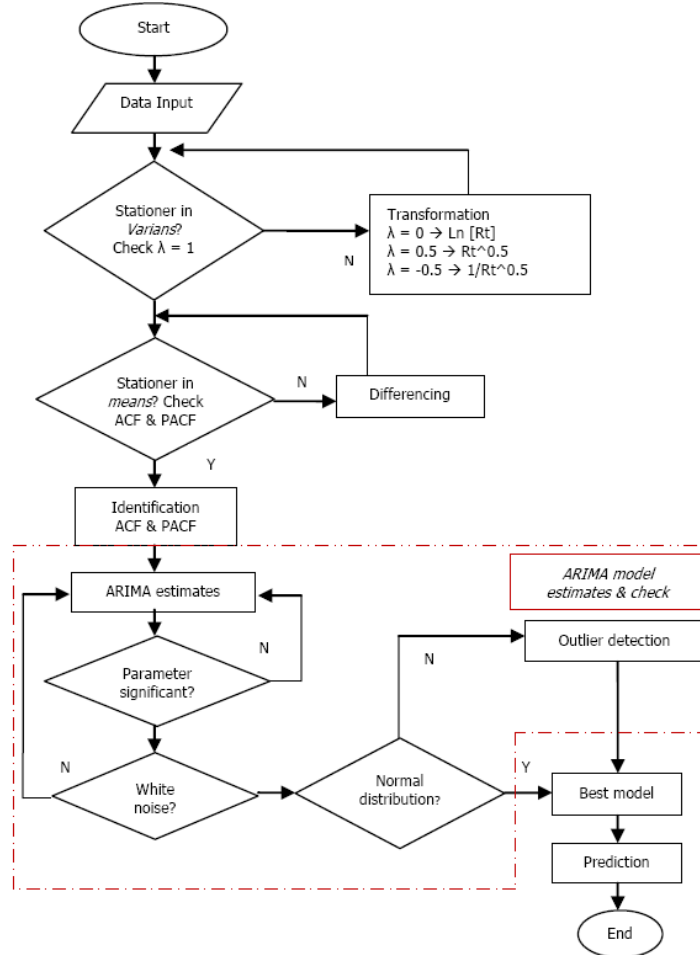
Expression as  $MA(\infty)$ :

$$X_t = \frac{(1 + \theta_1 B + \dots + \theta_q B^q)}{(1 - \phi_1 B - \dots - \phi_p B^p)}Z_t = (1 + \psi_1 B + \psi_2 B^2 + \dots)Z_t$$

### ***Outlier Types***

- Additive Outlier (AO): It affects only one period  
Transfer function:  $X_t = \mathbf{1}_{t=T_0}(t)$
- Transitory Change (TC): It affects on one period and its effect decreases in the next periods  
Transfer function:  $X_t = \delta^{t-T_0} \mathbf{1}_{t \geq T_0}(t)$ , with  $\delta = 0.7$
- Level Shift (LS): It affects on one period and its effect remains in the next periods  
Transfer function:  $X_t = \mathbf{1}_{t \geq T_0}(t)$

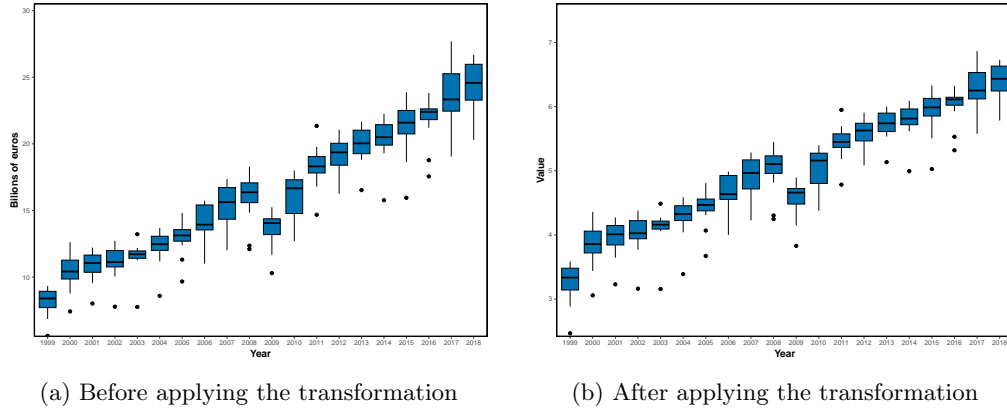
The majority of the plots in this study were created using the *ggplot* package [5].



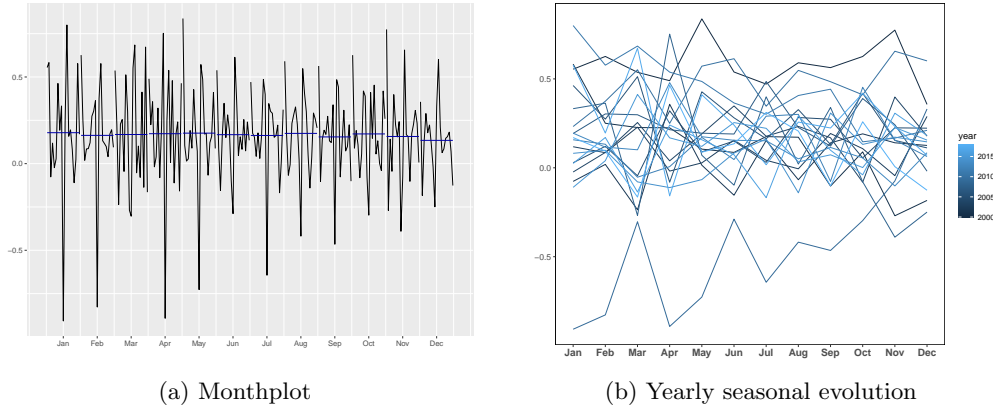
**Fig. 5.1:** Box-Jenkins Steps for time series analysis

Variance comparison	
Series	Variance
$W'_t$	0.0755
$(1 - B)W'_t$	0.0657
$(1 - B)^2W'_t$	0.2075

**Table 7:** Comparison of variance



**Fig. 5.2:** Boxplot of the series



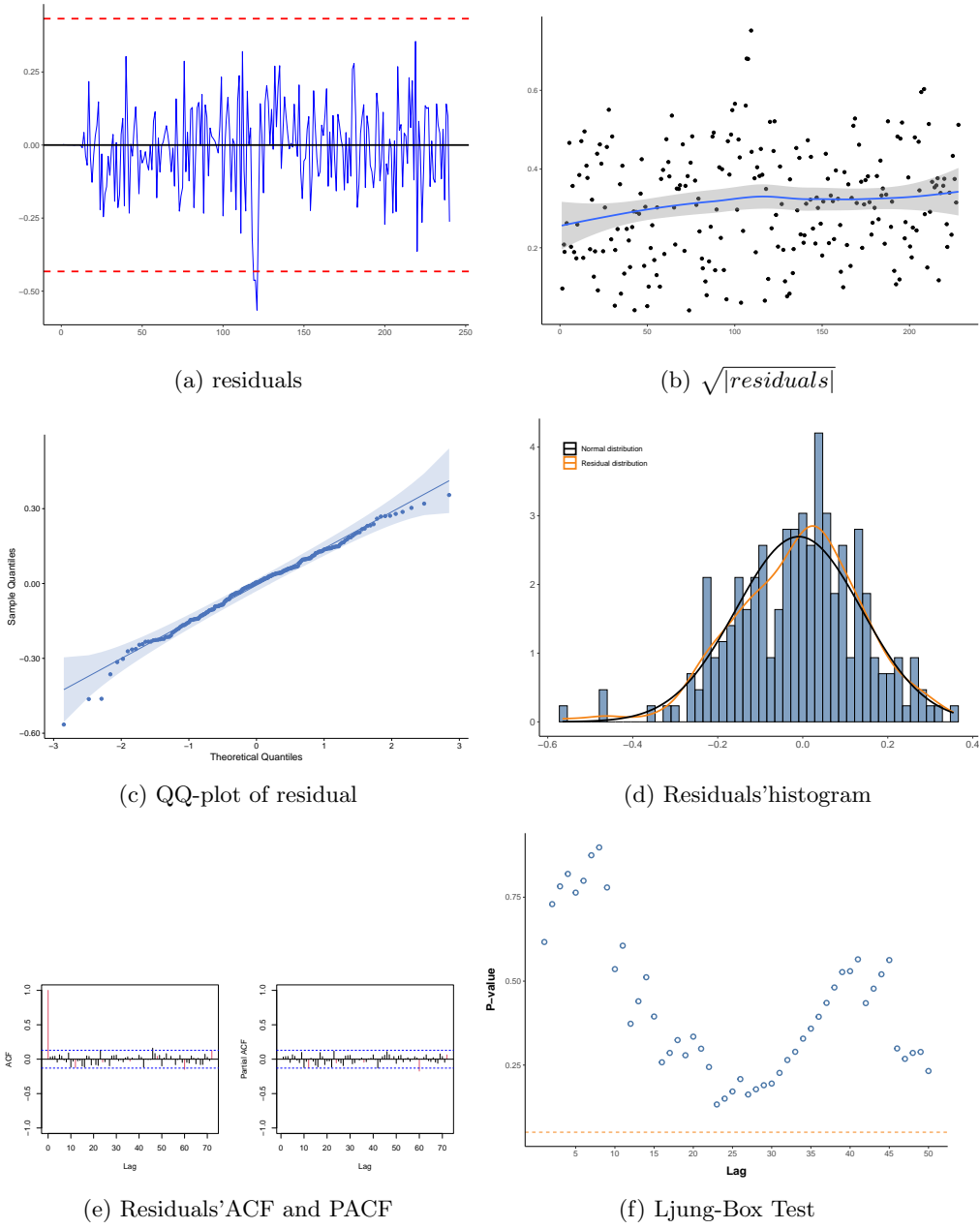
**Fig. 5.3:** Seasonality plots of the differentiated series

$m$ coefficients	-0.471	-0.176	-0.757	-0.193	0.028
$m_{12}$ coefficients	-0.455	-0.173	-0.762	-0.201	0.028
Difference ( $m - m_{12}$ )	0.015	0.003	0.005	0.008	0
T-ratios ( $m_1$ )	-7.15	-2.64	-14.23	-6.75	14.22
T-ratios ( $m_{12}$ )	-6.77	-2.55	-14.01	-6.75	13.72

**Table 8:** Comparison of coefficients for ARIMA(2,1,0)(0,1,6) with  $X_t^*$

	without Ea	with Ea
without TD	-191.7	-226.9
with TD	-266.9	-305.2

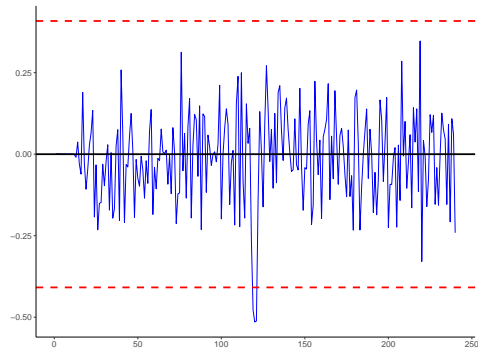
**Table 9:** AIC obtained for the models with the different calendar effects



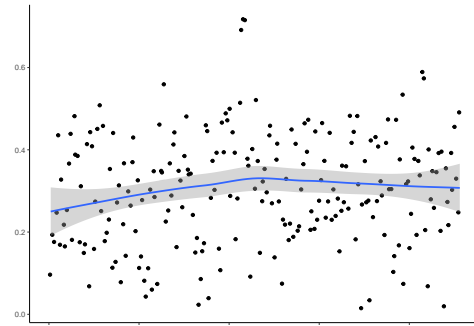
**Fig. 5.4:** Residual analysis for ARIMA(2,1,0)(4,1,0) with outlier treatment

Obs	type_detected	W_coef	ABS.L.Ratio	Fecha
23	TC	0.25	3.09	Nov 2000
108	AO	-0.26	3.34	Dic 2007
114	AO	-0.29	3.62	Jun 2008
118	LS	-0.27	3.23	Oct 2008
120	LS	-0.58	6.73	Dic 2008
136	AO	-0.23	3.15	Abr 2010
171	AO	0.23	3.13	Mar 2013
176	AO	0.21	2.93	Ago 2013

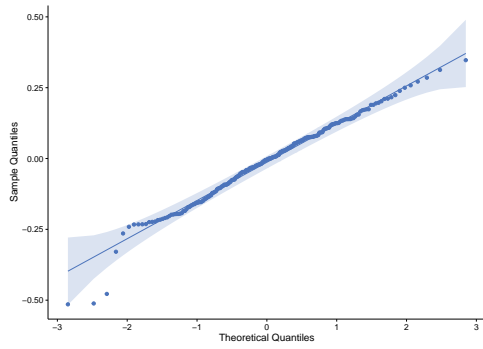
**Table 10:** Detected outliers



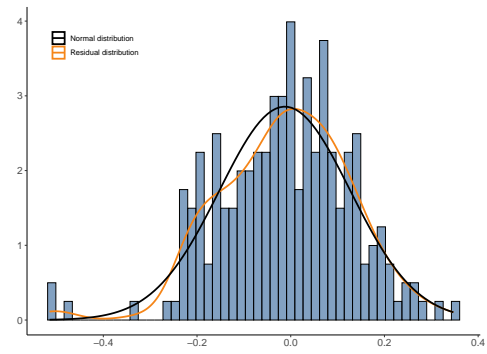
(a) residuals



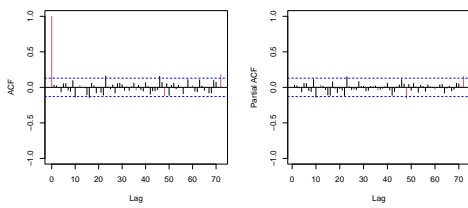
(b)  $\sqrt{|residuals|}$



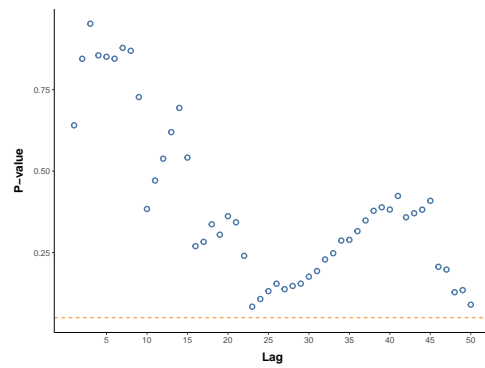
(c) QQ-plot of residual



(d) Residuals' histogram

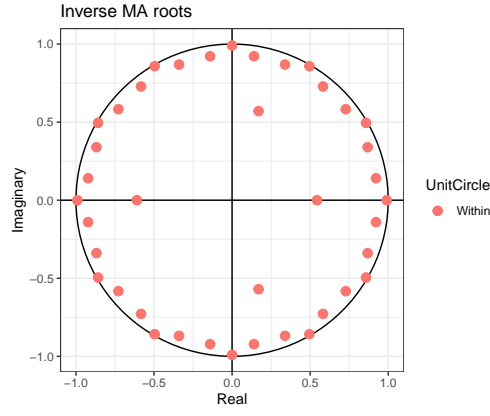


(e) Residuals' ACF and PACF

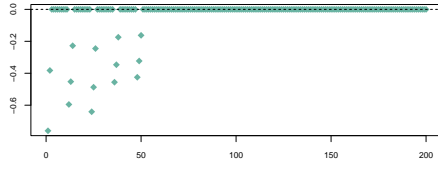


(f) Ljung-Box Test

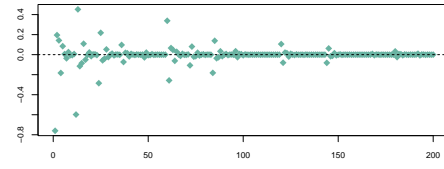
**Fig. 5.5:** Residual analysis for ARIMA(2,1,0)(0,1,6)



(a) Inverse AR roots

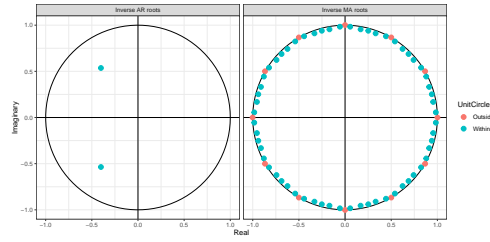


(b)  $\pi$ -weights

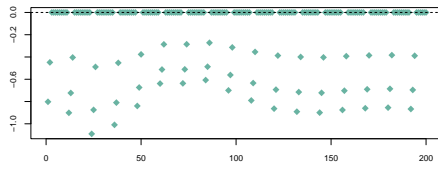


(c)  $\psi$ -weights

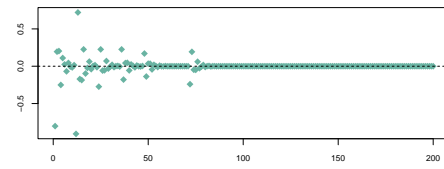
**Fig. 5.6:** Cheking invertibility and causality properties for ARIMA(2,1,0)(4,1,0)



(a) Inverse AR and MA roots

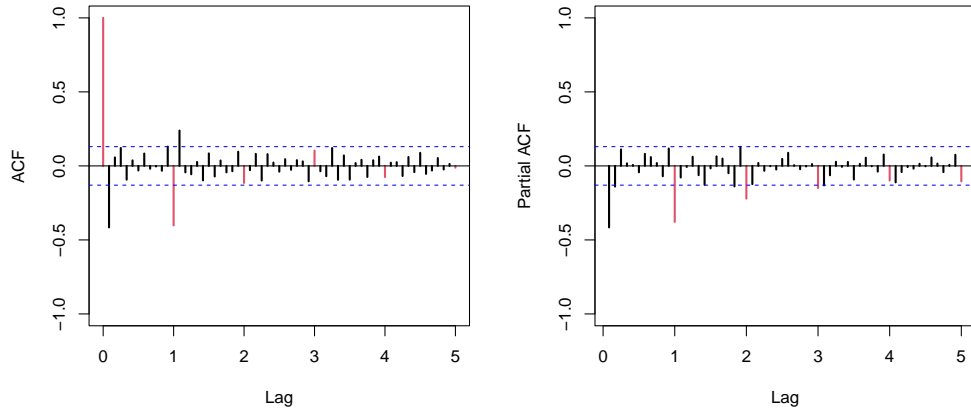


(b)  $\pi$ -weights

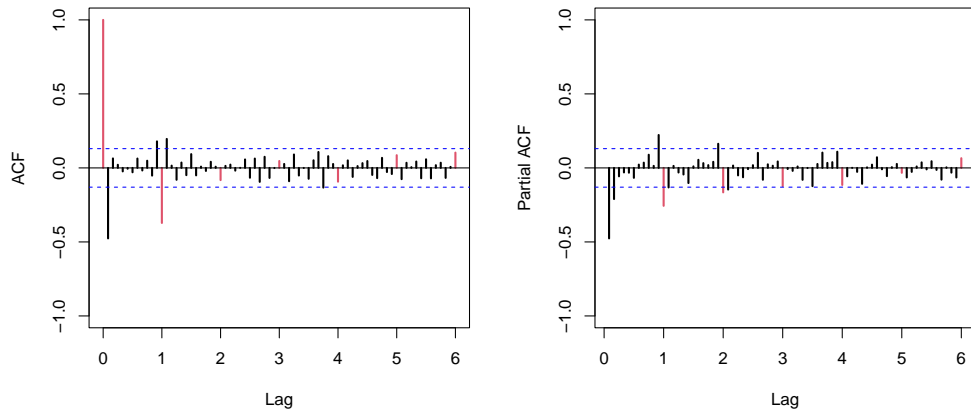


(c)  $\psi$ -weights

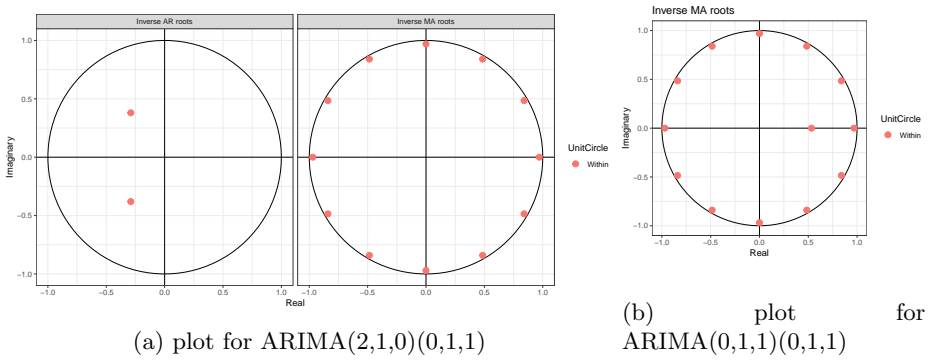
**Fig. 5.7:** Cheking invertibility and causality properties for ARIMA(2,1,0)(0,1,6)



**Fig. 5.8:** ACF of the stationary series without calendar effects

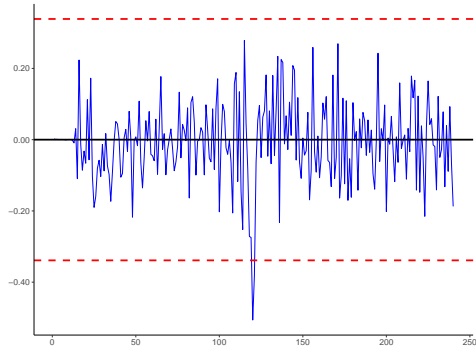


**Fig. 5.9:** ACF and PACF of the stationary series without outliers

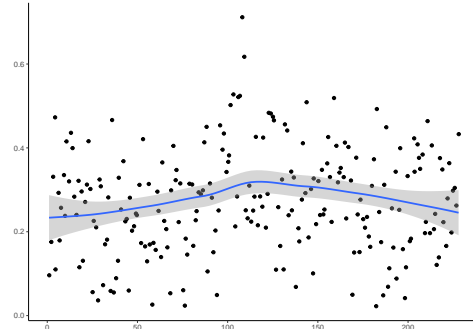


**Fig. 5.10:** Inverse AR and MA roots for the models with outlier treatment

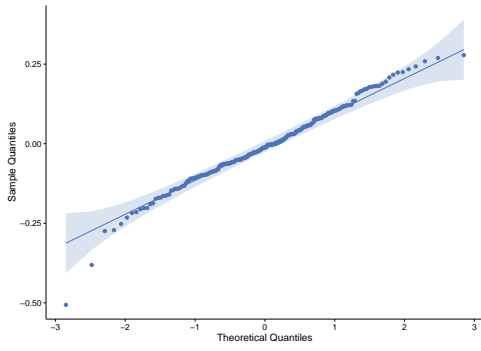




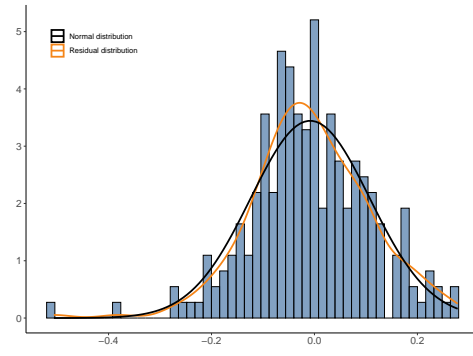
(a) residuals



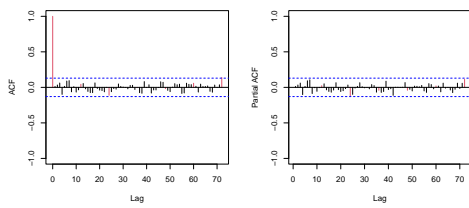
(b)  $\sqrt{|residuals|}$



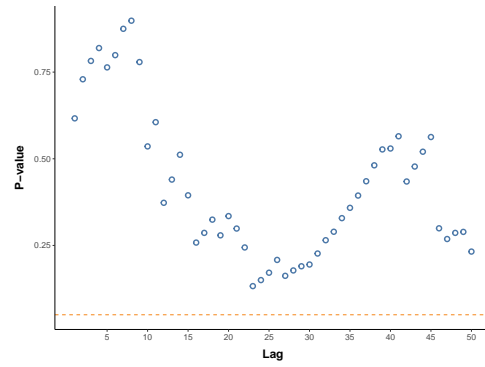
(c) QQ-plot of residual



(d) Residuals' histogram

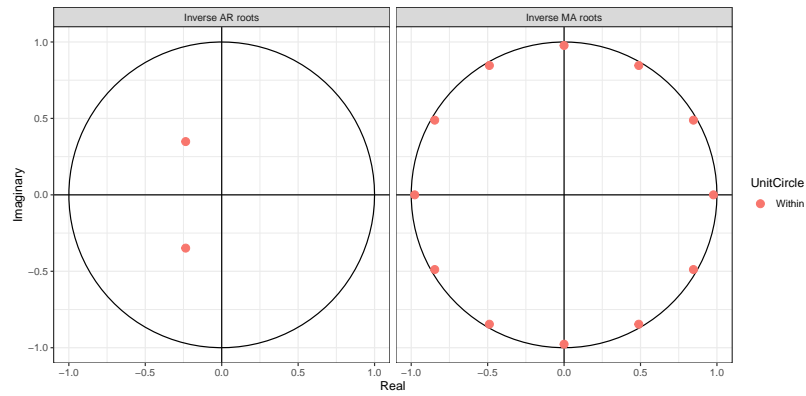


(e) Residuals' ACF and PACF

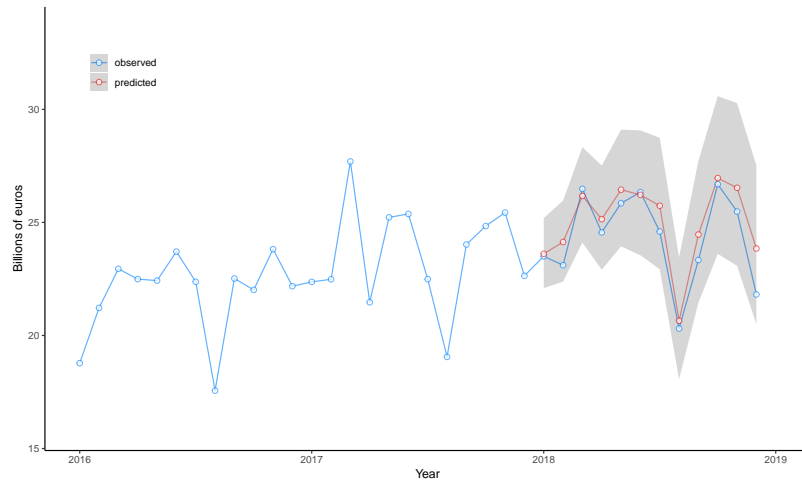


(f) Ljung-Box Test

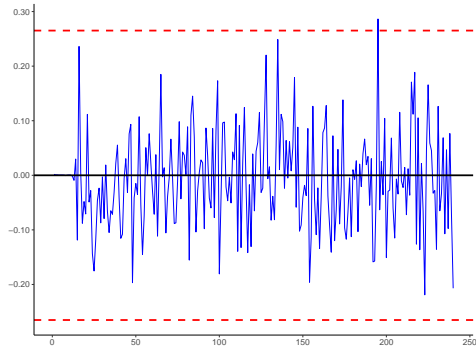
**Fig. 5.11:** Residual analysis for  $X_t^*$



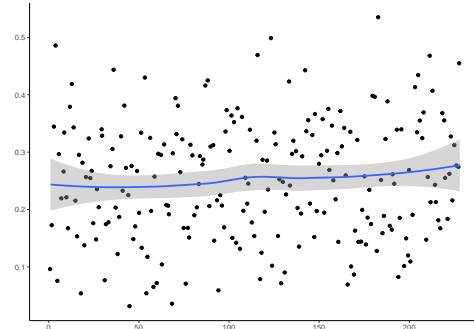
**Fig. 5.12:** Inverse AR and MA roots for the model with calendar effects



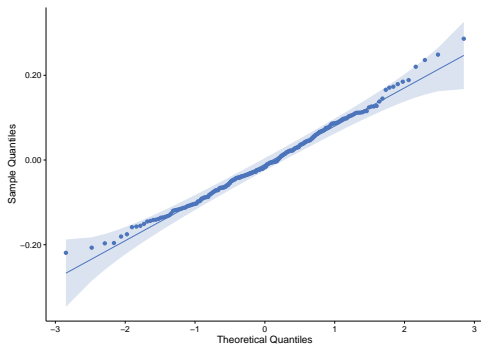
**Fig. 5.13:** Predictions for  $\text{ARIMA}(2,1,0)(0,1,1)$  with  $X_t^*$



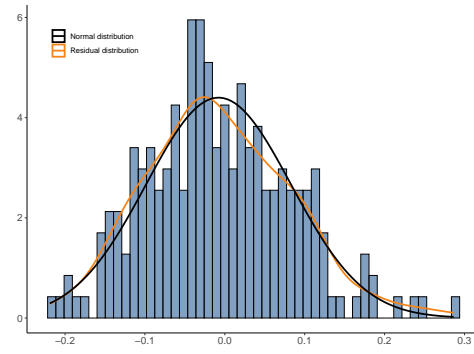
(a) residuals



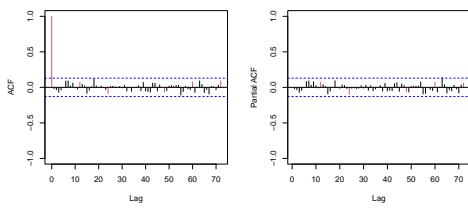
(b)  $\sqrt{|residuals|}$



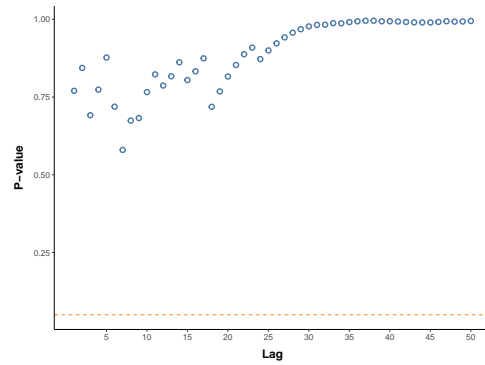
(c) QQ-plot of residual



(d) Residuals' histogram

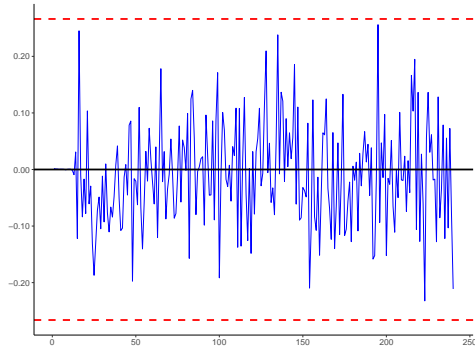


(e) Residuals' ACF and PACF

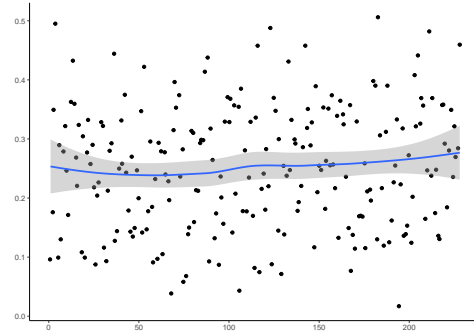


(f) Ljung-Box Test

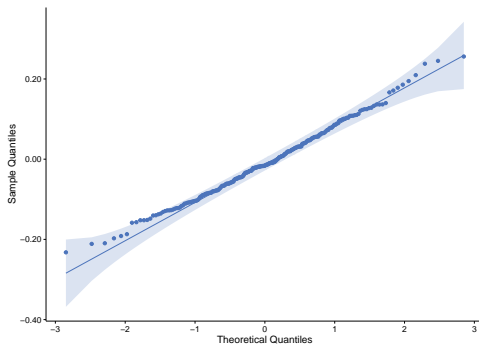
**Fig. 5.14:** Residual analysis for ARIMA(2,1,0)(0,1,1) with outlier treatment



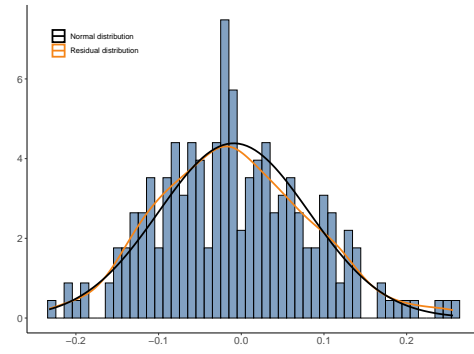
(a) residuals



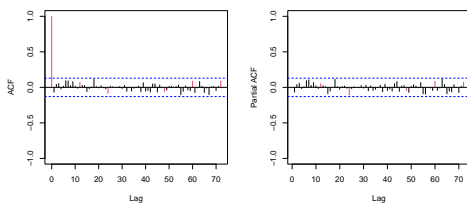
(b)  $\sqrt{|residuals|}$



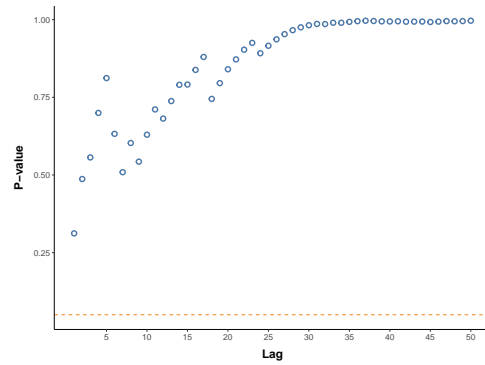
(c) QQ-plot of residual



(d) Residuals' histogram



(e) Residuals' ACF and PACF



(f) Ljung-Box Test

**Fig. 5.15:** Residual analysis for ARIMA(0,1,1)(0,1,1) with outlier treatment

## References

- [1] Ministry of Industry, Commerce and Tourism. Time series of total exportations in spain. Retrieved from <https://sedeaplicaciones.minetur.gob.es/Badase/BadasiUI/lstSeriesInformesPostBack.aspx>
- [2] E. Mundo. El crash del 2008. URL <https://www.elmundo.es/especiales/2008/10/economia/crisis2008/queestapasando/index.html>
- [3] ABC. Eyjafjallajökull, el volcán de islandia que sembró el caos en 2010 con el cierre del espacio aéreo de media europa. URL <https://www.abc.es/sociedad/eyjafjallajokull-volcan-islandia-sembro-caos-2010-cierre-20220803195547-nt.html>
- [4] J. Sanchez. Data analysis (2023). Class lecture, Polytechnic University of Catalonia.
- [5] ggplot2, function reference (2018). URL <https://ggplot2.tidyverse.org/reference/>. Last accessed 1 March 2023