



Linear Time Series Project

AUTHORS

CIANFARANI Ana-Sofia
BENABDESADOK Nayel

May 2025

Contents

1	Data	1
1.1	What does the chosen series represent	1
1.2	Making the series stationary	1
1.3	Graphic Representation	2
2	ARMA Models	2
2.1	Choosing the Appropriate ARMA(p, q) Model	2
2.2	Expression of X_t	3
3	Prediction	4
3.1	Forecasting X_{T+1} and X_{T+2}	4
3.2	Main hypotheses	5
3.3	Graphic Representation	5
3.4	Improving the prediction of X_{T+1}	6
A	Study of the original series	7
B	Seasonality	8
C	Stationarity tests	10
D	Choose p and q	11
E	Outlier Types and Their Treatment in Time Series Model	12
F	Calculus	14
G	Code	15

1 Data

1.1 What does the chosen series represent

The chosen series is a French Industrial Production Index (IPI) Series from INSEE¹. It represents construction as a sector (Section F of the NAF classification). Its base year is 2021. Data is available on a monthly basis from January 1990 to February 2025, thus containing 422 observations. The series is corrected from seasonal variations and working days (CVS-CJO). The goal is to facilitate conjunctural analysis of industrial production over time (see Appendix B for further detail).

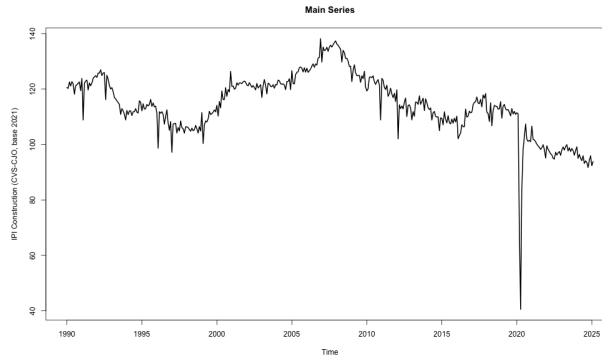


Figure 1: Plot of the Original Series (as a function of time)

The series evolves in an irregular manner, with three main phases. From 1990 to 2000, it decreases, with the level of production reducing by about 20 basis points over this period. This decrease reflects a period of economic difficulties for the construction sector as a whole, which struggled against poor economic conditions such as rising unemployment, high real interest rates, and public debt reduction to prepare for the introduction of the euro. Moreover, a wave of office construction in the early 1990s set off a real estate crisis due to the saturation of the market[INSEE, 2000].

Then, from 2000 to 2008, a net increase of industrial production can be seen, gaining about 30 basis points, thanks to generally favorable economic conditions. From 2008 onwards, the series has decreased regularly.

A slight increase was seen from early 2017, however these gains were wiped out in early 2020, and production has continued to decrease since, with production now at its lowest recorded level (excluding the Covid-era restrictions period) : it has lost 40 basis points since 2008. March & April 2020 stand out as remarkable values in our series, as construction was halved in the span of a month : this was due to the partial shutdown of economic activity linked to the pandemic. This outlier value particularly stands out; however, one can also notice many smaller outliers throughout the series, which will require specific treatment in our analysis so as to ensure they do not overly influence our model.

1.2 Making the series stationary

The series seems to have a deterministic trend and suffer from heteroscedasticity, which is confirmed by analysis of the ACF and the decomposition of the series into trend-cycle, seasonal variations and residuals. Classic stationarity tests also confirm this intuition. The detailed results are available in Appendixes A and C.

As the series has already been corrected, there should not be any seasonality : this is confirmed by a number of graphic tests presented in Appendix B.

To correct for non-stationarity, we take the first difference of the series. It now appears stationary, and this is confirmed by all three classic stationary tests (see Appendix C). We conclude that the series is of integration order 1, $I(1)$.

¹Full series available at : <https://www.insee.fr/fr/statistiques/serie/010767635>.

1.3 Graphic Representation

These graphs represent the original time series X_t and its stationary version $Y_t = X_t - X_{t-1}$.

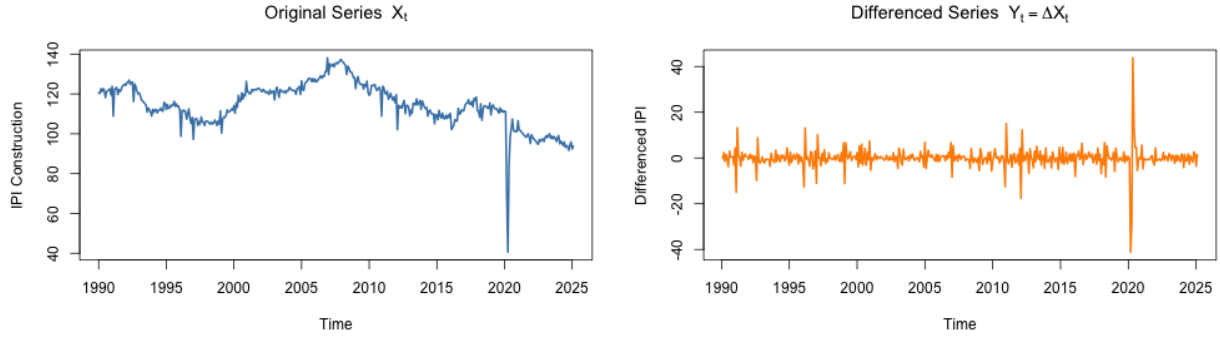


Figure 2: Plot of the Original Series (left) and the Differenced Series (right)

2 ARMA Models

2.1 Choosing the Appropriate ARMA(p, q) Model

After transforming the original series into a stationary one through first-order differencing (i.e., $Y_t = X_t - X_{t-1}$), we aim to identify the optimal AR and MA orders p and q for the ARMA(p, q) model. This step relies on examining the autocorrelation function (ACF) and the partial autocorrelation function (PACF).

From Figure 3, we can assess the statistical significance of the autocorrelation and partial autocorrelation coefficients. Values lying outside the bounds $\pm 1.96/\sqrt{n}$ exceed the 95% confidence interval under the null hypothesis of zero correlation, and are thus considered statistically significant.

- The **ACF** is used to identify the MA order q : for an MA(q) process, autocorrelations drop to zero after lag q .
- The **PACF** is used to identify the AR order p : for an AR(p) process, partial autocorrelations drop to zero after lag p .

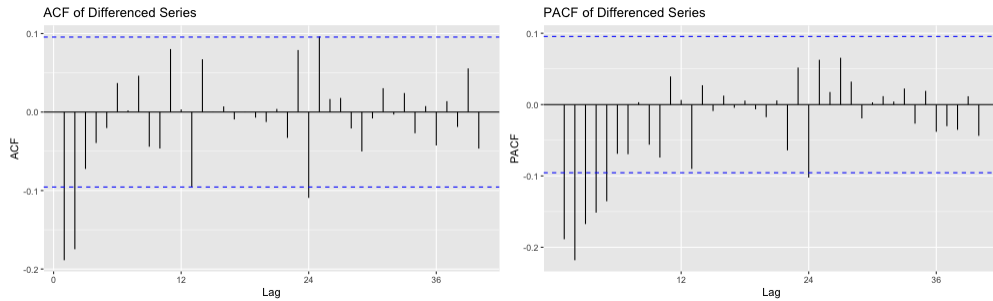


Figure 3: ACF (left) and PACF(right) of the differenced series Y_t .

Note: The dashed blue line represents the 95% significance bounds.

Based on Figure 3, we identify a maximum MA order $q_{\max} = 2$ (cutoff after lag 2) and a maximum AR order $p_{\max} = 5$ (cutoff after lag 5). Some higher-order lags exceed the confidence bounds, but they occur after long sequences of non-significant lags and are therefore disregarded according to standard practice².

²Some lags beyond p_{\max} and q_{\max} exceed the confidence bounds, but we ignore them since they appear after many insignificant lags.

We now aim to select an ARMA(p, q) model such that $p \leq 5$ and $q \leq 2$. The selection process follows these steps:

Significance of Coefficients We first assess the statistical significance of the estimated AR coefficients ϕ_i and MA coefficients θ_j . This involves testing the null hypothesis H_0 : “the coefficient equals zero” against the alternative hypothesis H_1 : “the coefficient is significantly different from zero.”

If any coefficient in a given ARMA(p, q) model is not statistically significant (typically, if its p-value exceeds 5%), this may indicate model overfitting. In such cases, a more parsimonious model—excluding the non-significant terms—is generally preferred, both for interpretability and model performance.

After this step, the following candidate models remain: ARMA(0,1), ARMA(0,2), ARMA(1,0), ARMA(1,1), ARMA(2,0), ARMA(3,0), ARMA(4,0), and ARMA(5,0).

Ljung-Box Test The Ljung-Box test is used to assess whether the residuals from a model are free of autocorrelation. The null hypothesis H_0 states that “the residuals are not autocorrelated” (i.e., they are white noise), whereas the alternative hypothesis H_1 states that “the residuals are autocorrelated.” We reject any model for which the p-value of the Ljung-Box test is below 5%, as this indicates statistically significant evidence of autocorrelation in the residuals, thereby violating the white noise assumption. In other words, we retain only those models for which the Ljung-Box test fails to reject the null hypothesis H_0 : the residuals are independently distributed. At this stage, only the following models remain viable: ARMA(0,2), ARMA(1,1), ARMA(4,0), and ARMA(5,0). These models satisfy both the requirement of coefficient significance and the Ljung-Box diagnostic for uncorrelated residuals.

Model Selection via AIC and BIC Minimization To finalize our choice, we compare the remaining models using the Akaike Information Criterion (AIC) and the Bayesian Information Criterion (BIC), both of which reward goodness of fit while penalizing model complexity. The ARMA(1,1) model exhibits the lowest values for both AIC and BIC, indicating the best trade-off between accuracy and simplicity. We therefore select ARMA(1,1), with $p = 1$ and $q = 1$.

A comprehensive summary Table 3 in the Appendix reports the results of coefficient significance tests, Ljung-Box statistics, and the corresponding AIC and BIC values for each ARMA model fitted to the series Y_t .

2.2 Expression of X_t

As concluded in the previous section, the process X_t can be modeled as an ARIMA(1,1,1). Here, $\epsilon_t \sim \text{WN}(0, \sigma^2)$. The standard formulation of such a process is given by:

$$\phi(B)(1 - B)X_t = \theta(B)\epsilon_t, \quad (1)$$

where B denotes the backshift operator ($BX_t = X_{t-1}$). The autoregressive and moving average polynomials are defined as:

$$\begin{aligned} \phi(B) &= 1 - 0.4764B, \\ \theta(B) &= 1 + 0.8387B. \end{aligned}$$

However, as highlighted in Section 1.1, several outliers are present in the series X_t and must be explicitly accounted for. Table 4 lists the identified outliers along with their type, estimated magnitude, and associated t -statistic.

A total of 10 outliers were detected using our R procedure, as shown in Figure 10. However, the last one corresponds to extreme values observed in 2020, which likely reflect the extraordinary shock caused by the COVID-19 pandemic. As these may not be representative of the usual dynamics of the series, we restrict our analysis to the first nine outliers. In particular, the last outlier does not appear significant when visually inspecting Figure 10.

Using the specification framework provided in Appendix E, we can now incorporate these shocks directly into the ARIMA(1,1,1) model. Starting from equation (1), and applying the appropriate structure for each outlier type as listed in Table 4, while adjusting the estimated coefficients ϕ_1 and θ_1 , the resulting model becomes:

$$Y_t - 0.1781 \cdot Y_{t-1} = \epsilon_t - 0.7285 \cdot \epsilon_{t-1} - \sum_{i=1}^7 \omega_i \cdot \mathbf{1}_{\{t=t_i\}} - \omega_8 \cdot \delta^{t-t_{\text{Feb.2016}}} \cdot \mathbf{1}_{\{t \geq t_{\text{Feb.2016}}\}} - \omega_9 \cdot \delta^{t-t_{\text{Mar.2020}}} \cdot \mathbf{1}_{\{t \geq t_{\text{Mar.2020}}\}}. \quad (2)$$

X_t is therefore :

$$X_t = 1.1781 \cdot X_{t-1} - 0.1781 \cdot X_{t-2} + \epsilon_t - 0.7285 \cdot \epsilon_{t-1} - \sum_{i=1}^7 \omega_i \cdot \mathbf{1}_{\{t=t_i\}} - \omega_8 \cdot \delta^{t-t_{\text{Feb.2016}}} \cdot \mathbf{1}_{\{t \geq t_{\text{Feb.2016}}\}} - \omega_9 \cdot \delta^{t-t_{\text{Mar.2020}}} \cdot \mathbf{1}_{\{t \geq t_{\text{Mar.2020}}\}}.$$

3 Prediction

3.1 Forecasting X_{T+1} and X_{T+2}

We aim to forecast the values of X_{T+1} and X_{T+2} , where $T = 422$ corresponds to the final observed time point.

Note that regardless of the value chosen for the damping parameter δ associated with transient change (TC) outliers, we have $\delta^{t-t_0} \rightarrow 0$ as $t \rightarrow \infty$. Therefore, the effect of such outliers becomes negligible in the long run ($t_0 = 314$ or $t_0 = 363$).

Under this assumption, for $t \geq T$, the dynamics of the series Y_t effectively reduce to a standard ARMA(1,1) model, which can be written as:

$$Y_t - \phi_1 \cdot Y_{t-1} = \epsilon_t - \theta_1 \cdot \epsilon_{t-1} \text{ where the estimated coefficients are } \phi_1 = 0.1781, \theta_1 = -0.7285.$$

According to the theory (Chapter 2), for any forecast horizon $h > 0$, $\mathbb{E}[\epsilon_{T+h} \mid Y_T, Y_{T-1}, \dots, Y_0] = 0$ as the residuals ϵ_{T+h} are orthogonal to the $Y_t, t \leq T$.

Using the equation of Y_t and recursively substituting Y_{T+1} into the expression for Y_{T+2} , we obtain the following system:

$$\begin{cases} Y_{T+1} = \phi_1 \cdot Y_T + \epsilon_{T+1} - \theta_1 \cdot \epsilon_T \\ Y_{T+2} = \phi_1^2 \cdot Y_T - \phi_1 \theta_1 \cdot \epsilon_T + (\phi_1 - \theta_1) \cdot \epsilon_{T+1} + \epsilon_{T+2} \end{cases}$$

From linearity of the conditional expectation and using the result above :

$$\begin{cases} \hat{Y}_{T+1|T} = \phi_1 \cdot \mathbb{E}[Y_T \mid Y_t, t \leq T] + \mathbb{E}[\epsilon_{T+1} \mid Y_t, t \leq T] - \theta_1 \cdot \mathbb{E}[\epsilon_T \mid Y_t, t \leq T] \\ \hat{Y}_{T+2|T} = \phi_1^2 \cdot \mathbb{E}[Y_T \mid Y_t, t \leq T] - \phi_1 \theta_1 \cdot \mathbb{E}[\epsilon_T \mid Y_t, t \leq T] + (\phi_1 - \theta_1) \cdot \mathbb{E}[\epsilon_{T+1} \mid Y_t, t \leq T] + \mathbb{E}[\epsilon_{T+2} \mid Y_t, t \leq T] \end{cases}$$

We obtain the following point forecasts based on the ARMA(1,1) model:

$$\begin{cases} \hat{Y}_{T+1|T} = \phi_1 \cdot Y_T - \theta_1 \cdot \epsilon_T, \\ \hat{Y}_{T+2|T} = \phi_1 \cdot (\phi_1 Y_T - \theta_1 \cdot \epsilon_T). \end{cases}$$

Our goal is now to derive the confidence region at level $1 - \alpha$ for the future values (X_{T+1}, X_{T+2}) . We are particularly interested in the forecast error vector, denoted $X - \hat{X} \in \mathbb{R}^2$. By reintroducing the original series X_t using the relation $Y_t = X_t - X_{t-1}$, and recalling that:

$$\begin{cases} \hat{Y}_{T+1|T} = \hat{X}_{T+1|T} - X_T, \\ \hat{Y}_{T+2|T} = \hat{X}_{T+2|T} - \hat{X}_{T+1|T}, \end{cases}$$

We calculate : $\begin{cases} X_{T+1} - \hat{X}_{T+1|T} = Y_{T+1} - \hat{Y}_{T+1|T} = \epsilon_{T+1}, \\ X_{T+2} - \hat{X}_{T+2|T} = Y_{T+2} - \hat{Y}_{T+2|T} + Y_{T+1} - \hat{Y}_{T+1|T} = (1 + \phi_1 - \theta_1) \cdot \epsilon_{T+1} + \epsilon_{T+2}. \end{cases}$

We thus obtain the forecast error vector:

$$X - \hat{X} = \begin{pmatrix} X_{T+1} - \hat{X}_{T+1|T} \\ X_{T+2} - \hat{X}_{T+2|T} \end{pmatrix} = \begin{pmatrix} \epsilon_{T+1} \\ (1 + \phi_1 - \theta_1) \cdot \epsilon_{T+1} + \epsilon_{T+2} \end{pmatrix}. \quad (3)$$

Assuming the normality of the residuals : $\epsilon_t \sim \mathcal{N}(0, \sigma^2)$, so $X - \hat{X} \sim \mathcal{N}(0, \Sigma)$ where Σ is the variance-covariance matrix of the forecast errors $\in \mathcal{M}_2(\mathbb{R})$. We assume that Σ is invertible.

The computations leading to the expression of Σ are provided in appendix F. We ultimately obtain:

$$\Sigma = \begin{pmatrix} 1 & 1 + \phi_1 - \theta_1 \\ 1 + \phi_1 - \theta_1 & 1 + (1 + \phi_1 - \theta_1)^2 \end{pmatrix} \sigma^2. \quad (4)$$

We are then interested in the test statistic $(X - \hat{X})^T \Sigma^{-1} (X - \hat{X}) \sim \chi^2(2)$. For a risk level α , and thus a confidence level $1 - \alpha$: $\mathbb{P} \left((X - \hat{X})^T \Sigma^{-1} (X - \hat{X}) \leq \chi_{1-\alpha}^2(2) \right) = 1 - \alpha$ where $\chi_{1-\alpha}^2(2)$ is the quantile of order $1 - \alpha$ of the chi-squared distribution with 2 degrees of freedom.

The $(1-\alpha)$ -level confidence region for the bivariate forecast vector X is given by (see [Banks and Fienberg, 2003]):

$$R_{1-\alpha}(X) = \left\{ x \in \mathbb{R}^2 \mid (x - \hat{X})^T \Sigma^{-1} (x - \hat{X}) \leq \chi_{1-\alpha}^2(2) \right\}$$

More concretely, we can derive the individual 95% confidence intervals for the two forecasts as follows:

$$\begin{cases} IC_{95\%}(X_{T+1}) = [\hat{X}_{T+1|T} - 1.96 \cdot \hat{\sigma}, \hat{X}_{T+1|T} + 1.96 \cdot \hat{\sigma}] \\ IC_{95\%}(X_{T+2}) = [\hat{X}_{T+2|T} - 1.96 \cdot \hat{\sigma} \cdot \sqrt{1 + (1 + \hat{\phi}_1 - \hat{\theta}_1)^2}, \hat{X}_{T+2|T} + 1.96 \cdot \hat{\sigma} \cdot \sqrt{1 + (1 + \hat{\phi}_1 - \hat{\theta}_1)^2}] \end{cases} \quad (5)$$

3.2 Main hypotheses

The main hypotheses used in estimating the confidence interval are:

- The consistency of all estimators, meaning our model is identifiable.
- Zero mean for future innovations (for any $h > 0$, $\mathbb{E}[\epsilon_{T+h} \mid Y_T, Y_{T-1}, \dots, Y_0] = 0$).
- The residuals are independent and identically distributed across a Gaussian distribution : $\epsilon_t \sim \mathcal{N}(0, \sigma^2)$. This also means the ϵ_t are strong white noise and homoskedastic.
- The variance-covariance matrix Σ is invertible and positive-definite ($\sigma_\epsilon^2 > 0$).
- For any damping parameter δ associated with transient change (TC) outliers, we have $\delta^{t-t_0} \rightarrow 0$ as $t \rightarrow \infty$.

3.3 Graphic Representation

The predicted variables for March and April 2025 are 93.8607 and 93.8662 respectively. Their respective confidence intervals are [85.7689; 101.9524] and [84.9941; 102.7382]. Figure 4 shows these predictions and confidence intervals in both a line plot (including observations over the last three years) and a joint confidence interval prediction ellipse.

The confidence intervals on our line plot seem quite large : our maximal and minimal values are, respectively, larger and smaller than any recent observations. The elongated shape of the confidence ellipse and its orientation show relatively strong positive correlation between our two predictions, which makes sense.

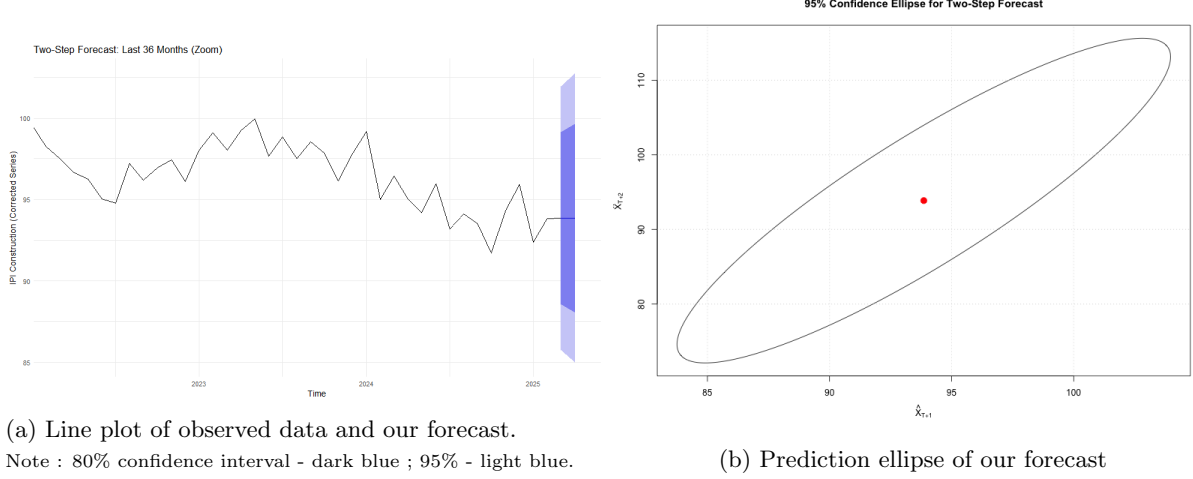


Figure 4: Forecasts and prediction intervals for March and April 2025

3.4 Improving the prediction of X_{T+1}

In this section, let Z_t a stationary time series available from $t = 1$ to T . We assume that Z_{T+1} is available faster than X_{T+1} . We will discuss how this information may allow us to improve the prediction X_{T+1} .

Saying Z_T is useful to predict X_T is equivalent to saying it causes X_T in the Granger sense [Granger, 1969]. Here, as we are specifically asking if Z_{T+1} improves the prediction of X_{T+1} , we are checking if there is instantaneous causality. This is true if and only if:

$$\hat{X}_{T+1}|\{X_u, Z_u, u \leq T\} \cup \{Z_{T+1}\} \neq \hat{X}_{T+1}|\{X_u, Z_u, u \leq T\}$$

This can also be represented in vector form.

Assume a bivariate vector autoregressive model of order p for the stationary time series (X_t, Z_t) :

$$\begin{pmatrix} X_{T+1} \\ Z_{T+1} \end{pmatrix} = \sum_{i=1}^p A_i \begin{pmatrix} X_{T+1-i} \\ Z_{T+1-i} \end{pmatrix} + \begin{pmatrix} \epsilon_{1,T+1} \\ \epsilon_{2,T+1} \end{pmatrix}, \quad \begin{pmatrix} \epsilon_{1,T+1} \\ \epsilon_{2,T+1} \end{pmatrix} \sim \mathcal{N}(\mathbf{0}, \Sigma),$$

where $\Sigma = \mathbb{E}[\epsilon_{T+1}\epsilon'_{T+1}]$ is the variance-covariance matrix of the innovation vector ϵ_{T+1} . We assume Σ is positive definite.

We have instantaneous causality from Z_{T+1} to X_{T+1} if and only if the covariance $\sigma_{12} \neq 0$, where $\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{pmatrix}$. For $T+1$, that means $\text{Cov}(\epsilon_{1,T+1}, \epsilon_{2,T+1}) \neq 0$.

To test this formally, one could use a **Wald Test** for the null hypothesis:

$$H_0 : \sigma_{12} = 0 \quad (\text{no instantaneous causality from } Z_t \text{ to } X_t).$$

As we are testing a single parameter, the Wald statistic takes the form :

$$W = \frac{(\hat{\sigma}_{12} - \sigma_{12})^2}{\text{Var}(\hat{\sigma}_{12})}$$

Under the null hypothesis, this follows a χ_1^2 distribution. We then reject H_0 at level α if:

$$W > \chi_1^2(1 - \alpha),$$

where $\chi_1^2(1 - \alpha)$ is the $(1 - \alpha)$ -quantile of χ_1^2 .

A Study of the original series

When decomposing the series into its trend, seasonal, and random components, we used a multiplicative model. The trend confirms the patterns discussed in Section 1.1. Residuals appear stationary apart from the outliers of early 2020, which we will treat later.

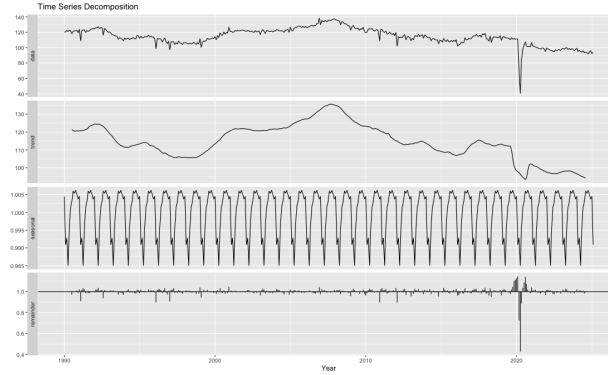


Figure 5: Decomposition of the original series into its trend / cycle / seasonality / residual components

We then study the ACF and PACF.

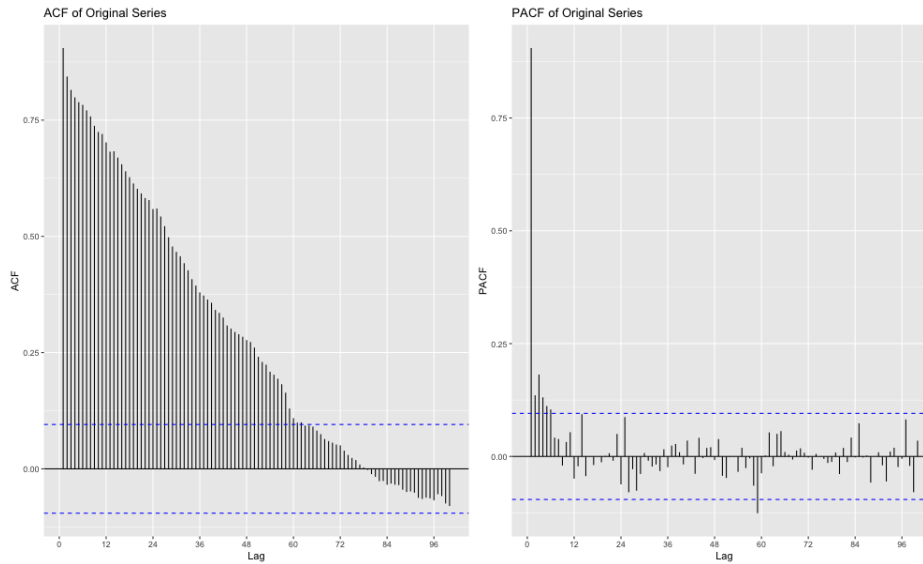


Figure 6: ACF (left) and PACF(right) of the Original Series

Note: The dashed blue line represents the 95% significance bounds.

The ACF decays slowly, and the PACF behaves like white noise : both of these suggest non-stationarity.

These hypotheses are confirmed by the three classic stationarity tests (see Appendix C for further detail).

B Seasonality

The original series is corrected from seasonal variations and working days (CVS-CJO) by INSEE.

The correction process seeks to break down the raw series into multiple components:

$$X_t = TC_t + S_t + WD_t + I_t$$

which correspond to the trend-cycle, the seasonal component, the calendar/working days effect, and the irregular component.

It then seeks to identify S_t and WD_t to correct the series of these variations. Today, the method used is JDemetra+, as recommended by Eurostat. It relies on an X13-ARIMA. The method itself is broken down into two principal modules³:

1. **RegArima** pre-adjusts the series by detecting additive outliers, transitory changes, level shifts, and seasonal outliers. It also corrects for working days thanks to regressors that reproduce the characteristics of the French calendar. Once identified, these effects are estimated by a Reg-Arima model to "linearise" the series and prolong it at the edges to facilitate seasonal adjustment.
2. **X11** does the actual seasonal adjustment through iterative smoothing using moving averages and decomposes the linearized series from the first part into orthogonal components: the trend-cycle, the seasonality, and the irregular component.

Coefficients are re-estimated each month, and model specifications each year.

Thus, our series should exhibit no seasonality. We verify this assumption graphically with the following plots.

As all the boxplots have similar profiles, no seasonality can be detected. As both the linear and polar plots of seasonality show similar behaviour for all months of the year (except the visible outlier of April 2020 due to Covid), we can once again conclude there is no seasonality.

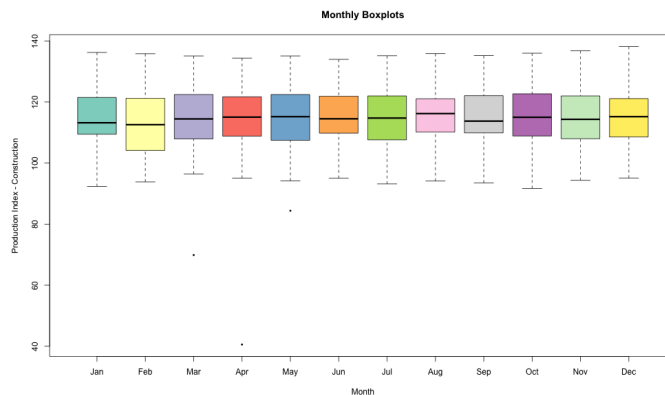


Figure 7: Monthly Boxplots of the Original Series

Note: Each colored box represents a different month. Both median (centre of box) and extreme values (ends of whiskers) are similar across each month.

³For more detail, see: <https://www.insee.fr/fr/statistiques/fichier/4186908/imet133-g.pdf>

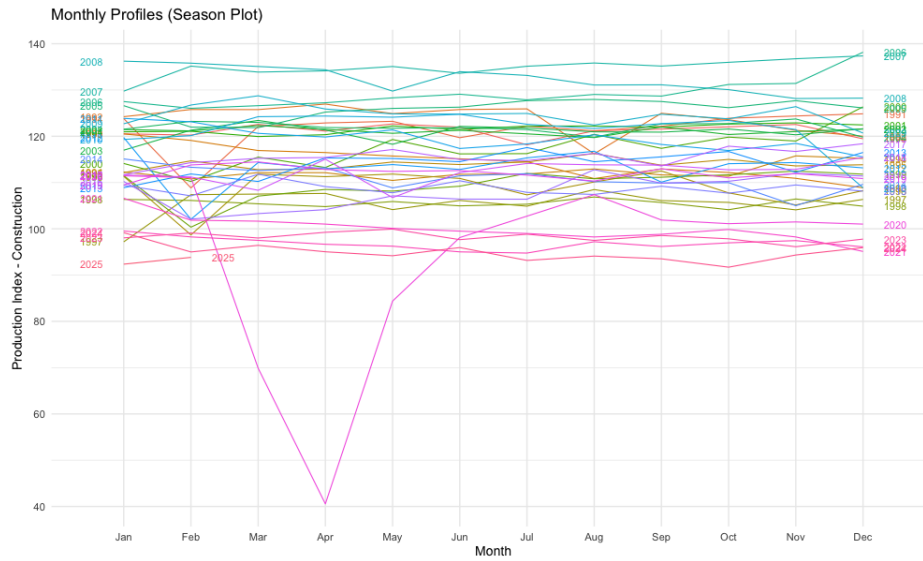


Figure 8: Seasonality of the Original Series

Note: The x-axis is months of the year, the y-axis the value of production.
 Each line represents a year : except changes in level, behaviour each year is very similar.
 No variations between months can be noted (except for the 2020 outlier).

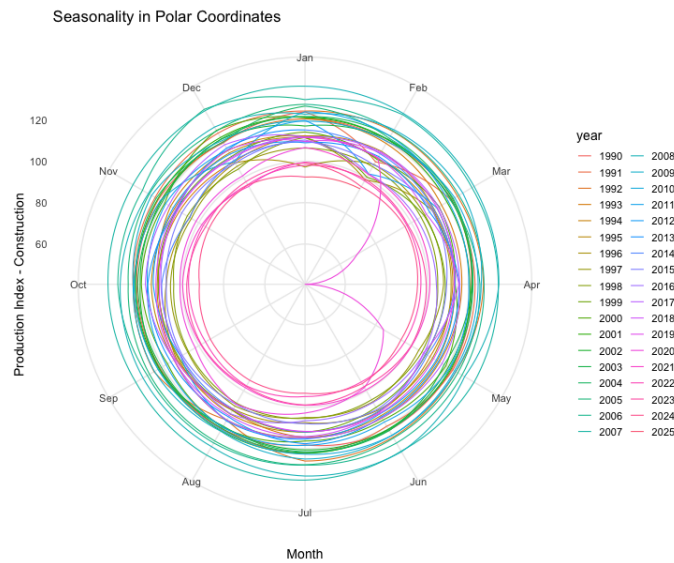


Figure 9: Seasonality of the Original Series in Polar Coordinates

Note: Each month represents a polar coordinate, and each line represents a year.
 No variations between months can be noted (except for the 2020 outlier).

C Stationarity tests

Three classic tests were used to ensure stationarity:

- Augmented Dickey-Fuller (ADF) based on the following regression:

$$\Delta X_t = \alpha + \beta t + \gamma X_{t-1} + \sum_{i=1}^p \delta_i \Delta X_{t-i} + \varepsilon_t$$

where the null hypothesis is $H_0 : \gamma = 0$ (non-stationarity).

- Phillips-Perron (PP), similar to the ADF but without lagged difference terms and using a semi-parametric model of the form:

$$\Delta X_t = \alpha + \beta t + \gamma X_{t-1} + \varepsilon_t$$

Here too, the null hypothesis is $H_0 : \gamma = 0$ (non-stationarity).

- Kwiatkowski, Phillips, Schmidt and Shin (KPSS), which assumes the series can be decomposed into a deterministic trend, a random walk, and a stationary error:

$$y_t = r_t + \beta t + \varepsilon_t, \quad \text{where } r_t = r_{t-1} + u_t$$

The null hypothesis is that the random walk component u_t has zero variance: $H_0 : \sigma_u^2 = 0$, implying stationarity.

The first two tests are unit root tests, meaning their null hypothesis is non-stationarity, whereas KPSS is a stationarity test, so its null hypothesis is stationarity.

Lag for the ADF test was chosen according to the method proposed by Ng and Perron (1995) : first, set an upper bound p_{max} for p , then estimate the ADF test regression with $p = p_{max}$. If the absolute value of the t-statistic for testing the significance of the last lagged difference is greater than 1.6, set $p = p_{max}$ and perform the unit root test. Otherwise, reduce the lag length by one and repeat the process. p_{max} is chosen according to Schwert's formula [Schwert, 1989]:

$$p_{max} = \lfloor 12 \cdot \left(\frac{T}{100}\right)^{0.25} \rfloor$$

For X_t , $T = 422$, so $p_{max} = 17$.

For the original series, this process brings us to lag 8, as this is the highest amount of lags under 17 for which the test statistic is above the 1.6 threshold in absolute values.

For the differenced series, as the absolute value for the test with lag 17 is $5.42 > 1.6$, we can indeed conduct our ADF test with lag 17.

For our original series, two out of three tests (ADF and KPSS) indicate non-stationarity. Although the PP test suggests stationarity, the overall evidence leans toward non-stationarity of the series. Furthermore, Phillips-Perron performs less well in finite samples than other unit root tests [Schwert, 1989].

Test Type	Test Statistic	Lag Used	p-value
ADF (Augmented Dickey-Fuller)	-1.387	8	0.7114
PP (Phillips-Perron)	-30.185 (Z-alpha)	5	< 0.01
KPSS	2.2324	5	0.01

Table 1: Unit root test results for the original series

We can easily conclude our differenced series is stationary, as all three tests concur:

Test Type	Test Statistic	Lag Used	p-value
ADF (Augmented Dickey-Fuller)	-5.42	17	< 0.01
PP (Phillips-Perron)	-386.85 (Z-alpha)	5	< 0.01
KPSS	0.049	5	> 0.1

Table 2: Unit root test results for the differenced series

D Choose p and q

This table summarises the steps involved in our selection of our ARMA(p, q) model:

- The outputs linked to our initial test of the significance of the coefficients are contained in the columns "Std Error", "t-value", "p-value S" : we reject the null hypothesis of $\beta = 0$ if our p-value is below 0.05. We are running a classic Student's test.
- The next two columns, "Ljung-Box Stat" and "p-value L" contain the outputs linked to our Ljung-Box tests on the residuals: we reject the null hypothesis of residual autocorrelation if the p-value is below 0.05.
- The two final columns, "AIC" and "BIC", contain the AIC and BIC used to select an optimal model : we want to minimize both.

Table 3: Overall Summary of the selection of p and q

Model	Coeff	Value	Std Error	t-value	p-value S	Ljung-Box stat	p-value L	AIC	BIC
ARMA(0,1)	ma1	-0.4008	0.0739	-5.4260	0.0000	40.0360	0.0000	2466.548	2478.676
ARMA(0,2)	ma1	-0.3565	0.0452	-7.8945	0.0000	6.8582	0.7388	2433.739	2449.910
	ma2	-0.2863	0.0448	-6.3896	0.0000				
ARMA(1,0)	ar1	-0.1883	0.0478	-3.9387	0.0001	36.0879	0.0001	2481.413	2493.541
ARMA(1,1)	ar1	0.4764	0.0649	7.3378	0.0000	5.5339	0.8528	2432.691	2448.861
	ma1	-0.8387	0.0380	-22.0435	0.0000				
ARMA(1,2)	ar1	0.2917	0.1386	2.1049	0.0359	2.8734	0.9842	2431.968	2452.181
	ma1	-0.6218	0.1385	-4.4895	0.0000				
	ma2	-0.1515	0.0879	-1.7236	0.0855				
ARMA(2,0)	ar1	-0.2293	0.0475	-4.8249	0.0000	27.6341	0.0021	2462.963	2479.134
	ar2	-0.2174	0.0475	-4.5775	0.0000				
ARMA(2,1)	ar1	0.4667	0.0697	6.6960	0.0000	2.6036	0.9893	2431.644	2451.857
	ar2	-0.0987	0.0559	-1.7655	0.0782				
	ma1	-0.7945	0.0530	-14.9906	0.0000				
ARMA(2,2)	ar1	0.8687	0.4284	2.0277	0.0432	2.3787	0.9925	2433.201	2457.457
	ar2	-0.2963	0.2003	-1.4795	0.1398				
	ma1	-1.2002	0.4363	-2.7508	0.0062				
	ma2	0.3462	0.3650	0.9486	0.3434				
ARMA(3,0)	ar1	-0.2658	0.0480	-5.5351	0.0000	20.5814	0.0242	2453.081	2473.294
	ar2	-0.2556	0.0481	-5.3139	0.0000				
	ar3	-0.1664	0.0479	-3.4718	0.0006				
ARMA(3,1)	ar1	0.4462	0.0844	5.2838	0.0000	2.4840	0.9911	2433.424	2457.680
	ar2	-0.0954	0.0565	-1.6870	0.0923				
	ar3	-0.0281	0.0598	-0.4706	0.6381				
	ma1	-0.7759	0.0700	-11.0902	0.0000				

Model	Coeff	Value	Std Error	t-value	p-value S	Ljung-Box stat	p-value L	AIC	BIC
ARMA(3,2)*	ar1	-0.0544	NaN	NaN	NaN	2.5800	0.9897	2435.576	2463.874
	ar2	0.1423	NaN	NaN	NaN				
	ar3	-0.0636	NaN	NaN	NaN				
	ma1	-0.2741	NaN	NaN	NaN				
	ma2	-0.4070	NaN	NaN	NaN				
ARMA(4,0)	ar1	-0.2912	0.0482	-6.0468	0.0000	13.0093	0.2232	2445.394	2469.650
	ar2	-0.2943	0.0491	-5.9875	0.0000				
	ar3	-0.2064	0.0491	-4.2066	0.0000				
	ar4	-0.1504	0.0480	-3.1307	0.0019				
ARMA(4,1)	ar1	0.4053	0.1132	3.5804	0.0004	2.0884	0.9956	2435.002	2463.301
	ar2	-0.1106	0.0615	-1.7978	0.0729				
	ar3	-0.0299	0.0611	-0.4888	0.6252				
	ar4	-0.0409	0.0628	-0.6514	0.5152				
	ma1	-0.7356	0.1037	-7.0936	0.0000				
ARMA(4,2)	ar1	0.5478	0.8873	0.6173	0.5373	2.0743	0.9957	2436.974	2469.315
	ar2	-0.1767	0.4031	-0.4384	0.6613				
	ar3	-0.0167	0.1022	-0.1636	0.8701				
	ar4	-0.0386	0.0686	-0.5623	0.5742				
	ma1	-0.8782	0.8872	-0.9899	0.3228				
	ma2	0.1132	0.6895	0.1642	0.8697				
ARMA(5,0)	ar1	-0.3117	0.0483	-6.4566	0.0000	5.7229	0.8380	2439.674	2467.972
	ar2	-0.3224	0.0497	-6.4833	0.0000				
	ar3	-0.2462	0.0507	-4.8594	0.0000				
	ar4	-0.1894	0.0496	-3.8198	0.0002				
	ar5	-0.1342	0.0481	-2.7919	0.0055				
ARMA(5,1)	ar1	0.3736	0.1728	2.1621	0.0312	2.0252	0.9961	2436.921	2469.262
	ar2	-0.1209	0.0740	-1.6342	0.1030				
	ar3	-0.0415	0.0763	-0.5447	0.5863				
	ar4	-0.0449	0.0662	-0.6773	0.4986				
	ar5	-0.0197	0.0707	-0.2783	0.7810				
	ma1	-0.7042	0.1666	-4.2258	0.0000				
ARMA(5,2)	ar1	-0.5734	0.1132	-5.0658	0.0000	2.7268	0.9871	2438.360	2474.744
	ar2	0.2976	0.1420	2.0951	0.0368				
	ar3	-0.1338	0.0847	-1.5799	0.1149				
	ar4	-0.0689	0.0892	-0.7729	0.4400				
	ar5	-0.0300	0.0634	-0.4735	0.6361				
	ma1	0.2448	0.1032	2.3719	0.0182				
	ma2	-0.7329	0.1008	-7.2726	0.0000				

Note: *The NaNs returned by R for this model mean that the parameters are near the edge of the stationarity region: the standard errors cannot be computed because the sum of the coefficients is close to 1.

E Outlier Types and Their Treatment in Time Series Model

Outliers are observations that deviate significantly from the expected behavior of a time series. In time series analysis, especially when using ARIMA models, uncorrected outliers can bias parameter estimation and impair forecasting accuracy. It is therefore important to detect them and model their impact appropriately. Though INSEE's CVS-CJO treatment methodology should address most of the issues related to outliers (see Appendix C), we still found values that needed correction.

Let X_t be the observed time series, ω the magnitude of the outlier, and Z_t the adjusted series after removing the outlier's effect at time $t = t_{\text{outlier}}$. Four common types of outliers are typically distinguished in time series analysis, each associated with a different dynamic impact on the series.

1. **Additive Outlier (AO)**: A sudden, isolated disturbance that affects only a single time point without altering the underlying dynamics. These are typically associated with recording errors or temporary shocks. The corrected series is:

$$Z_t = X_t + \omega \cdot \mathbf{1}_{\{t=t_{outlier}\}}.$$

2. **Innovative Outlier (IO)**: A shock that enters through the innovation term and propagates according to the model dynamics. Unlike AO, the effect of an IO is filtered by the ARIMA process itself. The correction must take into account the full ARIMA structure:

$$Z_t = X_t + \psi(B)\omega \cdot \mathbf{1}_{\{t=t_{outlier}\}},$$

where $\psi(B) = \frac{\theta(B)}{\phi(B)}$ is the impulse response function of the ARIMA model.

3. **Level Shift (LS)**: A permanent change in the mean level of the series starting at time $t = S$. This could represent a structural break or regime change. The effect persists indefinitely:

$$Z_t = X_t + \omega \cdot \sum_{k=0}^{\infty} \mathbf{1}_{\{t=t_{outlier}+k\}} = X_t + \omega \cdot \mathbf{1}_{\{t \geq t_{outlier}\}}.$$

4. **Transient Change (TC)**: A temporary deviation that decays gradually over time. This type of outlier reflects a temporary disturbance that fades exponentially. The standard correction is modeled by a geometric decay:

$$Z_t = X_t + \omega \cdot \delta^{t-t_{outlier}} \cdot \mathbf{1}_{\{t \geq t_{outlier}\}},$$

where $0 < \delta < 1$ is a damping parameter governing the speed of decay.

Table 4: Detected Outliers in the Series X_t (the last one not taken in the adjusted series)

Observation	Time	Type	Estimate ($\hat{\beta}$)	t -statistic
14	February 1991	AO	-13.1768	-7.39
32	August 1992	AO	-8.3201	-4.66
74	February 1996	AO	-13.5640	-7.61
85	January 1997	AO	-10.0067	-5.61
110	February 1999	AO	-7.3368	-4.09
252	December 2010	AO	-13.7539	-7.72
266	February 2012	AO	-13.5115	-7.54
314	February 2016	TC	-7.9013	-4.36
363	March 2020	TC	-41.1391	-21.10
364	April 2020	AO	-41.0581	-21.08

Coefficients of the Outliers ω_i The estimated outlier effects ω_i incorporated in equation (2) are provided below:

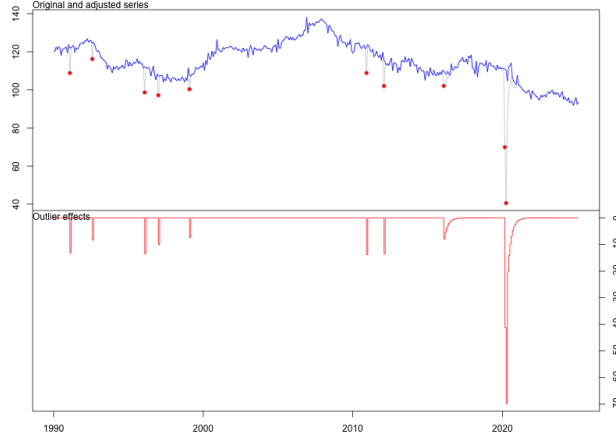


Figure 10: Outliers

Note: The downward red bars indicate substantial negative shocks to the series, which are often associated with exogenous events such as unexpected macroeconomic disturbances.

$\omega_1 = -13.1768$	(February 1991, AO)
$\omega_2 = -8.3201$	(August 1992, AO)
$\omega_3 = -13.5640$	(February 1996, AO)
$\omega_4 = -10.0067$	(January 1997, AO)
$\omega_5 = -7.3368$	(February 1999, AO)
$\omega_6 = -13.7539$	(December 2010, AO)
$\omega_7 = -13.5115$	(February 2012, AO)
$\omega_8 = -7.9013$	(February 2016, TC)
$\omega_9 = -41.1391$	(March 2020, TC)
$\delta \in]0, 1[$	(Damping parameter for TC effects)

F Calculus

$$\text{Var}(X_{T+1} - \hat{X}_{T+1|T}) = \text{Var}(\epsilon_{T+1}) = \sigma^2, \quad (\text{by assumption})$$

$$\begin{aligned} \text{Var}(X_{T+2} - \hat{X}_{T+2|T}) &= \text{Var}((1 + \phi_1 - \theta_1) \cdot \epsilon_{T+1} + \epsilon_{T+2}) \\ &= (1 + \phi_1 - \theta_1)^2 \cdot \text{Var}(\epsilon_{T+1}) + \text{Var}(\epsilon_{T+2}) \quad (\text{since the residuals are independent}) \\ &= (1 + (1 + \phi_1 - \theta_1)^2) \cdot \sigma^2, \end{aligned}$$

$$\begin{aligned} \text{Cov}(X_{T+1} - \hat{X}_{T+1|T}, X_{T+2} - \hat{X}_{T+2|T}) &= \text{Cov}((1 + \phi_1 - \theta_1) \cdot \epsilon_{T+1} + \epsilon_{T+2}, \epsilon_{T+1}) \\ &= (1 + \phi_1 - \theta_1) \cdot \text{Var}(\epsilon_{T+1}) + \text{Cov}(\epsilon_{T+2}, \epsilon_{T+1}) \\ &= (1 + \phi_1 - \theta_1) \cdot \sigma^2 \quad (\text{by bilinearity and independence of the residuals}). \end{aligned}$$

G Code

```
1 #####
2 ##### Linear Time Series #####
3 ##### BENABDESADOK Nayel, CIANFARANI Ana-Sofia #####
4 #####
5
6 ##### - Part 0 : Installations & Data Download - #####
7
8 # Load required packages
9 library(fUnitRoots)
10 library(zoo)
11 library(tseries)
12 library(forecast)
13 library(ggplot2)
14 library(scales)
15 library(RColorBrewer)
16 library(gridExtra)
17 library(astsa)
18 library(tsoutliers)
19 library(ellipse)
20
21 # # Load the data on personal computer (Sofia)
22 # path <- "C:/Users/cianf/OneDrive/Documents/STL"
23 # path <- "C:/Users/cianf/OneDrive/Documents/STL/data.csv"
24 # # Load the data on personal computer (Nayel)
25 path_donnees <- "/Users/nayelbenabdesadok/GitProjects/Time_Series_ENSAE/reprise/data.csv"
26 path <- "/Users/nayelbenabdesadok/GitProjects/Time_Series_ENSAE/reprise"
27 setwd(path)
28 getwd()
29
30
31 # Load the data
32 data <- read.csv(file = path_donnees, sep = ";")
33
34 # Clean the table
35 data_all <- data[4:425, ]
36 data_all$Codes <- NULL
37 colnames(data_all) <- c("Date_MY", "IPI")
38
39
40 # Converting date to the correct format
41 data_all$Date <- as.Date(paste(data_all$Date_MY, 1, sep = "-"), format = "%Y-%m-%d")
42 data_all$Date_MY <- NULL
43
44 # Extracting month and year from date
45 Year <- as.numeric(format(data_all$Date, format = "%Y"))
46 Month <- format(data_all$Date, format = "%m")
47 data_all <- cbind(data_all, Year, Month)
48 data_all$Year <- as.numeric(data_all$Year)
49 data_all$IPI <- as.numeric(data_all$IPI)
50 sort(data_all$IPI)
51 data_all <- data_all[order(data_all$Date), ]
52 rownames(data_all) <- seq(length = nrow(data_all))
53
54
55 # Create a time series for easier analysis, and plot
56 Xt.ts <- ts(data_all$IPI, start = c(1990, 1), end = c(2025, 2), frequency = 12)
57 par(mfrow = c(1, 1))
58 plot.ts(Xt.ts, xlab = "Years", ylab = "IPI Construction (CVS-CJ0, base 2021)")
59
60 png("Xt.png", width = 1000, height = 600)
61 plot.ts(Xt.ts,
62         main = "Main Series",
63         ylab = "IPI Construction (CVS-CJ0, base 2021)",
64         xlab = "Time",
65         lwd = 2)
```

```

66 dev.off()
67 #####
68 # The dates range from January 1990 (1990m1) to February 2025 (2025m2) as shown above.
69 # The series represents the Index of Industrial Production (IPI) in the construction sector,
70 # seasonally and working-day adjusted (CVS-CJ0), with a base of 100 in 2021.
71 # At first glance, the series appears non-stationary, with multiple structural breaks.
72 # A sharp drop is clearly visible around April 2020, likely linked to the COVID-19 crisis.
73 # Before applying any transformations, we will examine the overall trend and variability.
74 #####
75
76
77
78 ##### - Part I : The Data - #####
79
80 ##### - Part 1.1 : Trend/Cycle/Seasonal Decomposition (What does the chosen series
      represent ?)- #####
81 # Decomposition of the time series to analyze the trend component.
82 decompose <- decompose(Xt.ts,type="multiplicative")
83 png("decompose.png", width = 1000, height = 600)
84 autoplot(decompose, main = "Time Series Decomposition", xlab = "Year")
85 dev.off()
86
87
88 #####
89 # Trend Analysis: the series displays three distinct phases
90 # a decline from 1990 to 2000, a growth period from 2000 to 2008,
91 # and a significant downward trend from 2008 to 2025.
92 # A sharp outlier is observed in early 2020, likely due to the COVID-19 crisis.
93 #####
94
95 #####
96 # Is there Seasonality ? #
97 #####
98 n_years <- length(unique(floor(time(Xt.ts))))
99 colors <- scales::hue_pal()(n_years)
100
101 # Standard season plot (monthly profile)
102 png("season_plot.png", width = 1000, height = 600)
103 ggseasonplot(Xt.ts,
104             year.labels = TRUE,
105             year.labels.left = TRUE) +
106   ylab("Production Index - Construction") +
107   xlab("Month") +
108   ggtitle("Monthly Profiles (Season Plot)") +
109   scale_color_manual(values = colors) +
110   theme_minimal(base_size = 14)
111 dev.off()
112
113 # Polar season plot
114 png("season_polar_plot.png", width = 1000, height = 600)
115 ggseasonplot(Xt.ts,
116             polar = TRUE) +
117   ylab("Production Index - Construction") +
118   xlab("Month") +
119   ggtitle("Seasonality in Polar Coordinates") +
120   scale_color_manual(values = colors) +
121   theme_minimal(base_size = 14)
122 dev.off()
123
124 # Monthly boxplots
125 png("monthly_boxplots.png", width = 1000, height = 600)
126 cols <- RColorBrewer::brewer.pal(12, "Set3")
127 boxplot(Xt.ts ~ cycle(Xt.ts),
128         col = cols,
129         pch = 20,
130         cex = 0.5,
131         main = "Monthly Boxplots",
132         ylab = "Production Index - Construction",

```

```

133     xlab = "Month",
134     names = month.abb)
135 dev.off()
136 #=====#
137 # Regardless of the plot displayed, we observe no influence of the month
138 # on the behavior of the series.
139 # Absence of seasonality.
140 #=====#
141
142
143
144 ##### - Part 1.2 : Stationary #####
145 #=====#
146 # Is Xt.ts Stationary ?      #
147 #=====#
148 # ACF
149 acf_plot_orig <- ggAcf(Xt.ts, lag.max = 100, plot = TRUE) +
150   ggtitle("ACF of Original Series")
151 #PACF
152 pacf_plot_orig <- ggPacf(Xt.ts, lag.max = 100, plot = TRUE) +
153   ggtitle("PACF of Original Series")
154
155 png("ACF_PACF_Xt.png", width = 1000, height = 600)
156 gridExtra::grid.arrange(acf_plot_orig, pacf_plot_orig, ncol = 2)
157 dev.off()
158
159 #=====#
160 # ACF decays slowly + PACF behaves like white noise.
161 # Suggests non-stationarity      Let's run some tests: ADF, PP, KPSS.
162 #=====#
163
164 # Center the series
165 Xt_centered <- Xt.ts - mean(Xt.ts)
166 ggAcf(Xt_centered, lag.max=100, plot=T)
167
168 #=====#
169 # We decided to center the series but this is optional :
170 # the series X_t and X_t_centered are the same in our analysis !
171 #=====#
172
173
174 #=====#
175 # To determine the lag to use for our ADF test, we refer to Schwert (1989).
176 # The code below implements their lag selection algorithm.
177 # For more details see our report, Appendix : Stationarity tests.
178 # We set type = ct as there is a trend.
179 #=====#
180
181 # Initialize the lag value at pmax = 17
182 lag <- 17
183
184 # Store the ADF test result
185 adf_result <- NULL
186
187 # Loop until the t-statistic is above 1.6 in absolute value
188 while (lag > 0) {
189   # Perform the ADF test
190   adf_result <- adfTest(Xt_centered, lags = lag, type = "ct")
191
192   # Check the t-statistic
193   if (abs(adf_result@test$statistic) > 1.6) {
194     break
195   }
196
197   # Decrease the lag value
198   lag <- lag - 1
199 }
200 cat("Final lag value:", lag, "\n")

```

```

201 print(adf_result)
202
203 # Phillips-Perron Test (PP)
204 library(tseries)
205 cat("\nPhillips-Perron Test Results:\n")
206 print(pp.test(Xt_centered))
207
208 # KPSS Test
209 cat("\nKPSS Test Results:\n")
210 kpss_result <- kpss.test(Xt_centered)
211 print(kpss_result)
212
213 #=====#
214 # Stationarity Tests Summary:
215 # - ADF Test (lag = 8, type = "ct")      p-value = 0.7114      FAIL to reject H0
216 # - PP Test      p-value < 0.01      REJECT H0 (suggests stationarity but less reliable for
217 #   finite samples)
218 # - KPSS Test      p-value = 0.01      REJECT H0 of stationarity
219 #
220 # Conclusion: Two out of three tests (ADF and KPSS) indicate non-stationarity.
221 # Although the PP test suggests stationarity, the overall evidence leans toward
222 # non-stationarity of the centered series, likely due to trend or structural changes.
223 #=====#
224 #=====#
225 # Stationarity Tests on the Differenced Centered Series
226 #=====#
227
228 diff_Xt_centered <- diff(Xt_centered)
229
230 png("diff_Xt_centered.png", width = 1000, height = 600)
231 plot.ts(diff_Xt_centered,
232         main = "Differenced Centered Series",
233         ylab = "Differenced IPI",
234         xlab = "Time",
235         lwd = 2)
236 dev.off()
237
238 #=====#
239 # This time, our ADF test works with the max lag immediately.
240 # We set type = nc as there is the series is centered and stationary.
241 #=====#
242
243 # Augmented Dickey-Fuller Test (ADF)
244 adf_result_diff <- adfTest(diff_Xt_centered, lags = 17, type = "nc")
245 cat("ADF Test Results (Differenced Series):\n")
246 print(adf_result_diff)
247
248 # Phillips-Perron Test (PP)
249 cat("\nPhillips-Perron Test Results (Differenced Series):\n")
250 print(pp.test(diff_Xt_centered))
251
252 # KPSS Test
253 cat("\nKPSS Test Results (Differenced Series):\n")
254 kpss_result_diff <- kpss.test(diff_Xt_centered)
255 print(kpss_result_diff)
256
257 #=====#
258 ## Stationarity Tests Summary:
259 # - ADF Test (lag = 17, type = "nc")      p-value = 0.01      REJECT H0 - Stationarity
260 # - PP Test      p-value < 0.01      REJECT H0 - Stationarity
261 # - KPSS Test      p-value > 0.1      DO NOT REJECT H0 - Stationarity
262 #
263 # All three stationarity tests (ADF, PP, KPSS) applied to the differenced
264 # centered series confirm stationarity (p-values all in favor).
265 # Conclusion: The original series is integrated of order 1 (I(1)).
266 #=====#
267 ##### - Part 1.3 : Before/after Stationarity #####

```

```

268
269 png("Xt_vs_diffXt.png", width = 1000, height = 300)
270 par(mfrow = c(1, 2))
271
272 plot.ts(Xt.ts,
273         main = expression("Original Series " ~ X[t]),
274         ylab = "IPI Construction",
275         xlab = "Time",
276         lwd = 2,
277         col = "steelblue")
278
279 plot.ts(diff(Xt.ts),
280         main = expression("Differenced Series " ~ Y[t] == Delta * X[t]),
281         ylab = "Differenced IPI",
282         xlab = "Time",
283         lwd = 2,
284         col = "darkorange")
285
286 dev.off()
287 ##### - Part II : ARMA Models #####
288
289 ##### - Part 2.1 : Pick p and q #####
290
291
292 acf_diff_plot <- ggAcf(diff_Xt_centered, lag.max = 40, plot = TRUE) +
293   ggtitle("ACF of Differenced Series")
294 pacf_diff_plot <- ggPacf(diff_Xt_centered, lag.max = 40, plot = TRUE) +
295   ggtitle("PACF of Differenced Series")
296
297 png("ACF_PACF_diff_Xt_centered.png", width = 1000, height = 300)
298 gridExtra::grid.arrange(acf_diff_plot, pacf_diff_plot, ncol = 2)
299 dev.off()
300
301 #=====#
302 # Property of MA(2) models: if the autocorrelation function (ACF) becomes
303 # zero for all lags  $h > 2$ , the process can be identified as MA(2).
304 # Similarly, for AR(5) models: the partial autocorrelation function (PACF)
305 # cuts off after lag 5.
306 #
307 #  $p_{max}=5$  and  $q_{max}=2$ . Let's test all ARMA(p,q) such that  $p \leq p_{max}$  and  $q \leq q_{max}$ 
308 #=====#
309
310 # we call the differents models arma_p_q
311 #=====#
312 # We test all ARIMA(p,1,q) models for the construction series Xt.ts,
313 # where p = 5 and q = 2. We compute AIC and BIC for each combination.
314 #=====#
315 arma_0_1 <- sarima(Xt.ts, 0, 1, 1)
316 arma_0_2 <- sarima(Xt.ts, 0, 1, 2)
317
318 arma_1_0 <- sarima(Xt.ts, 1, 1, 0)
319 arma_1_1 <- sarima(Xt.ts, 1, 1, 1)
320 arma_1_2 <- sarima(Xt.ts, 1, 1, 2)
321
322 arma_2_0 <- sarima(Xt.ts, 2, 1, 0)
323 arma_2_1 <- sarima(Xt.ts, 2, 1, 1)
324 arma_2_2 <- sarima(Xt.ts, 2, 1, 2)
325
326 arma_3_0 <- sarima(Xt.ts, 3, 1, 0)
327 arma_3_1 <- sarima(Xt.ts, 3, 1, 1)
328 arma_3_2 <- sarima(Xt.ts, 3, 1, 2)
329
330 arma_4_0 <- sarima(Xt.ts, 4, 1, 0)
331 arma_4_1 <- sarima(Xt.ts, 4, 1, 1)
332 arma_4_2 <- sarima(Xt.ts, 4, 1, 2)
333
334 arma_5_0 <- sarima(Xt.ts, 5, 1, 0)
335 arma_5_1 <- sarima(Xt.ts, 5, 1, 1)

```

```

336 arma_5_2 <- sarima(Xt.ts, 5, 1, 2)
337
338
339 #=====#
340 # To eliminate unsuitable models, we first exclude models with at least
341 # one non-significant coefficient (p-value > 5%), based on hypothesis testing (H0 vs H1).
342 #=====#
343
344
345 arma_0_1$tttable
346 arma_0_2$tttable
347
348 arma_1_0$tttable
349 arma_1_1$tttable
350 arma_1_2$tttable # out
351
352 arma_2_0$tttable
353 arma_2_1$tttable # out
354 arma_2_2$tttable # out
355
356 arma_3_0$tttable
357 arma_3_1$tttable # out
358 arma_3_2$tttable # out
359
360 arma_4_0$tttable
361 arma_4_1$tttable # out
362 arma_4_2$tttable # out
363
364 arma_5_0$tttable
365 arma_5_1$tttable # out
366 arma_5_2$tttable # out
367
368 #=====#
369 # We still keep the excluded models ("out models") to compute the Ljung-Box
370 # statistics on their residuals.
371 # We will later exclude models whose Ljung-Box p-value is less than 5%.
372 #=====#
373 models <- c("arma_0_1", "arma_0_2",
374             "arma_1_0", "arma_1_1", "arma_1_2",
375             "arma_2_0", "arma_2_1", "arma_2_2",
376             "arma_3_0", "arma_3_1", "arma_3_2",
377             "arma_4_0", "arma_4_1", "arma_4_2",
378             "arma_5_0", "arma_5_1", "arma_5_2")
379
380 for (m in models) {
381   cat("\n--- Ljung-Box Test for", m, "---\n")
382   print(Box.test(residuals(get(m)$fit), lag = 10, type = "Ljung-Box"))
383 }
384 # The remaining models are ARMA(1,1), ARMA(0,2), and ARMA(5,0).
385 # To select the final model, we choose the one that minimizes both the AIC and BIC criteria.
386
387
388 info_criteria <- data.frame(
389   Model = character(),
390   AIC = numeric(),
391   BIC = numeric(),
392   stringsAsFactors = FALSE
393 )
394
395 for (m in models) {
396   fit <- get(m)$fit
397   model_aic <- AIC(fit)
398   model_bic <- BIC(fit)
399
400   info_criteria <- rbind(info_criteria, data.frame(
401     Model = m,
402     AIC = model_aic,
403     BIC = model_bic

```

```

404   ))
405 }
406 print(info_criteria)
407
408 # ARMA (1,1) is choosen !
409
410 ##### - Part 2.2 : Taking outliers into account #####
411
412 # Detect the outliers
413 arima111 <- arima(Xt.ts, order = c(1, 1, 1))
414 outlier_result <- tso(Xt.ts, types = c("AO", "LS", "TC", "IO"),
415                       tsmethod = "arima", args.tsmethod = list(order = c(1, 1, 1)))
416 all_outliers <- outlier_result$outliers
417 print(all_outliers)
418
419 #=====#
420 # A total of 10 outliers are detected: 8 additive outliers (AO) and 2 transient changes (TC)
421 #
422 # The last outlier, however, does not correspond to any visually identifiable shock in the
423 # series.
424 # Therefore, we retain only the first 9 outliers for the adjustment.
425 # We then correct the series accordingly using the corresponding estimated effects.
426 #=====#
427
428 selected_indices <- head(all_outliers$ind, 9)
429 adjusted_effects <- outlier_result$effects
430 keep_positions <- rep(FALSE, length(adjusted_effects))
431 keep_positions[selected_indices] <- TRUE
432 adjusted_effects[!keep_positions] <- 0
433
434 yadj_9out <- Xt.ts - adjusted_effects
435 # Adjusted coefficients : phi1 = 0.1781 , theta1 = -0.7285
436
437 final_model_9out <- arima(yadj_9out, order = c(1, 1, 1))
438 summary(final_model_9out)
439 # AIC = 2395.14 < 2432.691 (AIC of the initial model) --> improvement confirmed
440
441 png("outliers_plot_full.png", width = 800, height = 600)
442 plot(outlier_result)
443 dev.off()
444
445 ##### - Part III Forecasting #####
446
447 ##### - Part 3.1 : X_T+1 and X_T+2 #####
448
449 forecast_2 <- forecast(final_model_9out, h = 2)
450 print(forecast_2)
451
452 #=====#
453 # Point Forecast      80% CI      95% CI
454 #-----#
455 # Mar 2025:  93.86069   [88.56977 ; 99.15160]   [85.76894 ; 101.95240]
456 # Apr 2025:  93.86615   [88.06500 ; 99.66730]   [84.99406 ; 102.73820]
457 #=====#
458
459 #=====#
460 # Zoom Forecast Plot (Last 36 Months)
461 #=====#
462
463 start_zoom <- time(yadj_9out)[length(yadj_9out) - 35]
464 end_zoom <- time(forecast_2$mean)[2]
465
466 recent_values <- window(yadj_9out, start = start_zoom)
467 forecast_values <- forecast_2$mean
468 forecast_lower <- forecast_2$lower[,2] # 95% CI

```

```

470 forecast_upper <- forecast_2$upper[,2] # 95% CI
471
472 y_min <- min(recent_values, forecast_lower)
473 y_max <- max(recent_values, forecast_upper)
474
475 # Get the last observed point and first forecast point
476 # Then create a line between them, as autoplot does not do this automatically.
477 last_obs_time <- time(tail(recent_values, 1))
478 last_obs_value <- tail(recent_values, 1)
479 first_fc_time <- time(forecast_values)[1]
480 first_fc_value <- forecast_values[1]
481 connect_df <- data.frame(
482   Time = c(last_obs_time, first_fc_time),
483   Value = c(last_obs_value, first_fc_value)
484 )
485
486 # Plot, including the connecting line
487 png("forecast_zoom_adjusted.png", width = 1000, height = 600)
488 autoplot(forecast_2, series = "Forecast") +
489   autolayer(recent_values, series = "Observed", color = "black") +
490   geom_line(data = connect_df, aes(x = Time, y = Value), linetype = "solid", color = "black"
491   ) +
492   ggtitle("Two-Step Forecast: Last 36 Months (Zoom)") +
493   ylab("IPI Construction (Corrected Series)") +
494   xlab("Time") +
495   coord_cartesian(xlim = c(start_zoom, end_zoom), ylim = c(y_min, y_max)) +
496   scale_color_manual(name = "Series",
497     values = c("Observed" = "black", "Forecast" = "blue")) +
498   theme_minimal(base_size = 14)
499 dev.off()
500 ##### - Part 3.2 : Confidence Region #####
501
502 phi <- coef(final_model_9out)["ar1"]
503 theta <- coef(final_model_9out)["ma1"]
504 sigma2 <- final_model_9out$sigma2
505
506 sigma_g1 <- sqrt(sigma2)
507 sigma_g2 <- sqrt(sigma2 * (1 + (1 + phi - theta)^2))
508 rho <- sigma2 * (1 + phi - theta)
509
510 Sigma <- matrix(c(sigma_g1^2, rho,
511   rho, sigma_g2^2), nrow = 2)
512
513 centre <- c(forecast_2$mean[1], forecast_2$mean[2])
514
515 ell <- ellipse(Sigma, centre = centre, level = 0.95, npoints = 1000)
516
517 png("confidence_ellipse_forecast.png", width = 800, height = 600)
518 plot(ell,
519   type = 'l',
520   xlab = expression(hat(X)[T+1]),
521   ylab = expression(hat(X)[T+2]),
522   main = "95% Confidence Ellipse for Two-Step Forecast")
523 points(x = centre[1], y = centre[2], pch = 19, col = "red", cex = 1.5)
524 grid()
525 dev.off()

```


References

- [Banks and Fienberg, 2003] Banks, D. L. and Fienberg, S. E. (2003). *Encyclopedia of Physical Science and Technology (Third Edition)*, chapter II. Statistics, Multivariate. Academic Press.
- [Granger, 1969] Granger, C. (1969). Investigating causal relations by econometric models and cross-spectral methods. *Econometrica*.
- [INSEE, 2000] INSEE (2000). La construction entre 1990 et 1997 : une crise profonde avant la reprise. Technical report, INSEE.
- [Schwert, 1989] Schwert, W. (1989). Test for unit roots: A monte carlo investigation. *Journal of Business and Economic Statistics*.