

Time series Assignment

Plazzogna Marco, Yan Wenbin

Professor: Ngatchou Wandji Joseph

Contents

1 Data	1
2 ARMA Models	2
3 Predictions	3
Appendix	6
A.1 Custom Functions	6
A.2 Codes for Part I	10
A.3 Codes for Part II	10
A.4 Codes for Part III	11

May 21, 2024

1 Data

All the code in this section can be found in Appendix A.1 and A.2.

The seasonally and working-day adjusted (SA-WDA) industrial production index for the manufacture of weapons and ammunition (base 100 in 2021) is a crucial measure within the NAF rev. 2 classification, level Class, item 25.40. This index tracks monthly changes in production, serving as an essential indicator for monitoring the industry's business cycle and identifying early turning points. Specifically, it provides insights into the defense sector's production activities and economic impact. Published monthly, it ensures stakeholders have up-to-date and relevant information.

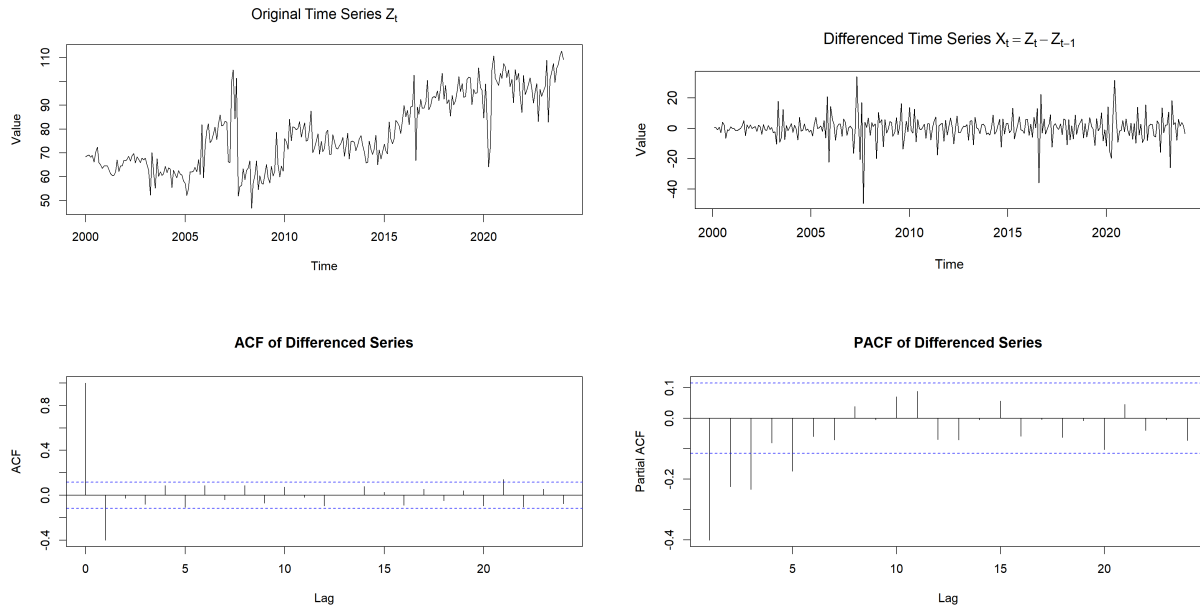


Figure 1: Data Representation

The time series plot (jan/2000-jan/2024) shows an overall increasing, although not linear trend, with values rising from around 50 in 2000 to over 110 by 2023. There are noticeable short-term fluctuations and variability, indicating irregular movements possibly due to seasonal effects or cyclical factors. High volatility is observed particularly between 2008-2010 and from 2020 onwards, likely due to macroeconomics effects, suggesting periods of increased activity. Since 2020, a sharper increase and higher volatility suggest a recent change in the underlying process affecting the variable. We implemented ADF and PP tests to check whether it has a trend. ADF prefers not to reject it while PP prefers to reject it. Since we also agree that there is a trend, we finally decided to difference it.

The differenced time series $X_t = \Delta Z_t := Z_t - Z_{t-1}$ of the original data Z_t appears to be stationary, with a relatively constant mean (most likely equal to 0) and variance over time, suggesting that differencing has successfully removed any trends and made the series more stable. Fluctuations around the zero value are more uniform compared to the original series, further supporting the notion of stationarity. Although there is some volatility, especially notable around 2008 and 2020, it is less pronounced compared to the original series, indicating that differencing has reduced some irregular movements.

The ACF plot of the differenced series X_t shows no significant peaks outside the 95% confidence intervals, indicating the absence of seasonality. If there were seasonality, we would expect to see significant autocorrelations at specific lags (e.g., every 12 months for yearly seasonality). The autocorrelations rapidly decay to near zero after the first lag, which is typical for a stationary series, supporting the conclusion that the differenced series does not have significant autocorrelation at higher lags. The PACF plot shows no significant peaks outside the 95% confidence intervals, suggesting no underlying autoregressive process influencing the series at specific lags, further supporting the absence of seasonality and indicating that the series does not have strong dependencies at particular lags. Minor significant lags within the first five periods are within the confidence intervals, implying weak autocorrelation that does not indicate a clear pattern or seasonal structure. For these reasons, we would prefer to use an ARMA model with parameters $p \leq 5$ and $q \leq 1$.

2 ARMA Models

All the code in this section can be found in Appendix A.1 and A.3.

In this part, in order to select and evaluate ARMA (AutoRegressive Moving Average) models for a given time series, we begin by defining two functions, `Ljung.Box2` and `Ljung.Box`, which perform the Ljung-Box test to check for residual autocorrelation in time series data. The `Ljung.Box2` function calculates the Ljung-Box test statistic, degrees of freedom, and p-value for a given number of lags, while the `Ljung.Box` function extends this test across multiple lag lengths, optionally plotting the p-values to visualize the results.

Next, the `modelchoice` function evaluates ARMA(p, q) models by fitting them to the data and checking the significance of the AR (AutoRegressive) and MA (Moving Average) coefficients. It also ensures the absence of residual autocorrelation using the Ljung-Box test. This function returns a summary of the model parameters, including whether the AR and MA coefficients are significant and whether the residuals show no autocorrelation.

The `armamodelchoice` function generates all possible combinations of AR and MA orders up to specified maximum values (pmax and qmax). It evaluates each combination using the `modelchoice` function, and returns a data frame of the results, including the significance of coefficients and residual diagnostics.

The main script then filters for well-adjusted and valid models from the results, identifying those that meet the criteria of having significant coefficients and no residual autocorrelation. For each selected ARMA model, it fits the model to the data using the `Arima` function and calculates their AIC (Akaike Information Criterion) and BIC (Bayesian Information Criterion) values. These criteria are used to determine the best-fitting models, as they balance model fit and complexity, with lower values indicating better models.

The script creates a comparison table of AIC and BIC values for each model, allowing for easy identification of the best models based on these criteria. It then identifies the models with the lowest AIC and BIC values as the best models. Summaries of the best models according to AIC and BIC are displayed, providing detailed information about the model parameters and fit.

Residual diagnostics are performed on the best AIC and BIC model to ensure its adequacy (in fact in our case they prefer the same model). This includes checking the residuals for any remaining autocorrelation using the Ljung-Box test. The test results are plotted to visualize the p-values across different lags, with a horizontal line indicating the 0.05 significance level. If the residuals show no significant autocorrelation, it confirms the model's adequacy.

Finally, the best ARMA model for the processed data based on the AIC criterion is selected and returned, which in our case also minimizes the BIC criterion, resulting in an ARIMA(0,0,1) model.

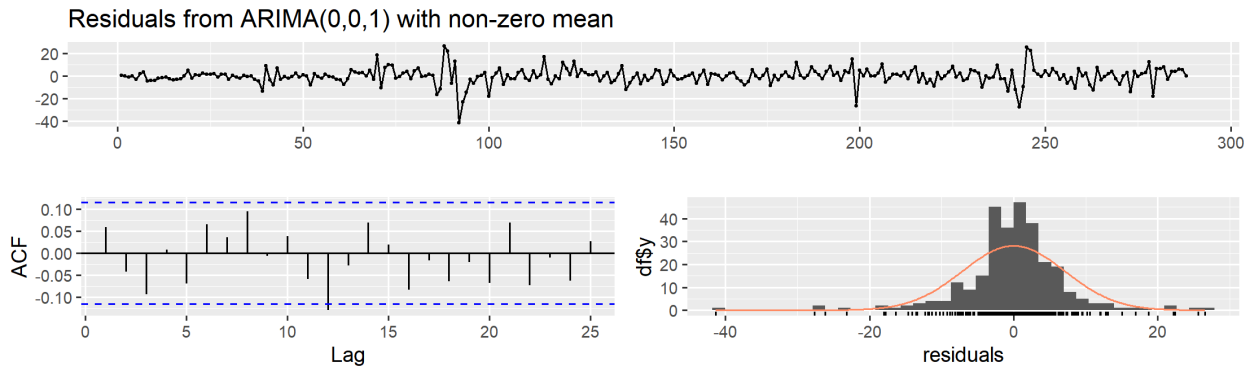


Figure 2: Residuals Analysis for ARIMA(0,0,1) Model

Figure 2 shows the residuals from the ARIMA(0,0,1) model. The top plot displays the residual time series, indicating that the residuals are roughly centered around zero with no apparent trends or patterns. This suggests that the model has adequately captured the underlying process. The bottom-left plot shows the ACF (AutoCorrelation Function) of the residuals. Most of the autocorrelations fall within the 95% confidence intervals (dashed blue lines), indicating that there is no significant autocorrelation left in the residuals, which is a desirable property for a well-fitted model. The bottom-right plot is a histogram of the residuals with a superimposed normal density curve. The residuals appear to follow a roughly normal distribution, with a mean close to zero, suggesting that the assumption of normally distributed errors is reasonable for this model.

However, it is important to note that the results might be influenced by external factors, particularly around the years 2008 and 2020. The global financial crisis in 2008 and the COVID-19 pandemic in 2020 could introduce irregularities and unexpected variations in the data, which the model may not fully account for, potentially affecting its performance.

Table 1 presents the results of the Ljung-Box test for the residuals from the ARIMA(0,0,1) model across lags 1 to 24. The p-values for all lags are above the 0.05 significance level, indicating that we fail to reject the null hypothesis of no autocorrelation in the residuals. This further supports the adequacy of the model.

Lag	Statistic	DF	p-value	Lag	Statistic	DF	p-value	Lag	Statistic	DF	p-value
1	0.00	0	NaN	9	9.90	8	0.2718746	17	20.54	16	0.1968713
2	1.55	1	0.2135640	10	10.36	9	0.3218545	18	21.79	17	0.1930638
3	4.11	2	0.1282974	11	11.41	10	0.3265800	19	21.92	18	0.2355422
4	4.13	3	0.2480566	12	16.49	11	0.1239577	20	23.35	19	0.2223964
5	5.54	4	0.2366010	13	16.73	12	0.1600526	21	24.87	20	0.2063134
6	6.84	5	0.2331901	14	18.20	13	0.1500613	22	26.53	21	0.1868712
7	7.21	6	0.3016077	15	18.31	14	0.1929653	23	26.56	22	0.2283307
8	9.89	7	0.1948043	16	20.45	15	0.1553145	24	27.81	23	0.2229840

Table 1: Ljung-Box Test of Residuals from ARIMA(0,0,1) for Lag Range 1-24

Also note that, since we differenced the original data, this effectively means we implemented the ARIMA(0,1,1) model to the original data in our analysis:

$$X_t = \Delta Z_t = \theta_0 + \epsilon_t + \theta_1 \epsilon_{t-1}, \quad (1)$$

where in our case, $\theta_0 = 0.1422$ and $\theta_1 = -0.6543$.

In summary, the ARIMA(0,0,1) model is selected as the best-fitting model for the corrected series $X_t = \Delta Z_t$ and the ARIMA(0,1,1) model as the best-fitting model for the chosen series Z_t , based on both AIC and BIC criteria. The model chosen effectively captures the underlying structure of the time series data, with well-behaved residuals indicating that the model assumptions are met. The Ljung-Box test further validates the model by showing no significant autocorrelation in the residuals across various lags.

One final note regarding the behaviour of the series: we previously observed two pivotal points likely corresponding to the 2008 financial crisis and the COVID-19 pandemic. During these periods, there were significant changes in the trend of the series. Consequently, it may be beneficial to segment the data into three distinct periods and estimate separate models for each: from 2000 to 2008/2009, from 2008/2009 to 2020, and from 2020 to the end in 2024.

3 Predictions

All the code in this section can be found in Appendix A.1 and A.4.

Denote T as the length of the series. Assume the series' residuals are Gaussian, that is, $\epsilon_t \sim \mathcal{N}(0, \sigma^2)$ with $\sigma^2 > 0$. Given $\{X_t, \epsilon_t : 1 \leq t \leq T\}$, our best expectation of X_{T+1} and X_{T+2} would be:

$$\begin{cases} X_{T+1|T} = \theta_0 + \theta_1 \epsilon_T, \\ X_{T+2|T} = \theta_0. \end{cases}$$

Hence, the error would be respectively:

$$\begin{cases} X_{T+1} - X_{T+1|T} = \epsilon_{T+1}, \\ X_{T+2} - X_{T+2|T} = \epsilon_{T+2} + \theta_1 \epsilon_{T+1}. \end{cases}$$

Since the series' residuals are Gaussian, it is then easy to see that:

$$\begin{cases} X_{T+1} - X_{T+1|T} \sim \mathcal{N}(0, \sigma^2), \\ X_{T+2} - X_{T+2|T} \sim \mathcal{N}(0, (1 + \theta_1^2)\sigma^2). \end{cases}$$

Hence, the α -confidence interval for them would be:

$$\begin{cases} \tilde{X}_{T+1} \in [\theta_0 + \theta_1 \epsilon_T - \sigma q_{1-\frac{\alpha}{2}}, \theta_0 + \theta_1 \epsilon_T + \sigma q_{1-\frac{\alpha}{2}}], \\ \tilde{X}_{T+2} \in [\theta_0 - \sqrt{1 + \theta_1^2 \sigma} q_{1-\frac{\alpha}{2}}, \theta_0 + \sqrt{1 + \theta_1^2 \sigma} q_{1-\frac{\alpha}{2}}]. \end{cases}$$

We should also notice the key hypotheses used to derive this confidence region:

1. The differentiated time series follows an ARMA(0, 0, 1) model.
2. The residuals of the series follow a Gaussian distribution centered at zero with a known variance.

By implementing the results above, we can obtain the future projection with 95% confidence bounds as below,

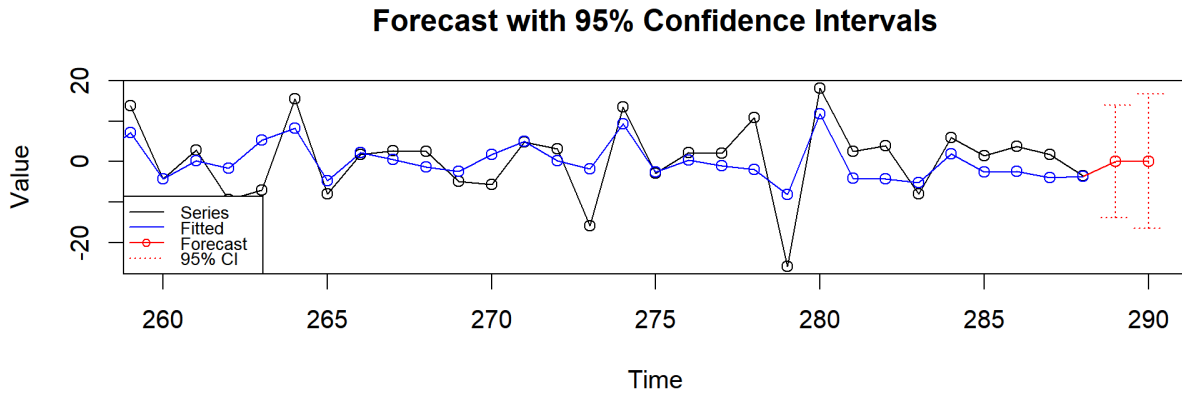


Figure 3: Future Projection with 95% Confidence Bounds Based on Historical Data

Figure 3 shows the last 30 observations of the differenced series along with the corresponding fitted values from the model. This focus on the final portion of the series is due to the need for clear visualization of the two new predicted values and their respective 95 % confidence intervals, which was difficult to achieve when displaying the entire series.

The two prediction are consistent with the type of model that we found for the differenced series, that is as MA(1). In fact, as the first prediction $T+1$ is still influenced by the observation at time T , then from that point on the predictions are equal to the mean estimated by the model that is 0.1422.

The two predictions align with the MA(1) model identified for the differenced series. Specifically, the first prediction at $T + 1$ is influenced by the observation at time T . From that point onward, the predictions converge to the model's estimated mean of 0.1422, which is consistent with the theoretical behavior of MA(1) models.

Open Question:

To determine if the additional information Y_t allows for an improved prediction of X_{T+1} , we need to evaluate the correlation and causality between the two time series. Firstly, there should be a significant correlation between Y_t and X_t , indicating that changes in Y_t are associated with changes in X_t . More importantly, we must establish that Y_t Granger-causes X_t and, as a consequence, verify that there is no spurious correlation between them. Granger causality implies that past values of Y_t contain information that helps predict future values of X_t beyond the information already contained in past values of X_t . To test for Granger causality, we can perform a Granger causality test, which involves estimating a regression model where X_t is regressed on its own past values and past values of Y_t . If the coefficients on the lagged values of Y_t are statistically significant, we conclude that Y_t Granger-causes X_t . This causality test helps identify whether incorporating Y_t as an exogenous variable in the forecasting model for X_t can improve the model's predictive performance.

Finally, having the information provided by Y_{T+1} sooner, can enhance the prediction of X_{T+1} . An illustrative example could be related to news from war zones impacting the French industry's production. For instance, an announcement by President Emmanuel Macron regarding increased defense spending or new defense contracts could provide early

indications of how the production index will evolve. By incorporating such timely information from Y_t , which might capture these announcements and their immediate effects, we can improve the accuracy of our predictions for the industrial production index X_{T+1} .

Appendix

A.1 Custom Functions

Load necessary libraries

```

1 library(tseries) # For time series analysis
2 library(zoo)     # For working with indexed totally ordered observations
3 library(forecast) # For forecasting functions
4 require(fUnitRoots) # For unit root tests
5 require(urca)     # For unit root and cointegration tests

```

pierce.test function

```

1 # This function implements the Pierce (or portmanteau) test to check for residual autocorrelation.
2 # It accepts four arguments:
3 # - x: The time series data or residuals (as a zoo object) to analyze.
4 # - m: The number of lags to include in the test.
5 # - s: The step size (typically 1).
6 # - para: The number of model parameters to exclude from the degrees of freedom.
7 #
8 # The function:
9 # 1. Computes the autocorrelation function (ACF) up to m*s lags.
10 # 2. Calculates the test statistic as per the Pierce test formula.
11 # 3. Returns a list containing:
12 #   - `statistic`: The calculated test statistic.
13 #   - `p.value`: The p-value indicating the significance of the test statistic.
14 pierce.test <- function(x = NULL, m = NULL, s = NULL, para = NULL) {
15   if (is.null(x) || is.null(m) || is.null(s) || is.null(para)) {
16     stop("All parameters (x, m, s, para) must be provided")
17   }
18   x <- coredat(x) # Extract numeric data from zoo object
19   kk <- acf(x, m * s, plot = FALSE) # Calculate autocorrelation function
20   n <- length(x) # Number of observations
21   statistics <- 0
22   for (i in 1:m) {
23     statistics <- statistics + (1 / (n - (i * s))) * kk$acf[(i * s) + 1]^2
24   }
25   statistics <- statistics * n * (n + 2)
26   p.value <- 1 - pchisq(statistics, m - para)
27   return(list(statistic = statistics, p.value = p.value)) # Return values as a list
28 }

```

Ljung.Box2 function

```

1 # This function performs the Ljung-Box test to check residual autocorrelation.
2 # It accepts three parameters:
3 # - x: The time series data or residuals (as a zoo object) to analyze.
4 # - lag: The number of lags to include in the test.
5 # - param: The number of model parameters to adjust for.
6 #
7 # The function:
8 # 1. Computes the ACF up to the specified `lag` (ignoring the first autocorrelation).
9 # 2. Calculates the test statistic using the Ljung-Box formula.
10 # 3. Returns a list containing:
11 #   - `statistic`: The test statistic.
12 #   - `gdl`: The degrees of freedom (lag minus the number of parameters).
13 #   - `pvalue`: The p-value indicating the significance of the test.
14 Ljung.Box2 <- function(x, lag = 15, param = 1) {
15   if (NCOL(x) > 1) stop("x is not a vector or univariate time series")
16   x <- coredat(x) # Extract numeric data from zoo object
17   cor <- acf(x, lag.max = lag, plot = FALSE)
18   n <- length(x)
19   PARAMETER <- lag - param
20   obs <- cor$acf[2:(lag + 1)]
21   STATISTIC <- n * (n + 2) * sum((1 / seq(n - 1, n - lag)) * obs^2)
22   PVAL <- 1 - pchisq(STATISTIC, lag - param)
23   return(list(statistic = STATISTIC, gdl = PARAMETER, pvalue = PVAL))
24 }

```

Ljung.Box function

```

1 # The Ljung.Box function provides an extended Ljung-Box test across multiple lag lengths.
2 # It accepts four parameters:
3 # - x: The time series data or residuals (as a zoo object) to analyze.
4 # - maxlag: The maximum number of lags to include in the analysis.
5 # - par: The number of model parameters to adjust for.
6 # - all: If TRUE, computes the test across all lags up to `maxlag` and plots the p-values.
7 #
8 # The function:
9 # 1. Calls `Ljung.Box2` for each lag up to `maxlag`, excluding the specified `par` parameters.
10 # 2. Constructs a matrix containing the test statistic, degrees of freedom, and p-value at each lag.
11 # 3. Optionally plots the p-values as a bar chart if `all` is TRUE.
12 # 4. Returns a matrix containing the statistics for each lag.
13 Ljung.Box <- function(x, maxlag = 24, par = 0, all = TRUE) {
14   if (all == TRUE) {
15     LB <- matrix(0, nrow = maxlag, ncol = 3)
16     dimnames(LB) <- list(NULL, c("statistic", "gdl", "pvalue"))
17     for (i in (par + 1):maxlag) {
18       lbi <- Ljung.Box2(x, lag = i, param = par)
19       LB[i, ] <- c(round(lbi$statistic, 2), lbi$gdl, lbi$pvalue)
20     }
21     plot(LB[, 3], type = "h", lwd = 4, ylim = c(0, 1), ylab = "P-Value", xlab = "Lag")
22     abline(h = 0.05, lty = 2, col = "red")
23     abline(h = 0)
24   }
25   if (all == FALSE) LB <- Ljung.Box2(x, lag = maxlag, param = par)
26
27   return(LB)
28 }

```

stat.mod function

```

1 # The stat.mod function provides comprehensive model diagnostics for a fitted time series model.
2 # It takes three arguments:
3 # - fit.y: The fitted model object (e.g., ARIMA or other time series model).
4 # - maxlag: The maximum number of lags for the Ljung-Box test to check residual autocorrelation.
5 # - ordini: The number of model parameters.
6 #
7 # The function calculates and prints:
8 # 1. Coefficients, standard errors, t-statistics, and p-values of the model parameters.
9 # 2. Model diagnostics, such as the residual variance (sigma2), Akaike Information Criterion (AIC), and
10 #    log-likelihood.
11 # 3. Results from the Ljung-Box test to check for autocorrelation in the residuals up to `maxlag` lags.
12 # The results are formatted and displayed for clear interpretation of the model's fit and adequacy.
13 stat.mod=function(fit.y,maxlag=15,ordini=0)
14 {
15   Coef=fit.y$coef
16   Std.Err=sqrt(diag(fit.y$var.coef))
17   tstat=fit.y$coef/sqrt(diag(fit.y$var.coef))
18   pval=2*pnorm(abs(tstat),lower.tail=F)
19   tabella=data.frame(Coef,Std.Err,tstat,pval)
20   sigma2=fit.y$sigma2
21   names(sigma2)="sigma2"
22   AIC=fit.y$aic
23   names(AIC)="AIC"
24   loglik=fit.y$loglik
25   names(loglik)="loglik"
26   cc1=Ljung.Box(fit.y$resid,maxlag,ordini)
27   riga0="----- Parameters estimation -----"
28   names(riga0)=" "
29   riga="-----"
30   names(riga)=" "
31   print(riga0, quote=F)
32   print(tabella)
33   print(riga, quote=F)
34   riga2="----- Ljung-Box Test -----"
35   names(riga2)=" "
36   print(c(sigma2,AIC,loglik))
37   print(riga2, quote=F)
38   print(t(cc1))
39 }

```


Qtests function

```

1 # This function performs a series of Ljung-Box tests up to lag `k` for the given series.
2 # It accepts three arguments:
3 # - series: The time series data to analyze.
4 # - k: The maximum lag to test.
5 # - fitdf: The number of model parameters to adjust for.
6 #
7 # The function:
8 # 1. Applies the Ljung-Box test at each lag up to `k`.
9 # 2. Returns a matrix containing the lag and the corresponding p-value.
10 Qtests <- function(series, k, fitdf=0) {
11   pvals <- apply(matrix(1:k), 1, FUN=function(l) {
12     pval <- if (l<=fitdf) NA else Box.test(series, lag=l, type="Ljung-Box", fitdf=fitdf)$p.value
13     return(c("lag"=l,"pval"=pval))
14   })
15   return(t(pvals))
16 }

```

adfTest_valid function

```

1 # This function performs the Augmented Dickey-Fuller (ADF) test for stationarity, increasing lags until
  # no autocorrelation in residuals.
2 # It accepts three arguments:
3 # - series: The time series data to analyze.
4 # - kmax: The maximum number of lags to test.
5 # - type: The type of test ("c" for constant, "ct" for constant and trend, etc.).
6 #
7 # The function:
8 # 1. Performs the ADF test with increasing lags until residuals show no autocorrelation.
9 # 2. Prints whether the residuals are autocorrelated or not at each step.
10 # 3. Returns the final ADF test result.
11 adfTest_valid <-
12   function(series,kmax,type){ #ADF tests until no more autocorrelated residuals
13     k <- 0
14     noautocorr <- 0
15     while (noautocorr==0){
16       cat(paste0("ADF with ",k, " lags: residuals OK? "))
17       adf <- adfTest(series,lags=k,type=type)
18       pvals <- Qtests(adf@test$lm$residuals,24,fitdf=length(adf@test$lm$coefficients))[,2]
19       if (sum(pvals<0.05,na.rm=T) == 0) {
20         noautocorr <- 1; cat("OK \n")}
21       else cat("NO \n")
22       k <- k + 1
23     }
24     return(adf)
25   }pvals))
26 }

```

modelchoice function

```

1 # This function fits an ARIMA model and checks the significance of AR and MA terms and the residuals'
  # autocorrelation.
2 # It accepts four arguments:
3 # - p: The AR order.
4 # - q: The MA order.
5 # - data: The time series data.
6 # - k: The maximum number of lags for the Ljung-Box test.
7 #
8 # The function:
9 # 1. Tries to fit an ARIMA model with the given orders.
10 # 2. Checks the significance of AR and MA terms and residuals' autocorrelation.
11 # 3. Returns a list with the model's coefficients, standard errors, and Ljung-Box test results.
12 modelchoice <- function(p, q, data, k = 24) {
13   estim <- try(arima(data, order = c(p, 0, q), optim.control = list(maxit = 20000)), silent = TRUE)
14
15   if (class(estim) == "try-error") {
16     return(c("p" = p, "q" = q, "arsignif" = NA, "masignif" = NA, "resnocorr" = NA, "ok" = NA))
17   }

```

modelchoice function

```

18  coefs <- coef(estim)
19  se <- sqrt(diag(vcov(estim)))
20
21  arsignif <- if (p == 0) NA else {
22    ar_coefs <- coef(estim)[1:p]
23    ar_se <- sqrt(diag(vcov(estim)))[1:p]
24    sum(2 * (1 - pnorm(abs(ar_coefs / ar_se))) <= 0.05) == p
25  }
26
27  masignif <- if (q == 0) NA else {
28    ma_coefs <- coef(estim)[(p + 1):(p + q)]
29    ma_se <- sqrt(diag(vcov(estim)))[(p + 1):(p + q)]
30    sum(2 * (1 - pnorm(abs(ma_coefs / ma_se))) <= 0.05) == q
31  }
32
33  resnocorr <- sum(Box.test(estim$residuals, lag = k, type = "Ljung-Box")$p.value <= 0.05, na.rm = TRUE) == 0
34  checks <- c(arsignif, masignif, resnocorr)
35  ok <- as.numeric(sum(checks, na.rm = TRUE) == (3 - sum(is.na(checks))))
36
37  ljung_box <- Ljung.Box(estim$residuals, maxlag = k)
38
39  # Print coefficients and their standard errors
40  cat(paste0("ARMA(", p, ",", q, ") coefficients:\n"))
41  print(coefs)
42  cat("Standard errors:\n")
43  print(se)
44
45  # Print Ljung-Box test result
46  cat("Ljung-Box test result:\n")
47  print(ljung_box)
48
49  return(list(p = p, q = q, arsignif = arsignif, masignif = masignif, resnocorr = resnocorr, ok = ok,
50    coefficients = coefs, se = se, ljung_box = ljung_box))
51 }

```

armamodelchoice function

```

1  # This function finds the best ARMA model within given p and q ranges by checking each model's validity
2  # It accepts three arguments:
3  # - pmax: The maximum AR order to test.
4  # - qmax: The maximum MA order to test.
5  # - y_df: The differenced time series data.
6  #
7  # The function:
8  # 1. Tests all combinations of AR and MA orders up to pmax and qmax.
9  # 2. Stores results for each model's coefficients, standard errors, and diagnostic checks.
10 # 3. Returns a data frame with the comparison of all tested models.
11 armamodelchoice <- function(pmax, qmax, y_df) {
12   pqs <- expand.grid(p = 0:pmax, q = 0:qmax)
13   results <- list()
14
15   for (i in 1:nrow(pqs)) {
16     p <- pqs$p[i]
17     q <- pqs$q[i]
18     cat(paste0("Computing ARMA(", p, ",", q, ") \n"))
19     result <- modelchoice(p, q, coredat(y_df))
20     results[[i]] <- result
21   }
22
23   comparison <- do.call(rbind, results)
24   comparison <- as.data.frame(comparison, stringsAsFactors = FALSE)
25   comparison[, c("p", "q", "arsignif", "masignif", "resnocorr", "ok")] <-
26     lapply(comparison[, c("p", "q", "arsignif", "masignif", "resnocorr", "ok")], as.numeric)
27
28   return(comparison)
29 }

```

A.2 Codes for Part I

Codes for Part I

```

1  ## Set working directory to the location of the data file
2  # setwd("C:/Users/Marco Plazzogna/Desktop/MAGISTRALE/PRIMO ANNO/SECONDO SEMESTRE/TIME SERIES/Assignment
   ")
3
4  # Read data from a CSV file, skipping the first 4 lines and assuming ";" as the separator
5  data <- read.csv("valeurs_mensuelles.csv", sep = ";", skip = 4, header = FALSE)
6
7  # Extract the second column (assuming it contains the time series data)
8  serie <- data[, 2]
9
10 # Convert the series into a zoo object starting from January 2000 with monthly frequency
11 y <- zoo(serie, order.by = as.Date(seq(from = as.Date("2000-01-01"), by = "month", length.out = length(
   serie))))
12
13 # Plot the original time series
14 plot(y, main = "Original Time Series", ylab = "Value", xlab = "Time")
15
16 # Perform Augmented Dickey-Fuller test to check for stationarity
17 adf <- adfTest_valid(y,24,"ct")
18 print(adf)
19
20 adf <- adfTest_valid(y,24,"c")
21 print(adf)
22
23 # Perron-Phillips test to check for stationarity
24 pp.test(y)
25 y.pp = ur.pp(y, type = "Z-tau", model = "trend")
26 summary(y.pp)
27
28 # Optionally, add ACF and PACF plots for the series to inspect autocorrelations
29 acf(coredata(y), main = "ACF of Series")
30 pacf(coredata(y), main = "PACF of Series")
31
32 # Difference the series to achieve stationarity
33 y_df = diff(y)
34
35 # Plot the differenced series
36 plot(y_df, main = "Differenced Time Series", ylab = "Value", xlab = "Time")
37
38 # Plot ACF and PACF for the differenced series
39 acf(coredata(y_df), main = "ACF of Differenced Series")
40 pacf(coredata(y_df), main = "PACF of Differenced Series")
41
42 # Set maximum AR and MA orders for model selection
43 pmax = 5
44 qmax = 1

```

A.3 Codes for Part II

Codes for Part II

```

1  # Find the best ARMA models within the specified p and q ranges
2  armamodels <- armamodelchoice(pmax,qmax, y_df)
3
4  # Select models that are well-adjusted and valid
5  selec <- armamodels[armamodels[, "ok"] == 1 & !is.na(armamodels[, "ok"]), ] #models well adjusted and valid
6  selec[, c("p", "q", "arsignif", "massignif", "resnocorr", "ok")]
7
8  # Fit several ARMA models
9  models <- list()
10 aic_values <- numeric()
11 bic_values <- numeric()
12
13 # Find the best model with AIC and BIC
14 for (i in 1:nrow(selec)) {
15   p <- selec$p[i]
16   q <- selec$q[i]
17   fit <- tryCatch({
18     Arima(coredata(y_df), order = c(p, 0, q))
19   }, error = function(e) NULL)

```

Codes for Part II

```

20   if (!is.null(fit)) {
21     models[[paste0("ARMA(", p, ",", q, ")")] <- fit
22     aic_values <- c(aic_values, fit$aic)
23     bic_values <- c(bic_values, BIC(fit))
24   }
25 }
26
27 # Create a data frame to compare models
28 comparison <- data.frame(
29   Model = names(models),
30   AIC = aic_values,
31   BIC = bic_values
32 )
33
34 # Print the model comparison table
35 print(comparison)
36
37 # Find the best models based on AIC and BIC
38 best_aic_model_index <- which.min(comparison$AIC)
39 best_bic_model_index <- which.min(comparison$BIC)
40
41 best_aic_model_name <- names(models)[best_aic_model_index]
42 best_bic_model_name <- names(models)[best_bic_model_index]
43
44 best_aic_model <- models[[best_aic_model_name]]
45 best_bic_model <- models[[best_bic_model_name]]
46
47 # Display summaries of both models
48 cat("Model with the best AIC:\n")
49 summary(best_aic_model)
50
51 cat("\nModel with the best BIC:\n")
52 summary(best_bic_model)
53
54 best_model = best_aic_model
55
56 # Check residuals of the best AIC model
57 cat("\nResidual diagnostics for the best AIC model:\n")
58 checkresiduals(best_model)
59
60 # Perform Ljung-Box test on residuals of the best AIC model
61 cat("\nLjung-Box test for residuals (best AIC model):\n")
62 Ljung.Box(best_model$residuals, maxlag = 24, par=1)

```

A.4 Codes for Part III

Codes for Part III

```

1   h <- 2
2
3   # Generate forecasts
4   forecasts <- forecast(best_model, h = h, level = 0.95)
5
6   # Extract the last 30 observations from the original series
7   last_30_index <- tail(index(forecasts$x), 30)
8   last_30_data <- window(forecasts$x, start = last_30_index[1])
9
10  # Extract forecast values and confidence intervals
11  forecast_values <- forecasts$mean
12  lower_bound <- forecasts$lower # 95% lower bound
13  upper_bound <- forecasts$upper # 95% upper bound
14
15  # Assuming you have the index of the last fitted value and the first forecast value
16  last_fitted_time <- tail(time(forecasts$fitted), 1)
17  first_forecast_time <- head(time(forecasts$mean), 1)
18
19  last_fitted_value <- tail(forecasts$fitted, 1)
20  first_forecast_value <- head(forecasts$mean, 1)
21
22  # Plot the last 30 observations with forecasts and confidence intervals
23  plot(last_30_data, main = "Forecast with 95% Confidence Intervals", xlab = "Time", ylab = "Value", type
      = "l", xlim=c(260,290))
24  lines(forecasts$fitted, col = "blue")
25  points(last_30_data, col="black")

```

Codes for Part III

```
26 points(forecasts$fitted, col="blue")
27 points(time(forecasts$mean), forecasts$mean, col = "red")
28 arrows(time(forecasts$mean), lower_bound, time(forecasts$mean), upper_bound, code = 3, angle = 90,
29         length = 0.1, col = "red", lty=3)
30 lines(forecasts$mean, col="red")
31
32 # Draw a link between the last blue point and the first red point
33 segments(last_fitted_time, last_fitted_value, first_forecast_time, first_forecast_value, col="red")
34
35 # Add legend
36 legend("bottomleft", legend = c("Series", "Fitted", "Forecast", "95% CI"), col = c("black", "blue", "red", "red"), lty = c(1, 1, 1, 3), pch = c(NA, NA, 1, NA), bty = "o", cex = 0.7)
```