

République Algérienne Démocratique et Populaire
Ministère de l'Enseignement Supérieur et de la Recherche
Scientifique
UNIVERSITÉ KASDI MERBAH OUARGLA
FACULTÉ DE MATHÉMATIQUES ET SCIENCES DE LA
MATIÈRE
DÉPARTEMENT DE MATHÉMATIQUES



Mémoire présenté En Vue De L'obtention Du

DIPLÔME DE MASTER

EN MATHÉMATIQUES

Option : Probabilité et Statistique

Par

Chekima Romayssa

Intitulé

Modélisation des séries chronologiques

Membres du jury

AMARA Abdelkader	M. C. A	UKMO	Président
MEDDI Fatima	M. C. B	UKMO	Examineur
ARBIA Hanane	M. C. B	UKMO	Rapporteur

2021-2022

Table des matières

Dédicace	iv
Remerciements	v
Table des figures	vii
Notations et abriviations	viii
Introduction	1
1 Généralités sur les séries chronologiques	3
1.1 Définition d'une série chronologique	3
1.2 Composantes d'une série chronologique	4
1.2.1 Tandanse (T_t)	4
1.2.2 Saisonnalité(S_t)	5
1.2.3 Résidus ou les variations accidentelles (ε_t) . . .	5
1.2.4 Cycle (C_t)	5
1.3 Les schémas de composition	6
1.4 Notion d'opérateur retard L et de différenciation Δ . .	7
1.4.1 Opérateur de retard	7
1.4.2 Opérateurs de différenciation	9
1.4.3 Propriétés de l'opérateur retard	9

1.5	Processus stochastique	9
1.5.1	Processus bruit blanc	10
1.6	La Stationnarité	10
1.7	La Non Stationnarité	10
1.7.1	Non stationnarité déterministe	11
1.7.2	Non stationnarité stochastique	11
1.8	Autocorrélations simple et partielle	12
1.8.1	La fonction d'autocovariance	12
1.8.2	La fonction d'autocorrélation	13
1.8.3	La fonction d'autocorrélation partielle	13
2	Les modèles des séries chronologiques	16
2.1	Les modèles autorégressifs	16
2.1.1	Représentation stationnaire	17
2.1.2	Fonction d'autocorrélation d'un $AR(p)$	17
2.1.3	Fonction d'autocorrélation partielle d'un $AR(p)$	18
2.1.4	Estimation	19
2.2	Les modèles moyennes mobiles, MA (Moving Average)	21
2.2.1	Autocorrélations d'un $MA(q)$	21
2.2.2	Représentation stationnaire	22
2.3	Les modèles $ARMA(p,q)$	23
2.3.1	La stationnarité	23
2.3.2	Autocorrélation	23
2.4	Les modèles $ARIMA(p,d,q)$	25
2.5	Méthode de Box & Jenkins	26
2.5.1	Identification du modèle	26
2.5.2	Estimation des paramètres	27
2.5.3	Validation du modèle	30
2.5.4	La prévision	33

3	Application : Modélisation par la méthode de Box-Jenkins	34
3.1	Etude de la stationnarité	35
3.1.1	Stationnarisation de la série	40
3.1.2	Modélisation de la série	46
3.2	Estimation	48
3.3	Validation du modèle	48
3.3.1	Test de significativité des paramètres	49
3.3.2	Test de normalité des résidus	49
3.4	La prévision	51
	Conclusion	52
	Bibliographie	53

Dédicace

Je dédie ce modeste travail

A ma très cher Mère et mon très cher Père

A mon cher mari

A tous mes frères

A toutes la familles : Chekima

A ma deuxième chère famille : Kadri

A ceux qui m'ont toujours encouragé pour que je réussisse dans mes études

A tout ce qui est m'encouragé dans la réalisation de ce modeste travail, A mes amis

A mon encadreur ARBIA Hanane, en espérant qu'elle trouve dans ce travail le témoignage de ma profonde gratitude

A tout mes enseignants, A Tous qui m'on aider de proche ou de loin.

Remerciements

Tout d'abord, je remercie ALLAH Le Tout Puissant qui m'a accordée la volonté et le courage pour réaliser ce mémoire.

Mes plus vifs remerciements et ma profonde gratitude vont à Dr.Arbia Hanane qui a accepté de diriger ce travail. Sa grande disponibilité et ses encouragements ont joué un rôle important dans la réalisation de ce mémoire.

Je remercie sincèrement, les membres de jury d'avoir bien voulu accepter de faire partie de la commission d'examineur.

Enfin, je ne peux terminer cette partie sans exprimer notre gratitude à mes parents, mon mari, mes frères et mes grands parents, et ma gentille deuxième famille qui j'ont toujours soutenue, encouragé et stimulée pendant mes études.

Merci à tous et à toutes.

Liste des tableaux

3.1	L'évolution annuelle des cas de paludisme.	35
3.2	Comparaison entre les critères des modèles.	48
3.3	Prévision du cas de paludisme.	51

Table des figures

3.1	Représentation graphique de la série brute.	36
3.2	Représentation graphique de la série stationnaire. . . .	46
3.3	Normalité des résidus du modèle AR(1).	50

Notations et abriviations

ADF	Test de Dickey Fuller Augmenté
AIC	Critère d'information d'Akaike
AR	Autorégressif
ARIMA	Autorégressif moyenne mobile intégré
ARMA	Autorégressif moyenne mobile
BB	Bruit blanc
BIC	Critère de Schwarz
DS	Non stationnarité stochastique
EMV	Estimateur du maximum de vraisemblance
JB	Test de Jarque & Bira
MA	Moyennes Mobiles (Moving Average).
MAE	Erreur absolue moyenne
MAPE	Ecart absolu moyen en pourcentage

MSE	Erreur quadratique moyenne
RMSE	Racine carrée de l'erreur quadratique moyenne
TS	Non stationnarité déterministe

Introduction

L'analyse des séries chronologiques est un outil couramment utilisé de nos jours pour mieux comprendre et décrire l'évolution d'un phénomène présenté par la série. L'idée est de prendre un échantillon de données et de construire un modèle qui correspond le mieux à ces données. Ce modèle nous permet de tirer certaines conclusions sur la série, pour la prédiction de données futures.

La théorie des séries chronologiques trouve des applications dans de nombreux domaines : économétrie (évolution d'indicateurs boursiers, des productions agricoles ou industrielles, des prix), Assurance (analyse des sinistres), Médecine, Biologie, finance, astronomie...etc. Cette dernière a connue un grand développement depuis la publication du livre de *Box-Jenkins* [2] dans lequel les principales propriétés des processus stationnaires autorégressif moyenne mobile (ARMA) sont décrites avec les méthodes d'identification, d'estimation, validation et de prévision.

Ce mémoire s'organise en trois chapitres :

Le premier chapitre est consacré pour les préliminaires ou nous mentionnons quelques notions et définitions sur les séries chronologiques et leurs caractéristiques.

Ensuite, le deuxième chapitre présente différents modèles : *AR*, *MA*, *ARMA* et *ARIMA*, leurs propriétés (la stationnarité, les fonctions d'autocorrélation simples et partielles,...), et la méthode de Box & Jenkins.

Finalement, le troisième chapitre est consacré à une application de la méthodologie de Box-Jenkins avec le programme Eviews sur des données réelles (le nombre des cas de paludisme tirées de direction de la santé et de la population de Ouargla (2000-2021)) pour obtenir le meilleur modèle qui représente cette série afin de prévoir.

Chapitre 1

Généralités sur les séries chronologiques

La théorie des séries chronologiques est appliquée dans des domaines aussi variés que démographie, finance, économétrie, Energie, Médecine, Traitement de signal.

Dans ce chapitre, plusieurs concepts importants liés à l'analyse des séries chronologiques seront abordés. Parmi ceux-ci, on retrouve les notions de stationnarité et d'autocorrélation et bruit blanc.

1.1 Définition d'une série chronologique

Une série chronologique ou chronique est définie par : [1], [3]

Définition 1.1 *est un ensemble des observations d'une variable statistique d'un même phénomène, ordonnés dans le temps (année, trimestre, mois, jour,...).*

1.2 Composantes d'une série chronologique

Une série X_t est composée des quatre parties comme suit : [3]

1.2.1 Tandanse (T_t)

Il s'agit d'un terme de la série qui traduit l'évolution à moyen terme du phénomène. Elle sera estimée sous forme paramétrique (linéaire, polynomiale, logarithmique, exponentielle,... etc)

Tendance lineaire

La tendance la plus simple est linéaire on peut estimer les paramètres au moyen de la méthode des moindre carés c'est une régression simple :

$$T_t = a + bt$$

Tendance quadratique

La forme de la tendance quadratique

$$T_t = a + bt + ct^2$$

Tendance polynomiale d'ordre q

On peut ajuster la série par un polynôme d'ordre q. Les paramètres peut être estimer au moyen de la méthode de moindre carées c'est une régression avec q variables explicatives :

$$T_t = b_0 + b_1t + b_2t^2 + \dots + b_qt^q$$

Tendance logistique

La fonction logistique permet de modéliser des processus ne pouvant dépasser une certaine valeur c :

$$T_t = \frac{c}{1 + be^{-at}}$$

1.2.2 Saisonnalité(S_t)

Elle représente des effets périodiques de période connue p qui se reproduisent de façon plus ou moins identique d'une période à l'autre elle est notée par S_t $t = 1, \dots, T$.

Elle est généralement supposée périodique :

$$S_{t+p} = S_t$$

d'une période p .

1.2.3 Résidus ou les variations accidentelles (ε_t)

Les résidus noté ε_t est la partie non structurée du phénomène. Elle est modélisée par une suite de variables aléatoires ε_t , centrées, non corrélées et de même variance.

1.2.4 Cycle (C_t)

C'est une succession de mouvements persistant des variations de mouvement ascendant (période de prospérité) et de mouvement descendant (période de dépression). Pour de longues séries, un mouvement cyclique peut se superposer à la tendance.

1.3 Les schémas de composition

Pour pouvoir séparer les composantes servant à décrire la série observée, il est nécessaire de préciser leur mode d'interaction. Il existe trois types de schémas ou de modèles : [9]

1-Schéma additif

Dans ce modèle, la série chronologique, s'écrit de la façon suivante :

$$X_t = Z_t + S_t + C_t + \varepsilon_t, \quad t = 1 \dots T,$$

La tendance, la Saisonnalité, les variations accidentelles ou Résidus et cyclique,

ont un effet additif dans ce modèle de série chronologique.

2-Schéma multiplicatif

Dans ce modèle, la série chronologique, s'écrit de la façon suivante :

$$X_t = Z_t \times S_t \times C_t \times \varepsilon_t, \quad t = 1 \dots T.$$

Dans ce modèle, la tendance, la Saisonnalité, les variations accidentelles ou Résidus et cyclique, ont un effet multiplicatif dans ce modèle de série chronologique.

3-Schéma mixte

Dans ce modèle les deux opération addition et multiplication sont utilisées. On donne par exemple :

$$X_t = (Z_t + C_t)S_t + \varepsilon_t, \quad t = 1 \dots T,$$

1.4 Notion d'opérateur retard L et de différenciation Δ

La manipulation pratique ou théorique des séries chronologiques est effectuée à l'aide de l'opérateur de retard et de différenciation Δ . [3] , [7]

1.4.1 Opérateur de retard

On appellera opérateur retard L (L =lag, ou B =backward) l'opérateur linéaire défini par

$$L : X_t \longmapsto L(X_t) = LX_t = X_{t-1}$$

et opérateur avance F (F =forward)

$$F : X_t \longmapsto F(X_t) = FX_t = X_{t+1}$$

Remarque

$$L \circ F = F \circ L = I$$

(opérateur identité) et on notera par la suite

$$F = L^{-1} \text{ et } L = F^{-1}$$

Polynômes d'opérateurs L (i) Il est possible de composer les opérateurs :

$$L^2 = L \circ L$$

et plus généralement

$$L^p = \underbrace{L \circ L \circ \dots \circ L}_{P \text{ fois}} \quad \text{où } p \in \mathbb{N}$$

avec la convention $L^0 = I$. On notera que

$$L^p(X_t) = X_{t-p}$$

(ii) Soit A le polynôme

$$A(z) = a_0 + a_1z + a_2z^2 + \dots + a_pz^p$$

On notera $A(L)$ l'opérateur

$$A(L) = a_0I + a_1L + a_2L^2 + \dots + a_pL^p = \sum_{k=0}^p a_kL^k$$

Soit (X_t) une série temporelle. La série (Y_t) définie par $Y_t = A(L)X_t$ vérifie

$$Y_t = A(L)X_t = \sum_{k=0}^p a_kX_{t-k}$$

1.4.2 Opérateurs de différenciation

Définition 1.2 On définit l'opérateur de différenciation Δ par :

$$\Delta X_t = X_t - X_{t-1} , \forall t \geq 2$$

Définition 1.3 On définit l'opérateur linéaire de différenciation d'ordre k par la formule de récurrence :

$$\Delta^{(k)} X_t = \Delta(\Delta^{(k-1)} X_t)$$

1.4.3 Propriétés de l'opérateur retard

1. $L^0 X_t = X_t$
2. $La = a$, l'opérateur d'une constante a est une constante.
3. $(L^i + L^j)X_t = L^i X_t + L^j X_t = X_{t-i} + X_{t-j}$.
4. $L^i(L^j X_t) = L^i X_{t-j} = X_{t-i-j}$, de même $L^i(L^j X_t) = L^{i+j} X_t = X_{t-i-j}$.
5. $L^{-i} X_t = X_{t+i}$.

1.5 Processus stochastique

On appelle processus stochastique ou processus aléatoire toute famille de variables aléatoires X_t . Cela signifie qu'à tout $t \in T$ est associée une variable aléatoire prenant ses valeurs dans un ensemble numérique E . On note le processus X_t . Si T est dénombrable, On dit que le processus est discret. [8]

1.5.1 Processus bruit blanc

Le processus ε_t est un bruit blanc si :

- i) $E(\varepsilon_t) = 0$ pour tout t .
- ii) $var(\varepsilon_t) = E(\varepsilon_t^2) = \sigma^2$ pour tout t .
- iii) $cov(\varepsilon_t, \varepsilon_{t-h}) = E(\varepsilon_t \varepsilon_{t-h}) = 0$ pour tout t et $h \neq 0$.

1.6 La Stationnarité

Une des grandes questions dans l'étude de séries chronologiques est de savoir si celles-ci suivent un processus stationnaire. [3], [7]

Définition 1.4 Soit un processus temporel à valeurs réelles et en temps discret Z_1, Z_2, \dots, Z_t . Il est dit stationnaire au sens fort si pour toute fonction f mesurable :

$$f(Z_1, Z_2, \dots, Z_t) \text{ et } f(Z_{1+k}, Z_{2+k}, \dots, Z_{t+k})$$

ont même loi.

Définition 1.5 Soit un processus temporel à valeurs réelles et en temps discret Z_1, Z_2, \dots, Z_t . Il est dit stationnaire au sens faible (ou de second ordre) si

- $E(Z_i) = \mu$ (ne dépend pas de t) $\forall i = 1 \dots t$.
- $Var(Z_i) = \delta \neq \infty$ $\forall i = 1 \dots t$.
- $Cov(Z_i, Z_{i-k}) = \rho_k$ (ne dépend pas de t) $\forall i = 1 \dots t, \forall k = 1 \dots t$.

1.7 La Non Stationnarité

La plupart des séries sont non stationnaires, c'est-à-dire que le processus qui les décrit ne vérifie pas au moins une des conditions de

la définition d'un processus stationnaire du second ordre. Ceci nous conduit à définir deux types de non stationnarité : non stationnarité déterministe et non stationnarité stochastique. [5]

1.7.1 Non stationnarité déterministe

On dit que le processus Y_t est caractérisé par une non stationnarité déterministe, ou encore que le processus Y_t est *TS* (Trend stationary) s'il peut s'écrire :

$$Y_t = f(t) + Z_t$$

ou $f(t)$ est une fonction qui dépend du temps et Z_t est un processus stationnaire.

Ainsi, ce processus est rendu stationnaire en lui enlevant sa tendance déterministe :

$$Y_t - f(t) = Z_t$$

stationnaire.

1.7.2 Non stationnarité stochastique

On dit que le processus Y_t est caractérisé par une non stationnarité stochastique, ou encore que le processus Y_t est *DS* (Difference stationary) si le processus différencié une fois $(1 - L)Y_t$ est stationnaire. On parle aussi de processus intégré d'ordre 1, on note $Y_t \sim I(1)$:

De manière générale, on dit que le processus Y_t est un processus intégré d'ordre d , avec d le degré d'intégration, si le processus différencié d fois $(1 - L)^d Y_t$ est stationnaire. On note $Y_t \sim I(d)$:

$$(1 - L)^d Y_t = Z_t$$

1.8 Autocorrélations simple et partielle

Les principales caractéristiques temporelles d'un processus sont données par l'autocorrélation (simple) et l'autocorrélation partielle [3].

1.8.1 La fonction d'autocovariance

La fonction d'autocovariance $\gamma(h)$ mesure la covariance entre une variable et cette même variable à des dates différentes, pour un délai h :

$$\gamma(h) = Cov(Y_t, Y_{t-h}) = E[(Y_t - E(Y_t))(Y_{t-h} - E(Y_{t-h}))]$$

Ainsi

$$\begin{aligned}\gamma(0) &= Var(Y_t) \\ &= E[(Y_t - E(Y_t))^2]\end{aligned}$$

Elle fournit une information sur la variabilité de la série et sur les liaisons temporelles qui existent entre les diverses composantes de la série Y_t .

Remarque 1.1 *La fonction d'autocovariance d'un processus stationnaire est une fonction paire :*

$$\gamma(-h) = \gamma(h) \quad \forall h$$

1.8.2 La fonction d'autocorrélation

La fonction d'autocorrélation est définie par :

$$\rho(h) = \frac{\gamma(h)}{\gamma(0)}, \quad h \in \mathbb{Z}$$

avec $\rho(0) = 1$ et $|\rho(h)| < 1$.

On appelle coefficient d'autocorrélation d'ordre 1 (resp. d'ordre k) le coefficient de corrélation linéaire $\rho(1)$ (resp. $\rho(k)$) calculé entre la série et cette série décalée d'une période (resp. k périodes).

On définit la matrice de corrélation (de dimension m) de la manière suivante :

$$R(m) = \begin{pmatrix} 1 & \rho_1 & \rho_2 & \cdot & \cdot & \cdot & \rho(m-1) \\ \rho_1 & 1 & \rho_1 & \cdot & \cdot & \cdot & \rho(m-2) \\ & & \cdot & & & & \\ & & \cdot & & & & \\ & & \cdot & & & & \\ \rho(m-1) & \rho(m-2) & \cdot & \cdot & \cdot & \rho_1 & 1 \end{pmatrix}$$

1.8.3 La fonction d'autocorrélation partielle

Elle mesure la liaison (linéaire) entre Y_t et Y_{t-h} une fois retirés les liens transitant par les variables intermédiaires $Y_{t-1}, \dots, Y_{t-h+1}$.

Le coefficient d'autocorrélation partielle d'ordre h , noté $r(h)$, est le coefficient de corrélation entre :

- $Y_t - E(Y_t/Y_{t-1}, \dots, Y_{t-h+1})$ et
- $Y_{t-h} - E(Y_{t-h}/Y_{t-1}, \dots, Y_{t-h+1})$

On a donc :

$$r(h) = \text{cov}(Y_t, Y_{t-h}/Y_{t-1}, \dots, Y_{t-h+1}).$$

C'est donc le coefficient de Y_{t-h} dans la régression de Y_t sur $Y_{t-1}, \dots, Y_{t-h+1}, Y_{t-h}$.

Le coefficient d'autocorrélation partielle d'ordre h d'un processus stationnaire se calcule de la manière suivante :

$$r(h) = \frac{|R(h)^*|}{|R(h)|}$$

avec

$$R(h) = \begin{pmatrix} 1 & \rho(1) & \cdot & \cdot & \cdot & \rho(h-1) \\ \rho(1) & 1 & \cdot & \cdot & \cdot & \rho(h-2) \\ \cdot & & & & & \cdot \\ \cdot & & & & & \cdot \\ \cdot & & & & & \cdot \\ \rho(h-1) & \rho(h-2) & \cdot & \cdot & \cdot & 1 \end{pmatrix}$$

et $R(h)^*$ la matrice $R(h)$ dans laquelle on a remplacé la colonne h

par $\begin{pmatrix} \rho(1) \\ \rho(2) \\ \cdot \\ \cdot \\ \cdot \\ \rho(h) \end{pmatrix}$, soit :

$$R(h)^* = \begin{pmatrix} 1 & \rho(1) & \cdot & \cdot & \cdot & \rho(1) \\ \rho(1) & 1 & \cdot & \cdot & \cdot & \rho(2) \\ \cdot & & & & & \cdot \\ \cdot & & & & & \cdot \\ \cdot & & & & & \cdot \\ \rho(h-1) & \rho(h-2) & \cdot & \cdot & \cdot & \rho(h) \end{pmatrix}$$

Ainci,

$$r(1) = \rho(1) \quad r(2) = \frac{\rho(2) - \rho(1)^2}{1 - \rho(1)^2} \dots$$

Chapitre 2

Les modèles des séries chronologiques

Dans ce chapitre, nous discutons de certains modèles linéaires fréquemment utilisés pour une série temporelle : le modèle autorégressif (*AR*) et le modèle à moyenne mobile (*MA*). Nous pour suivons avec le modèle *ARMA* qui permet de combiner ces deux modèles et le modèle *ARIMA*. Nous discutons également d'une méthode de prévision des séries temporelles, telle que la méthode de Box et Jenkins.

2.1 Les modèles autorégressifs

Définition 2.1 *Un processus autorégressif d'ordre p , noté $AR(p)$ est donné par : [10], [4], [6]*

$$X_t = c + \varphi_1 X_{t-1} + \varphi_2 X_{t-2} + \dots + \varphi_p X_{t-p} + \varepsilon_t$$

$\varphi_1, \dots, \varphi_p$ sont les paramètres du modèle, c est une constante, ε_t un bruit blanc.

En utilisant L l'opérateur des retards, on peut écrire le modèle :

$$(1 - \varphi_1 L - \varphi_2 L^2 - \dots - \varphi_P L^P) X_t = c + \varepsilon_t.$$

2.1.1 Représentation stationnaire

Ce processus est actuellement défini sous une forme implicite et en particulier il n'est pas certain que cette dernière équation admette toujours une solution stationnaire.

Si le polynôme Φ a des racines de module différent de 1, alors on peut inverser l'opérateur $\Phi(B)$. Nous concluons que l'équation accepte une solution unique, en écrivant :

$$X_t = \Phi(L)^{-1} \varepsilon_t = \sum_{i=-\infty}^{+\infty} h_i \varepsilon_{t-i}.$$

Nous pouvons alors montrer que l'on a $\sum_{i=-\infty}^{+\infty} |h_i| < +\infty$ donc la représentation est stationnaire.

2.1.2 Fonction d'autocorrélation d'un $AR(p)$

Soit un processus stationnaire $AR(p)$ défini par :

$$\Phi(L)X_t = X_t + \sum_{i=1}^p \phi_i X_{t-i} = \varepsilon_t,$$

dont les racines sont de module supérieur 1.

Alors la fonction d'autocorrélation :

$$\rho(h) + \sum_{j=1}^p \phi_j \rho(h-j) = 0 \quad h = 1, 2, \dots$$

est donc définie comme la solution d'une suite récurrente d'ordre p , donc racines non-nulles de l'équation :

$$\lambda^p + \sum_{j=1}^p \phi_j \lambda^{p-j} = 0.$$

Remarquons également qu'on a la relation matricielle (équation de Yule-Walker) :

$$\begin{pmatrix} -\phi_1 \\ \cdot \\ \cdot \\ \cdot \\ -\phi_p \end{pmatrix} = [\rho(i-j)]^{-1} \begin{pmatrix} \rho(1) \\ \cdot \\ \cdot \\ \cdot \\ \rho(p) \end{pmatrix}$$

Ce qui donne un algorithme d'estimation des coefficients d'un processus $AR(p)$ en remplaçant $\rho(h)$ par leur autocorrélation empirique. Remarquons également que la variance de ce processus vaut :

$$\gamma(0) = -\sum_{j=1}^p \phi_j \gamma(j) + \sigma^2 = \frac{\sigma^2}{1 + \sum_{j=1}^p \phi_j \rho(j)}.$$

2.1.3 Fonction d'autocorrélation partielle d'un $AR(p)$

Soit $X_t \sim AR(p)$:

$$\Phi(L)X_t = X_t + \sum_{i=1}^p \phi_i X_{t-i} = \varepsilon_t,$$

dont les racines sont à l'extérieur du disque unité.

On a que l'autocorrélation partielle $r(h) = 0$ si $h \geq p + 1$ car la projection de X_t sur $M(X_{t-1}, \dots, X_{t-h})$ est $-\sum_{j=1}^p \phi_j X_{t-j}$ et le coefficient associé à X_{t-h} est nul.

Cette propriété est très utile en pratique lorsque l'on cherche à identifier l'ordre d'un processus *AR*. On peut ainsi calculer les autocorrélation partielles empiriques et regarder quand celles-ci sont négligeables (non-significativement différentes de 0).

2.1.4 Estimation

Les équations de Yule-Walker sont utilisées pour estimer les coefficients. Les coefficients sont exprimés en fonction des auto-corrélations, et les paramètres estimés sont trouvés en fonction des auto-corrélations estimées.

On écrit l'équation aux différences :

$$\rho(h) = \phi_1 \rho(h-1) + \dots + \phi_p \rho(h-p).$$

sous forme matricielle :

$$\left\{ \begin{array}{l} \phi_1\rho(0) + \phi_2\rho(1) + \dots + \phi_p\rho(p-1) = \rho(1), h = 1 \\ \phi_1\rho(1) + \phi_2\rho(0) + \dots + \phi_p\rho(p-2) = \rho(2), h = 2 \\ \vdots \\ \vdots \\ \phi_1\rho(p-1) + \phi_2\rho(p-2) + \dots + \phi_p\rho(0) = \rho(p), h = p \end{array} \right.$$

Les équations de Yule Walker s'écrivent

$$R_p\phi = p.$$

Alors

$$\widehat{\phi} = \widehat{R}_p^{-1}\widehat{\rho},$$

avec

$$\widehat{R}_p = \begin{pmatrix} 1 & \widehat{\rho}(1) & \widehat{\rho}(2) & \dots & \widehat{\rho}(p-1) \\ \widehat{\rho}(1) & 1 & \widehat{\rho}(1) & \dots & \widehat{\rho}(p-2) \\ \widehat{\rho}(2) & \widehat{\rho}(1) & 1 & \dots & \widehat{\rho}(p-3) \\ \vdots & \vdots & \vdots & \dots & \vdots \\ \vdots & \vdots & \vdots & \dots & \vdots \\ \widehat{\rho}(p-1) & \widehat{\rho}(p-2) & \widehat{\rho}(p-3) & \dots & 1 \end{pmatrix},$$

$$\widehat{\rho} = \begin{pmatrix} \widehat{\rho}(1) \\ \widehat{\rho}(2) \\ \vdots \\ \vdots \\ \widehat{\rho}(p) \end{pmatrix}, \text{ et } \widehat{\phi} = \begin{pmatrix} \widehat{\phi}_1 \\ \widehat{\phi}_2 \\ \vdots \\ \widehat{\phi}_p \end{pmatrix}.$$

2.2 Les modèles moyennes mobiles, MA (Moving Average)

Définition 2.2 On appelle processus moyenne mobile d'ordre q , noté $MA(q)$ pour Moving Average, un processus $\{Y_t\}_{t \in \mathbb{Z}}$ défini par : [3]

$$Y_t = \varepsilon_t - \sum_{i=1}^q \theta_i \varepsilon_{t-i} \quad , \forall t \in \mathbb{Z}$$

θ_i sont des réels, $\theta_q \neq 0$, et $\varepsilon_t \sim BB(0, \sigma_\varepsilon^2)$.

Cette relation équivalent à :

$$\begin{aligned} Y_t &= (1 - \theta_1 B - \dots - \theta_q B^q) \varepsilon_t \\ \iff Y_t &= \Theta(B) \varepsilon_t. \end{aligned}$$

Proposition 2.1 un modèle moyenne mobile d'ordre q est inversible si les racines de

$$1 - \theta_1 Z - \theta_2 Z^2 - \dots - \theta_q Z^q = 0,$$

sont en module strictement supérieures à 1.

2.2.1 Autocorrélations d'un MA(q)

La fonction d'autocovariance d'un $MA(q)$ est :

$$\gamma(h) = E[X_t X_{t+h}].$$

Ce qui est plus adéquat ici, car il y a non-corrélation des ε_t avec le futur.

On a

$$X_t = \varepsilon_t - \theta_1 \varepsilon_{t-1} - \theta_2 \varepsilon_{t-2} - \dots - \theta_q \varepsilon_{t-q},$$

$$X_{t+h} = \varepsilon_{t+h} - \theta_1 \varepsilon_{t+h-1} - \theta_2 \varepsilon_{t+h-2} - \dots - \theta_q \varepsilon_{t+h-q}.$$

donc

$$\gamma(h) = \begin{cases} (-\theta_h + \theta_1 \theta_{h+1} + \dots + \theta_q \theta_{q+1}) \sigma_\varepsilon^2 & \text{si } 1 \leq h \leq q \\ (1 + \theta_1^2 + \theta_2^2 + \dots + \theta_q^2) \sigma_\varepsilon^2 & \text{si } h = 0 \\ 0 & \text{si } h > q \end{cases}$$

La fonction d'autocorrélation :

$$\rho(h) = \begin{cases} 1 & \text{si } h = 0 \\ \frac{-\theta_h + \theta_1 \theta_{h+1} + \dots + \theta_q \theta_{q+1}}{1 + \theta_1^2 + \theta_2^2 + \dots + \theta_q^2} & \text{si } 1 \leq h \leq q \\ 0 & \text{si } h > q \end{cases}$$

2.2.2 Représentation stationnaire

La définition d'un $MA(q)$ est explicite et ne pose donc pas de problème : le processus Y_t est parfaitement défini et est automatiquement stationnaire.

2.3 Les modèles ARMA(p,q)

Définition 2.3 Un modèle est dit autorégressif et moyenne mobile d'ordre p et q , noté $ARMA(p, q)$ si : [10], [5]

$$X_t - \sum_{j=1}^p \alpha_j X_{t-j} = c + \varepsilon_t + \sum_{k=1}^q \beta_k \varepsilon_{t-k} \quad (2.1)$$

$\varepsilon_t \sim BB(0, \sigma_\varepsilon^2)$, $\alpha_p \neq 0$, $\beta_q \neq 0$, c est une constante.

L'équation (2.1) est équivalente à l'écriture :

$$\Phi(L)X_t = c + \Theta(L)\varepsilon_t,$$

ou

$$\Phi(L) = 1 - \alpha_1 L - \alpha_2 L^2 - \dots - \alpha_p L^p,$$

et

$$\Theta_q(L) = 1 + \beta_1 L + \beta_2 L^2 + \dots + \beta_q L^q.$$

Les polynômes Φ et Θ n'ont pas de racines communes. On considère que Φ et Θ ont toutes les racines de module strictement supérieur à 1.

2.3.1 La stationnarité

Les conditions de stationnarité d'un processus $ARMA$ si toutes les racines de $\Phi(Z)$ en module supérieur à 1.

2.3.2 Autocorrélation

Soit (X_t) un processus $ARMA(p; q)$, alors l'autocovariance $\gamma(k)$ satisfait :

$$\begin{aligned}\gamma(k) &= E(X_t X_{t-k}) \\ &= E \left[\left(\sum_{j=1}^p \phi_j X_{t-j} + \varepsilon_t - \sum_{i=1}^q \theta_i \varepsilon_{t-i} \right) X_{t-k} \right].\end{aligned}$$

La covariance croisée entre X_{t-k} et ε_t est :

$$\gamma_{X_\varepsilon}(k) = E(\varepsilon_t X_{t-k}).$$

On a :

$$\begin{aligned}\gamma(k) &= \sum_{j=1}^p \phi_j E(X_{t-j} X_{t-k}) - \sum_{i=1}^q \theta_i (\varepsilon_{t-i} X_{t-k}) + \gamma_{X_\varepsilon}(k) \\ &= \sum_{j=1}^p \phi_j \gamma(k-j) - \sum_{i=1}^q \theta_i \gamma_{X_\varepsilon}(k-i) + \gamma_{X_\varepsilon}(k), \quad k \geq 0.\end{aligned}$$

ou

$$\begin{aligned}E[X_t \varepsilon_t] &= E \left[\left(\sum_{j=1}^p \phi_j X_{t-j} + \varepsilon_t - \sum_{i=1}^q \theta_i \varepsilon_{t-i} \right) \varepsilon_t \right] \\ &= E(\varepsilon_t \varepsilon_t) \\ &= \sigma_\varepsilon^2,\end{aligned}$$

et ε_t n'est pas corrélé avec le passé.

$$\left\{ \begin{array}{ll} \gamma_{X_\varepsilon}(k) = 0 & k > 0, \\ \gamma_{X_\varepsilon}(k) \neq 0 & k < 0, \\ \gamma_{X_\varepsilon}(k) = \sigma_\varepsilon^2 & k = 0. \end{array} \right.$$

La fonction $\gamma_{X_\varepsilon}(k)$ n'est pas paire.
 si $k > q$, tous les $\gamma_{X_\varepsilon}(k) = 0$, et on a :

$$\gamma(k) = \sum_{j=1}^p \phi_j \gamma(k-j),$$

et

$$\rho(k) = \sum_{j=1}^p \phi_j \rho(k-j).$$

Pour un $ARMA(p, q)$, la structure d'autocorrélation ne suit pas un schéma connu jusqu'au délai q mais ensuite le comportement est le même que celui d'un $AR(p)$.

2.4 Les modèles ARIMA(p,d,q)

Définition 2.4 un processus stationnaire X_t admet une représentation $ARIMA(p, d, q)$ minimale s'il satisfait : [8]

$$\Phi(L)(1-L)^d X_t = \Theta(L)\varepsilon_t, \forall t \in \mathbb{Z}.$$

avec les conditions suivantes :

1. $\phi_p \neq 0$ et $\theta_q \neq 0$.
2. Φ et Θ , polynômes de degrés resp p et q n'ont pas de racines communes et leurs racines sont de modules > 1 .
3. ε_t est un BB de variance σ^2 .

Un processus $ARIMA(p, d, q)$ convient pour modéliser une série temporelle comprenant une tendance polynômiale de degrés d , l'opérateur $(1-L)^d$ permettant de transformer un polynôme de degré d en une constante.

2.5 Méthode de Box & Jenkins

Box & Jenkins [2] ont élaboré une méthodologie pour identifier un modèle adéquat pour une série chronologique. Leur méthode est fondée sur les modèles *ARIMA*. Dans ce chapitre, les principales étapes de cette technique sont présentées. Pour les méthodes d'inférence présentées dans la suite, on supposera que T réalisations d'une série chronologique univariée, notées X_1, \dots, X_T ont été observées.

2.5.1 Identification du modèle

De façon générale, l'étape d'identification du modèle consiste à identifier le modèle qui représente au mieux la série étudiée. En d'autres mots, il s'agit de trouver un modèle stationnaire qui tient compte de la variabilité dans le temps et pour lequel il y a absence d'autocorrélation des résidus. Plus particulièrement, cette étape implique les méthodes d'estimation du paramètre d'intégration d , l'estimation des ordres p et q , les tests d'hétéroscédasticité, les tests de non stationnarité ou de racine unitaire.

Estimation du paramètre d'intégration

Lorsqu'un processus possède une non stationnarité de type DS, on parle alors d'un modèle intégré. Dans ce cas, il convient de déterminer l'ordre d'intégration d du processus filtré $(1 - L)^d X_t$ pour lequel le processus est stationnaire, d'où le processus $I(d)$.

Estimation des ordres p et q

Identification d'un processus *MA* Soit un processus stationnaire $(X_t)_{t \in \mathbb{Z}}$ satisfaisant une représentation $MA(q)$. Pour déterminer la va-

leur de q , on se base sur la fonction d'autocorrélation du processus MA . En fait, q correspondra au plus grand délai tel que l'autocorrélation n'est pas statistiquement égale à 0.

Identification d'un processus AR Soit un processus stationnaire $(X_t)_{t \in \mathbb{Z}}$ satisfaisant une représentation $AR(p)$. Globalement, l'identification d'un processus $AR(p)$ s'effectue de la même façon que celle d'un processus $MA(q)$. La seule différence réside dans le fait que c'est l'autocorrélation partielle, plutôt que l'autocorrélation, qui est utilisée.

Identification d'un processus $ARMA$ Il existe plusieurs méthodes pour déterminer les ordres p et q d'un processus $ARMA$. On peut également se baser sur les autocorrélations et les autocorrélations partielles. Pour obtenir l'ordre de la composante MA , il faut identifier l'autocorrélation significative dont l'ordre est le plus élevé; pour la composante AR , il faut identifier l'autocorrélation partielle significative dont l'ordre est le plus élevé.

Tests de non stationnarité

Il existe plusieurs tests de non stationnarité ou de racine unitaire, ces tests ayant pour objet de déterminer la présence de non stationnarité. Ces tests permettent aussi d'identifier le type de non stationnarité (TS ou DS) d'un processus. Parmi les tests de non stationnarité les plus populaires, on a ceux de Dickey & Fuller, Dickey Fuller augmenté et de Phillips & Perron.

2.5.2 Estimation des paramètres

L'estimation des paramètres d'un modèle $ARIMA(p, d, q)$ peut être effectuée lorsque'il est supposé que p , d et q peut être effectué de

différentes méthodes dans le domaine temporel parmi ces méthodes, nous avons :

1. Maximum de vraisemblance :

Une méthode populaire pour estimer les paramètres d'un modèle est le maximum de vraisemblance. La fonction de vraisemblance associée à un échantillon X_1, \dots, X_T *iid* d'une loi dont la densité est $f(x|\theta)$, avec $\theta = (\theta_1, \dots, \theta_k) \in \mathbb{R}^k$, est définie par :

$$L(\theta) = \prod_{t=1}^T f(X_t|\theta).$$

L'estimateur du maximum de vraisemblance (EMV) est la valeur $\widehat{\theta}_{EMV}$ qui maximise $L(\theta)$. Parfois, il est possible de déduire cet estimateur en dérivant $L(\theta)$ par rapport à chacun des paramètres $(\theta_1, \dots, \theta_k)$ et de résoudre le système à k équations

$$\frac{\partial L(\theta)}{\partial \theta_j} = 0, \quad \text{ou } j = 1, \dots, k.$$

2. Estimation de Yule Walker.

Dans le cas d'un $AR(p)$, on utilise les équations de Yule-Walker :

$$\begin{bmatrix} \rho(1) \\ \cdot \\ \cdot \\ \cdot \\ \rho(p) \end{bmatrix} = \begin{bmatrix} \rho(0) & \rho(1) & \cdot & \cdot & \cdot & \rho(p-1) \\ \rho(1) & \rho(0) & \cdot & \cdot & \cdot & \rho(p-2) \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \rho(p-1) & \rho(p-2) & \cdot & \cdot & \cdot & \rho(0) \end{bmatrix} \cdot \begin{bmatrix} \phi_1 \\ \phi_2 \\ \cdot \\ \cdot \\ \cdot \\ \phi_p \end{bmatrix}$$

pour déterminer $\hat{\phi}_i$ $i = 1, \dots, p$ en fonction de $\hat{\rho}(i)$ estimés.

De nombreux programmes informatiques mettent en oeuvre ces méthodes d'estimation du modèle ARIMA (Eviews, Spss,R,...), notamment les méthodes du maximum de vraisemblance .

Critères de choix des modèles

Souvent il n'est pas facile de déterminer un modèle unique . Le modèle qui est finalement choisi et celui qui minimise l'un des critères à partir T observations.

Critère standard 1. L'erreur absolue moyenne (Mean Absolute Error) :

$$MAE = \frac{1}{T} \sum_{t=1}^T |\varepsilon_t| .$$

2. L'erreur quadratique moyenne (Mean Squared Error) :

$$MSE = \frac{1}{T} \sum_{t=1}^T \varepsilon_t^2 .$$

3. La racine carrée de l'erreur quadratique moyenne (Root Mean Square Error) :

$$RMSE = \sqrt{\frac{1}{T} \sum_{t=1}^T \varepsilon_t^2} .$$

4. Ecart absolu moyen en pourcentage (Mean Absolute Percent Error) :

$$MAPE = \frac{100}{T} \sum_{t=1}^T \left| \frac{\varepsilon_t}{X_t} \right| .$$

Critère d'information

1. Akaike (1969) :

$$AIC(p, q) = \log(\hat{\sigma}_\varepsilon^2) + 2 \frac{p+q}{T} .$$

2. Schwarz (1977) :

$$BIC(p, q) = \log(\hat{\sigma}_\varepsilon^2) + (p+q) \frac{\log T}{T} .$$

3. Hannan_Quinn(1979) :

$$\varphi(p, q) = \log(\hat{\sigma}_\varepsilon^2) + (p+q)c \frac{\log(\log T)}{T} , \text{ avec } c > 2 .$$

2.5.3 Validation du modèle

Il s'agit de vérifier notamment que les résidus du modèle *ARMA* estimé, résidus notés $\hat{\varepsilon}_t$, vérifient les propriétés requises pour que l'estimation soit valide, à savoir qu'ils suivent un processus BB, non autocorrélé et de même variance, et qu'ils suivent une loi normale. Si ces hypothèses ne sont pas rejetées, on peut alors mener des tests sur les paramètres.

Significativité des paramètres

Les coefficients du modèle doivent s'écarter significativement de zéro, pour cela on utilise le test de student classique.

On rejette l'hypothèse nulle H_0 :

$$H_0 : \theta_j = 0, \text{ si } |tc| > |T_{T-q}^\alpha|, \text{ ou } |tc| = \left| \frac{\widehat{\theta}}{\widehat{\sigma_\theta}} \right|.$$

S'il s'avère qu'un ou plusieurs paramètres du modèle ne sont pas significativement différents de 0, on estime à nouveau le modèle.

Test sur les d'autocorrélations

Le test de Box Pierce aide à identifier les processus de bruit blanc. Ce test s'écrit :

$$\begin{cases} H_0 : \rho(1) = \rho(2) = \dots = \rho(h) = 0 , \\ H_1 : \rho(i) \neq 0 \text{ pour au moins } i , i \in \{1, \dots, h\} . \end{cases}$$

Une statistique de test proposée par Box et Pierce est

$$Q = T \sum_{k=1}^h \widehat{\rho}^2(k).$$

h : nombre de retards,

T : nombre d'observations,

$\widehat{\rho}(k)$: autocorrélation empirique d'ordre k .

La loi asymptotique des statistiques Q est la khi-carré à h degrés de liberté. On rejette donc H_0 , au seuil α , si la statistique Q est supérieure au khi-carré dans la table au seuil $(1 - \alpha)$ et h degrés de liberté.

Test de normalité des résidus

Plusieurs des modèles de séries chronologiques supposent que les termes d'innovation sont indépendants et distribués selon la loi Normale. Une façon de vérifier cette hypothèse consiste à étudier les résidus. Un des tests permettant de vérifier la normalité des résidus est celui de Jarque et Bera.

Les hypothèses à confronter sont

$$\begin{cases} H_0 : \varepsilon_t \sim N(0, 1), \\ H_1 : \varepsilon_t \not\sim N(0, 1). \end{cases}$$

Avant de décrire le test pour les résidus d'un modèle de séries chronologiques, prenons le cas de T observations X_1, \dots, X_T indépendantes et de même loi.

Dans ce cas, la statistique du test de Jarque et Bera, de loi asymptotique khi-carré à deux degrés de liberté, est définie par

$$JB = \frac{T}{6}\beta_1^2 + \frac{T}{24}(\beta_2 - 3)^2.$$

où

$$\beta_1 = \frac{1}{T} \sum_{i=1}^T \left(\frac{X_i - \bar{X}}{S} \right)^3 \quad \text{et} \quad \beta_2 = \frac{1}{T} \sum_{i=1}^T \left(\frac{X_i - \bar{X}}{S} \right)^4.$$

sont respectivement les coefficients d'asymétrie et d'aplatissement. Ici, \bar{X} et S sont respectivement la moyenne et l'écart-type empiriques. Sous l'hypothèse de normalité, on peut montrer que :

$$\frac{\sqrt{T}\beta_1}{\sqrt{6}} \rightarrow N(0, 1), \quad \frac{\sqrt{T}(\beta_2 - 3)}{\sqrt{24}} \rightarrow N(0, 1).$$

et que asymptotiquement, ces deux variables aléatoires sont indépendantes.

Il s'ensuit que

$$\frac{T\beta_1^2}{6} \rightarrow \chi_1^2 \text{ et } \frac{T(\beta_2 - 3)^2}{24} \rightarrow \chi_1^2.$$

où χ_v^2 représente la loi khi-carré à v degrés de liberté. Ainsi, la loi asymptotique de JB est la distribution khi-carré à deux degrés de liberté. On rejette donc l'hypothèse de normalité, au seuil α , si $JB > \chi_{\alpha, 2}^2$.

2.5.4 La prévision

Une fois qu'un modèle acceptable est trouvé pour la série chronologique à l'étude, les prévisions peuvent être calculées. On note \widehat{X}_{T+h} la prévision de X_{T+h} au temps $T+h$ où T est la taille de l'échantillon des observations X_t et h l'horizon de la prévision.

La prévision est calculée par la formule suivante :

$$\widehat{X}_{T+h} = E(X_{T+h}/X_T, X_{T-1}, \dots, X_1).$$

Chapitre 3

Application : Modélisation par la méthode de Box-Jenkins

Nous intéressons à modéliser par la méthode de Box-Jenkins les données réelles "l'évolution annuelle des cas de **paludisme** de direction de la santé et de la population de Ouargla -Monographie sanitaire".

Le tableau suivant présente les données :

Année	2000	2001	2002	2003	2004	2005	2006	2007
N des cas de paludisme	28	25	16	06	02	01	00	02
Année	2008	2009	2010	2011	2012	2013	2014	2015
N des cas de paludisme	00	01	03	00	07	10	04	23
Année	2016	2017	2018	2019	2020	2021		
N des cas de paludisme	13	13	28	14	19	18		

TAB. 3.1 – L'évolution annuelle des cas de paludisme.

3.1 Etude de la stationnarité

Le graphe de la série "les cas de paludisme" montre que la série est non stationnaire :

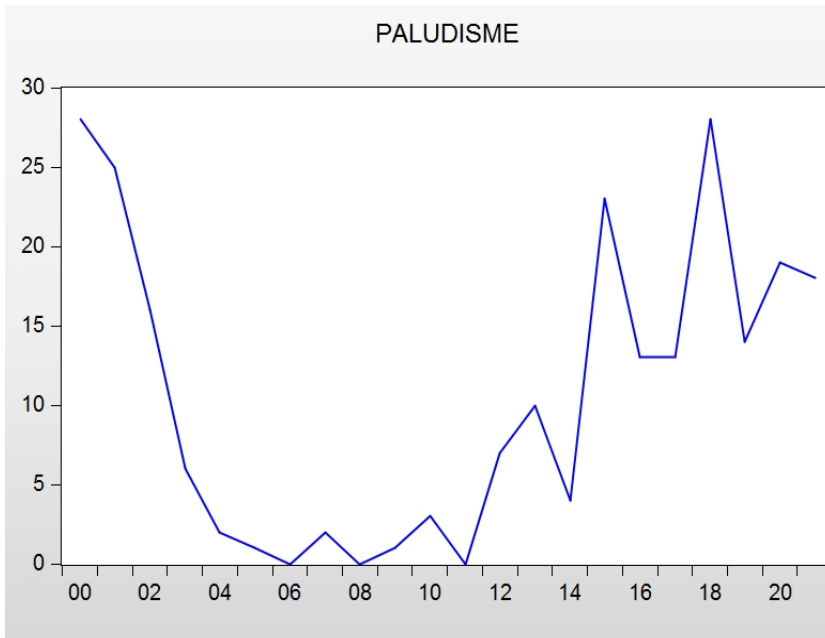


FIG. 3.1 – Représentation graphique de la série brute.

On utilise le test de racine unitaire (test ADF) pour confirmer la non stationnarité :

Modèle sans tendance et avec constante

Null Hypothesis: PALUDISME has a unit root				
Exogenous: Constant				
Lag Length: 3 (Automatic - based on SIC, maxlag=4)				
			t-Statistic	Prob.*
Augmented Dickey-Fuller test statistic			-2.278205	0.1887
Test critical values:	1% level		-3.857386	
	5% level		-3.040391	
	10% level		-2.660551	
*MacKinnon (1996) one-sided p-values.				
Warning: Probabilities and critical values calculated for 20 observations and may not be accurate for a sample size of 18				
Augmented Dickey-Fuller Test Equation				
Dependent Variable: D(PALUDISME)				
Method: Least Squares				
Date: 05/28/22 Time: 21:21				
Sample (adjusted): 2004 2021				
Included observations: 18 after adjustments				
Variable	Coefficient	Std. Error	t-Statistic	Prob.
PALUDISME(-1)	-0.358374	0.157305	-2.278205	0.0403
D(PALUDISME(-1))	-0.294981	0.177297	-1.663765	0.1201
D(PALUDISME(-2))	0.175520	0.160353	1.094585	0.2936
D(PALUDISME(-3))	0.734056	0.156170	4.700364	0.0004
C	3.729901	1.666240	2.238514	0.0433
R-squared	0.751163	Mean dependent var		0.666667
Adjusted R-squared	0.674598	S.D. dependent var		7.798944
S.E. of regression	4.448832	Akaike info criterion		6.053293
Sum squared resid	257.2973	Schwarz criterion		6.300619
Log likelihood	-49.47964	Hannan-Quinn criter.		6.087396
F-statistic	9.810765	Durbin-Watson stat		2.042698
Prob(F-statistic)	0.000697			

La p-value > 0.05 donc la série a une racine unitaire, et la constante est non significative (la probabilité critique affectées à la constante est supérieure à 0.05).

Modèle avec tendance et avec constante

Null Hypothesis: PALUDISME has a unit root
 Exogenous: Constant, Linear Trend
 Lag Length: 3 (Automatic - based on SIC, maxlag=4)

	t-Statistic	Prob.*
Augmented Dickey-Fuller test statistic	-2.297076	0.4144
Test critical values:		
1% level	-4.571559	
5% level	-3.690814	
10% level	-3.286909	

*Mackinnon (1996) one-sided p-values.

Warning: Probabilities and critical values calculated for 20 observations
 and may not be accurate for a sample size of 18

Augmented Dickey-Fuller Test Equation
 Dependent Variable: D(PALUDISME)
 Method: Least Squares
 Date: 05/28/22 Time: 21:23
 Sample (adjusted): 2004 2021
 Included observations: 18 after adjustments

Variable	Coefficient	Std. Error	t-Statistic	Prob.
PALUDISME(-1)	-0.712986	0.310388	-2.297076	0.0404
D(PALUDISME(-1))	-0.333079	0.174973	-1.903602	0.0812
D(PALUDISME(-2))	-0.100386	0.261701	-0.383590	0.7080
D(PALUDISME(-3))	0.396973	0.298286	1.330848	0.2080
C	-4.630776	6.569209	-0.704921	0.4943
@TREND("2000")	0.885977	0.674592	1.313352	0.2136
R-squared	0.782436	Mean dependent var		0.666667
Adjusted R-squared	0.691784	S.D. dependent var		7.798944
S.E. of regression	4.329753	Akaike info criterion		6.030100
Sum squared resid	224.9612	Schwarz criterion		6.326890
Log likelihood	-48.27090	Hannan-Quinn criter.		6.071023
F-statistic	8.631237	Durbin-Watson stat		1.684693
Prob(F-statistic)	0.001141			

La p-value > 0.05 donc la série a une racine unitaire, la tendance et la constante sont significative (la probabilité critique < 0.05).

Modèle sans tendance et sans constante

Null Hypothesis: PALUDISME has a unit root
 Exogenous: None
 Lag Length: 3 (Automatic - based on SIC, maxlag=4)

	t-Statistic	Prob.*
Augmented Dickey-Fuller test statistic	-0.757272	0.3741
Test critical values:		
1% level	-2.699769	
5% level	-1.961409	
10% level	-1.606610	

*MacKinnon (1996) one-sided p-values.

Warning: Probabilities and critical values calculated for 20 observations and may not be accurate for a sample size of 18

Augmented Dickey-Fuller Test Equation
 Dependent Variable: D(PALUDISME)
 Method: Least Squares
 Date: 05/28/22 Time: 21:25
 Sample (adjusted): 2004 2021
 Included observations: 18 after adjustments

Variable	Coefficient	Std. Error	t-Statistic	Prob.
PALUDISME(-1)	-0.085310	0.112654	-0.757272	0.4614
D(PALUDISME(-1))	-0.475096	0.179196	-2.651259	0.0190
D(PALUDISME(-2))	0.064474	0.172956	0.372774	0.7149
D(PALUDISME(-3))	0.651440	0.172117	3.784876	0.0020
R-squared	0.655247	Mean dependent var		0.666667
Adjusted R-squared	0.581371	S.D. dependent var		7.798944
S.E. of regression	5.046034	Akaike info criterion		6.268212
Sum squared resid	356.4745	Schwarz criterion		6.466073
Log likelihood	-52.41391	Hannan-Quinn criter.		6.295495
Durbin-Watson stat	1.751349			

La p-value > 0.05 donc la série a une racine unitaire.

D'après ce qu'on a vu précédemment, la série "les cas de paludisme" est non stationnaire.

3.1.1 Stationnarisation de la série

On utilise la première différence et pour confirmer l'élimination de la racine unitaire on utilise autre fois le test ADF.

$$d\text{paludisme} = \text{paludisme} - \text{paludisme}(-1).$$

1)

Null Hypothesis: DPALUDISME has a unit root
 Exogenous: Constant
 Lag Length: 2 (Automatic - based on SIC, maxlag=4)

		t-Statistic	Prob.*
Augmented Dickey-Fuller test statistic		-2.299028	0.1826
Test critical values:	1% level	-3.857386	
	5% level	-3.040391	
	10% level	-2.660551	

*MacKinnon (1996) one-sided p-values.
 Warning: Probabilities and critical values calculated for 20 observations and may not be accurate for a sample size of 18

Augmented Dickey-Fuller Test Equation
 Dependent Variable: D(DPALUDISME)
 Method: Least Squares
 Date: 05/28/22 Time: 22:06
 Sample (adjusted): 2004 2021
 Included observations: 18 after adjustments

Variable	Coefficient	Std. Error	t-Statistic	Prob.
DPALUDISME(-1)	-0.842574	0.366492	-2.299028	0.0374
D(DPALUDISME(-1))	-0.685297	0.281521	-2.434266	0.0289
D(DPALUDISME(-2))	-0.633630	0.170775	-3.710325	0.0023
C	0.786219	1.199203	0.655618	0.5227

R-squared	0.891771	Mean dependent var	0.500000
Adjusted R-squared	0.868579	S.D. dependent var	13.98844
S.E. of regression	5.071085	Akaike info criterion	6.278117
Sum squared resid	360.0226	Schwarz criterion	6.475977
Log likelihood	-52.50305	Hannan-Quinn criter.	6.305399
F-statistic	38.45192	Durbin-Watson stat	1.855955
Prob(F-statistic)	0.000001		

La constante n'est pas significative, en outre la série a un racine unitaire.

2)

Null Hypothesis: DPALUDISME has a unit root				
Exogenous: Constant, Linear Trend				
Lag Length: 2 (Automatic - based on SIC, maxlag=4)				
		t-Statistic	Prob.*	
Augmented Dickey-Fuller test statistic				
	1% level	-0.446291	0.9762	
Test critical values:	5% level	-4.571559		
	10% level	-3.690814		
		-3.286909		
*MacKinnon (1996) one-sided p-values.				
Warning: Probabilities and critical values calculated for 20 observations and may not be accurate for a sample size of 18				
Augmented Dickey-Fuller Test Equation				
Dependent Variable: D(DPALUDISME)				
Method: Least Squares				
Date: 05/28/22 Time: 22:14				
Sample (adjusted): 2004 2021				
Included observations: 18 after adjustments				
Variable	Coefficient	Std. Error	t-Statistic	Prob.
DPALUDISME(-1)	-0.267163	0.598629	-0.446291	0.6627
D(DPALUDISME(-1))	-1.120673	0.455459	-2.460536	0.0286
D(DPALUDISME(-2))	-0.861225	0.252904	-3.405345	0.0047
C	6.664909	5.021469	1.327283	0.2073
@TREND("2000")	-0.462004	0.383578	-1.204460	0.2499
R-squared	0.902637	Mean dependent var		0.500000
Adjusted R-squared	0.872679	S.D. dependent var		13.98844
S.E. of regression	4.991374	Akaike info criterion		6.283433
Sum squared resid	323.8796	Schwarz criterion		6.530758
Log likelihood	-51.55089	Hannan-Quinn criter.		6.317536
F-statistic	30.13007	Durbin-Watson stat		2.170227
Prob(F-statistic)	0.000002			

La constante et la tendance sont non significative, en outre la série a un racine unitaire.

3)

Null Hypothesis: DPALUDISME has a unit root				
Exogenous: None				
Lag Length: 2 (Automatic - based on SIC, maxlag=4)				
		t-Statistic	Prob.*	
Augmented Dickey-Fuller test statistic				
		-2.385162	0.0203	
Test critical values:				
	1% level	-2.699769		
	5% level	-1.961409		
	10% level	-1.606610		
*MacKinnon (1996) one-sided p-values.				
Warning: Probabilities and critical values calculated for 20 observations and may not be accurate for a sample size of 18				
Augmented Dickey-Fuller Test Equation				
Dependent Variable: D(DPALUDISME)				
Method: Least Squares				
Date: 05/28/22 Time: 22:17				
Sample (adjusted): 2004 2021				
Included observations: 18 after adjustments				
Variable	Coefficient	Std. Error	t-Statistic	Prob.
DPALUDISME(-1)	-0.856023	0.358895	-2.385162	0.0307
D(DPALUDISME(-1))	-0.673256	0.275530	-2.443491	0.0274
D(DPALUDISME(-2))	-0.630407	0.167428	-3.765235	0.0019
R-squared	0.888448	Mean dependent var		0.500000
Adjusted R-squared	0.873575	S.D. dependent var		13.96844
S.E. of regression	4.973773	Akaike info criterion		6.197246
Sum squared resid	371.0762	Schwarz criterion		6.345641
Log likelihood	-52.77522	Hannan-Quinn criter.		6.217708
Durbin-Watson stat	1.795065			

La série a un racine unitaire.

Donc on utilise autre fois le test ADF pour la deuxième différence.

$$ddpaludisme = dpaludisme - dpaludisme(-1)$$

1)

Null Hypothesis: DDPALUDISME has a unit root				
Exogenous: Constant				
Lag Length: 1 (Automatic - based on SIC, maxlag=4)				
			t-Statistic	Prob.*
Augmented Dickey-Fuller test statistic				
			-12.24728	0.0000
Test critical values:				
	1% level		-3.857385	
	5% level		-3.040391	
	10% level		-2.660551	
*MacKinnon (1996) one-sided p-values.				
Warning: Probabilities and critical values calculated for 20 observations and may not be accurate for a sample size of 18				
Augmented Dickey-Fuller Test Equation				
Dependent Variable: D(DDPALUDISME)				
Method: Least Squares				
Date: 05/28/22 Time: 22:20				
Sample (adjusted): 2004 2021				
Included observations: 18 after adjustments				
Variable	Coefficient	Std. Error	t-Statistic	Prob.
DDPALUDISME(-1)	-3.164689	0.258399	-12.24728	0.0000
D(DDPALUDISME(-1))	0.895040	0.144478	6.195009	0.0000
C	0.940536	1.357630	0.692778	0.4990
R-squared	0.955487	Mean dependent var		-0.277778
Adjusted R-squared	0.949552	S.D. dependent var		25.60056
S.E. of regression	5.750042	Akaike info criterion		6.487303
Sum squared resid	495.9448	Schwarz criterion		6.635699
Log likelihood	-55.38573	Hannan-Quinn criter.		6.507765
F-statistic	160.9907	Durbin-Watson stat		1.637707
Prob(F-statistic)	0.000000			

2)

Null Hypothesis: DDPALUDISME has a unit root
 Exogenous: Constant, Linear Trend
 Lag Length: 1 (Automatic - based on SIC, maxlag=4)

	t-Statistic	Prob.*
Augmented Dickey-Fuller test statistic	-14.77297	0.0001
Test critical values:		
1% level	-4.571559	
5% level	-3.690814	
10% level	-3.286909	

*MacKinnon (1996) one-sided p-values.

Warning: Probabilities and critical values calculated for 20 observations
 and may not be accurate for a sample size of 18

Augmented Dickey-Fuller Test Equation
 Dependent Variable: D(DDPALUDISME)
 Method: Least Squares
 Date: 05/28/22 Time: 22:22
 Sample (adjusted): 2004 2021
 Included observations: 18 after adjustments

Variable	Coefficient	Std. Error	t-Statistic	Prob.
DDPALUDISME(-1)	-3.275319	0.221710	-14.77297	0.0000
D(DDPALUDISME(-1))	0.958624	0.124087	7.725439	0.0000
C	8.421017	3.029036	2.780098	0.0147
@TREND("2000")	-0.598620	0.224434	-2.667243	0.0184
R-squared	0.970485	Mean dependent var		-0.277778
Adjusted R-squared	0.964161	S.D. dependent var		25.60056
S.E. of regression	4.846514	Akaike info criterion		6.187527
Sum squared resid	328.8418	Schwarz criterion		6.385387
Log likelihood	-51.68774	Hannan-Quinn criter.		6.214809
F-statistic	153.4464	Durbin-Watson stat		2.241931
Prob(F-statistic)	0.000000			

3)

Null Hypothesis: DDPALUDISME has a unit root				
Exogenous: None				
Lag Length: 1 (Automatic - based on SIC, maxlag=4)				
			t-Statistic	Prob.*
Augmented Dickey-Fuller test statistic			-12.44385	0.0001
Test critical values:	1% level		-2.699769	
	5% level		-1.961409	
	10% level		-1.606610	
*MacKinnon (1996) one-sided p-values.				
Warning: Probabilities and critical values calculated for 20 observations and may not be accurate for a sample size of 18				
Augmented Dickey-Fuller Test Equation				
Dependent Variable: D(DDPALUDISME)				
Method: Least Squares				
Date: 05/28/22 Time: 22:23				
Sample (adjusted): 2004 2021				
Included observations: 18 after adjustments				
Variable	Coefficient	Std. Error	t-Statistic	Prob.
DDPALUDISME(-1)	-3.162572	0.254147	-12.44385	0.0000
D(DDPALUDISME(-1))	0.896179	0.142101	6.306636	0.0000
R-squared	0.954063	Mean dependent var		-0.277778
Adjusted R-squared	0.951192	S.D. dependent var		25.60056
S.E. of regression	5.655822	Akaike info criterion		6.407687
Sum squared resid	511.8131	Schwarz criterion		6.506617
Log likelihood	-55.66918	Hannan-Quinn criter.		6.421328
Durbin-Watson stat	1.588443			

D'après le test ADF, la série obtenue est une série stationnaire (p-value<0.05).

Graphique de la série stationnaire

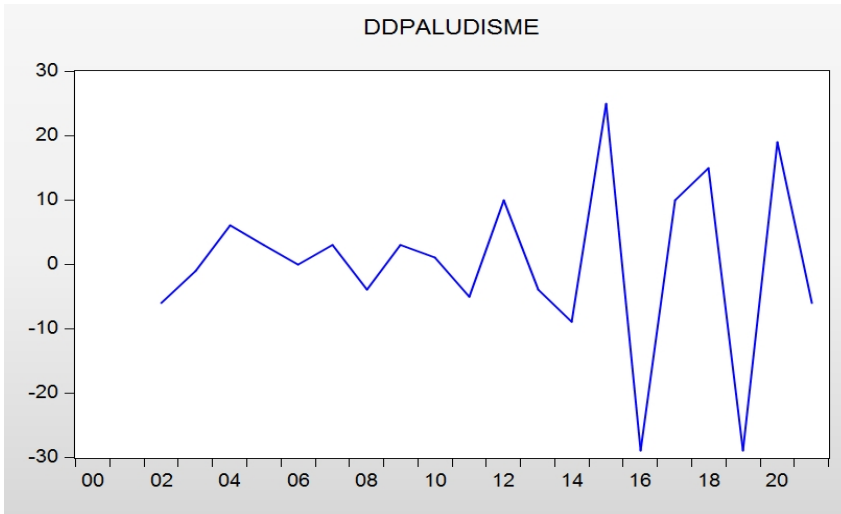


























FIG. 3.2 – Représentation graphique de la série stationnaire.

3.1.2 Modélisation de la série

A ce phase, la série ddpaludisme est stationnaire, on cherche un modèle $ARMA(p; q)$ qui représente cette série. Pour déterminer quel processus représente le mieux notre série, nous examinons les autocorrelations simples et partielles.

Corrélogramme de la série stationnaire

Date: 05/28/22 Time: 22:27
 Sample: 2000 2021
 Included observations: 20

Autocorrelation	Partial Correlation	AC	PAC	Q-Stat	Prob	
		1	-0.669	-0.669	10.368	0.001
		2	0.013	-0.786	10.372	0.006
		3	0.496	-0.064	16.752	0.001
		4	-0.600	-0.060	26.647	0.000
		5	0.341	0.033	30.049	0.000
		6	-0.007	-0.250	30.050	0.000
		7	-0.181	-0.006	31.159	0.000
		8	0.191	-0.009	32.491	0.000
		9	-0.088	0.104	32.803	0.000
		10	-0.031	-0.105	32.844	0.000
		11	0.092	-0.007	33.254	0.000
		12	-0.088	-0.078	33.679	0.001

L'examen de ce corrélogramme montre que l'on peut s'anticiper des modèles de type $MA(1)$, $MA(3)$, $MA(4)$, $AR(1)$, $AR(2)$, $ARMA(1, 1)$, $ARMA(1, 3)$, $ARMA(1, 4)$, $ARMA(2, 1)$, $ARMA(2, 3)$, $ARMA(2, 4)$.

3.2 Estimation

Pour choisir un bon modèle, nous construisons un tableau de comparaison des critères d'information et nous retenons le modèle qui minimise ces critères.

	AKAIKE	SCHWARZ
<i>ARMA</i> (1, 3)	7,1835	7,3827
<i>ARMA</i> (1, 1)	7,1843	7,3835
<i>ARMA</i> (1, 4)	7,2007	7,3998
<i>MA</i> (1)	7,3247	7,4740
<i>ARMA</i> (2, 1)	7,4194	7,6186
<i>MA</i> (4)	7,6101	7,7594
<i>AR</i> (1)	7,6797	7,8291
<i>ARMA</i> (2, 4)	7,7101	7,9092
<i>MA</i> (3)	7,7301	7,8795
<i>ARMA</i> (2, 3)	7,8029	8,0020
<i>AR</i> (2)	8,2639	8,4133

TAB. 3.2 – Comparaison entre les critères des modèles.

En principe, d'après le tableau nous pouvons retenir le modèle *ARMA*(1, 3) parce que c'est le modèle qui minimise les critères d'information de AKAIKE et SCHWARZ par rapport aux autres modèles.

3.3 Validation du modèle

Nous allons diagnostiquer notre modèle car il ne peut être mis à défaut, avec des tests suivants :

3.3.1 Test de significativité des paramètres

Ce test consiste à vérifier que les paramètres du modèle qui ont été estimés sont statistiquement différents de 0.

Les hypothèses du test sont :

$$\left\{ \begin{array}{l} H_0 : \theta = 0, \text{ le coefficient est non significatif;} \\ H_1 : \theta \neq 0, \text{ le coefficient est significatif.} \end{array} \right.$$

Ainsi, au risque de 5%, les paramètres du modèle ne sont pas statistiquement différent de zéro car pour la partie AR : $|t - stat| = 4.495655 > 2$, mais pour la partie MA : $|t - stat| = 0.000255 < 2$. Alors le modèle $ARMA(1, 3)$ n'est pas valide.

Nous pouvons retenir le modèle $AR(1)$ car $|t - stat| = 3,867604 > 2$. (le paramètre θ estimé est statistiquement significatif).

3.3.2 Test de normalité des résidus

On veut tester la normalité des résidus du modèle $AR(1)$. Le test que nous utilisons est de *Jarque et Bera*.

Ce test se formule avec les hypothèses suivantes :

$$\left\{ \begin{array}{l} H_0 : \text{les résidus sont normalement distribués;} \\ H_1 : \text{les résidus ne sont pas normalement distribués.} \end{array} \right.$$

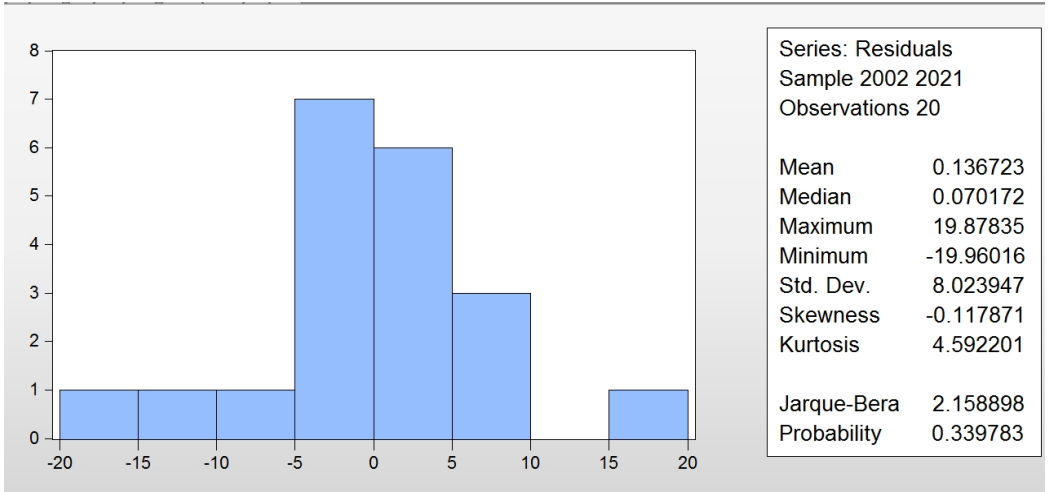


FIG. 3.3 – Normalité des résidus du modèle $AR(1)$.

La statistique $BJ < \chi_{1-\alpha}^2(2)$ ($2.158898 < 5,991$) avec $\alpha = 0,05$; alors on accepte l'hypothèse H_0 de normalité des résidus au seuil α du modèle $AR(1)$.

Alors notre série est modélisée par un $AR(1)$:

$$X_t = 0.35 + 0.65X_{t-1} + \varepsilon_t.$$

3.4 La prévision

Les résultats de prévision de cinq années sont dans le tableau suivant :

année	Les prévisions
2022	12
2023	8
2024	5
2025	3
2026	2

TAB. 3.3 – Prévision du cas de paludisme.

Par exemple on obtient la prévision de l'année 2022 :

$$\begin{aligned}\widehat{X}_{T+1} &= E[X_{T+1}/X_t] \\ &= E[0.35 + 0.65X_T + \varepsilon_{T+1}/X_T] \\ &= 0.35 + 0.65X_T \\ &\simeq 12\end{aligned}$$

Conclusion

Après avoir décrit divers modèles de séries temporelles. Nous nous intéressons à savoir comment sélectionner un modèle approprié qui peut produire des prévisions précises basées sur une description des données d'une série chronologique, et comment déterminer les ordres de modèle optimaux.

Les statisticiens *George Box et Wilym Jenkins [2]* ont proposé une approche de prévision de séries temporelles univariées, basée sur l'utilisation des modèles *ARIMA*. Leur concept est d'une importance fondamentale dans le domaine de l'analyse et de la prévision des séries chronologiques.

Le but de notre étude est de trouver le meilleur modèle pour prédire le nombre des cas de paludisme de Ouargla, on utilisant les modèles de type *ARIMA* à l'aide des tests statistiques.

Bibliographie

- [1] Arnaud,R.(2022). Séries chronologiques.Universété de Bourgogne.
- [2] Box, G. E. P., Jenkins, G. M. (1976). Time Series Analysis Forecasting and Control. Revised Edition.
- [3] Charpentier, A. (2004). Cours des séries temporelles. Théorie et Applications.Université de Paris.
- [4] Corinne, P. (2005). Magister d'Economie. Séries chronologiques : Quelques éléments du cours.Universété de Paris.
- [5] Bourbonnais, R. Terraza, M. (2004). Analyse des séries Temporelles. Dunod,Paris.
- [6] Dendouga, S.(2020).Mémoire de Master. Séries temporelles : Théorie et application. Universété de Biskra.
- [7] Gasmi, L. (2014). Mémoire de Magister : Application des AGs pour la détermination des paramètres du modèle de séries chronologiques. Université de Bichar.

- [8] Girard, Y. (2011). Série chronologiques à une et plusieurs variables : synthèse des méthodes classiques et modèles à base de copules. Université du Québec.

- [9] Rouba, S. (2019). Mémoire de Master : Modélisation Statistique d'une Série Chronologique. Université de Biskra.

- [10] Yannig, G. (2021). Les processus AR et MA : Séries chronologiques.

- [11] Zou,H et Yang, Y.(2004).Combining time series models for forecasting.

الملخص:

يعتمد عملنا في هاته المذكرة على نمذجة السلاسل الزمنية والتنبؤ بها من خلال المنهجية التي وضعها بوكس وجينكس، عن طريق النماذج الخطية.

اولا نقدم بعض المفاهيم الاساسية المتعلقة بالسلسلة الزمنية، يلي ذلك عرض لنماذج خطية وهي: AR, MA, ARMA, ARIMA, ثم نعرض طريقة بوكس جينكس، ننهي هذا العمل من خلال تطبيق هذه الطريقة على بيانات حقيقية " عدد حالات الملاريا المأخوذة من دائرة الصحة والسكان بورقلة " لتحديد النموذج الانسب للتنبؤ باستخدام برنامج الايفيوز.

الكلمات المفتاحية: النماذج الخطية، طريقة بوكس وجينكس، التنبؤ، السلسلة الزمنية، AR, AM, ARMA, ARIMA.

Résumé:

Notre travail repose dans ce mémoire sur la modélisation et la prévision des séries chronologique à travers la méthodologie élaborée par Box et Jenkins , au moyen des modèles linéaires .

Tout d'abord , nous introduisons quelques concepts de base liés aux séries chronologique . S'ensuit la présentation des modèles linéaires , à savoir le modèle AR , MA , ARMA , ARIMA . Puis nous présentons la méthode de Box et Jenkins . Nous terminons ce travail par une application de cette méthode à des données réelles " le nombre des cas de paludisme tirés de Direction de la santé et de la population de Ouargla " afin de déterminer le modèle le plus approprié pour prédire , en utilisant le logiciel Eviews .

Mot clés: modèle linéaire, méthode de Box et Jenkins, prévision, série chronologique, AR, AM, ARMA, ARIMA.

Abstract:

Our work is based in this memory on the modeling and forecasting of time series through the methodology developed by Box and Jenkins, using linear models.

First, we introduce some basic concepts related to time series. This is followed by the presentation of the linear models, namely the AR, MA, ARMA, ARIMA model. Then we present the Box and Jenkins method. We end this work by applying this method to real data "the number of malaria cases taken from the Department of Health and Population of Ouargla " in order to determine the most appropriate model to predict, using the Eviews software. .

Keywords: linear model, Box and Jenkins method, forecast, time series, AR, AM, ARMA, ARIMA.