# A Recursive Bayesian Approach to Estimation of Prevalence of High Blood Pressure Among Different Age Groups

R. K., Ogundeji [1*], S.O.N. Agwuegbo [2] And I.A. Adeleke [3]

1 *, Department of Mathematics, University of Lagos, Akoka, Nigeria.
2, Department of Statistics, Federal University of Agriculture, Abeokuta, Nigeria.
3, Department of Acturial Science, University of Lagos, Akoka, Nigeria.
Corresponding Author email: rogundeji@unilag.edu.ng

## Abstract

High blood pressure (or hypertension) is a major public health issue affecting old aged adults in many countries and is a major risk factor in the development of stroke, cardiovascular and chronic kidney disease. High blood pressure of recent is becoming an important area of research due to its high prevalence among lower aged groups (i.e. below 30 years old). Many studies have been conducted on prevalence of high blood pressures amongst adults at local and national levels, and at urban or rural areas and all pointing to the fact that there has been increasing prevalence of hypertension across the globe. Since blood pressure tends to rise with age, there is the need to investigate the prevalence of high blood pressures at different age groups in order to describe the process which generates a particular signal or set of observations. In this study, staff data from the University of Lagos medical Centre were used for the research. A recursive Bayesian approach to dynamic state space estimation was developed to model the prevalence of high blood pressure together with the use of analytic solutions based on the Kalman filter. Based on model diagnostic criteria adopted, the result generated a best fit autoregressive model for the number patients with hypertension. The Kalman filter provided an optimal estimate of the linear state space approach to modelling dynamic system. The state space approach to modelling dynamic system in this study focused on discrete time formulation by using difference equations to model the evolution of the system with time, and measurements assumed to be available at discrete time. It provided a generic and flexible framework for modelling the prevalence of high blood pressure.

## 1 Introduction

High blood pressure (or Hypertension) is a major public health issue affecting older adults in many countries due to its high prevalence all around the globe. High blood pressure symptoms are difficult

to identify and is a major risk factor for coronary heart disease, stroke, and chronic kidney disease. High blood pressure puts a greater risk for developing life changing and life threatening conditions and is a silent killer. Many studies have been conducted on prevalence of high blood pressures amongst adults at local and national levels, and at urban or rural areas and all pointing to the fact that there has been increasing prevalence of hypertension across the globe Adediran *et al.* [1]; Ekwunife and Aguwa [2]; Wang *et al.* [3].

According to Adeloyea *et al.* [4] and Rao *et al.* [5], the rate at which people are dying from this silent disease is quite alarming and suggested that hypertension is high in Nigeria, but the overall awareness of this silent killer disease (hypertension) cases is low in Nigeria. A number of researches have been carried out on hypertension at old age but only few of them considered having hypertension at early age (i.e. between 18 and 30 years old). Since blood pressure tends to rise with age, there is the need to investigate the prevalence of high blood pressures at different age groups in order to describe the process which generates a particular signal or set of observations.

The aim of this study is to model high blood pressure from a real data generation process of a system that changes over time. In order to analyze and make inference about prevalence of high blood pressure at different age groups, at least two models are required: firstly a model describing the evolution of the state with time (the system model), and secondly a model relating noisy measurements to the state (the measurement model). This study assumes that these models are available in a probabilistic form. Since hypertension is related to age, the probabilistic state space formulation and the requirement for the updating of information for new measurements are ideally suited within a Bayesian framework. This provides a rigorous general framework for dynamic state estimation problem and difference equations are used to model the evolution of the system with time, and measurements are assumed to be available at discrete time.

In Bayesian approach to dynamic state estimation, one attempts to construct the probability density function (pdf) of the state based on all available information including the set of received measurements. Since the pdf embodies all available statistical information, it may be said to be the complete solution to the estimation problem, Gordon *et al.* [6].

In principle, an optimal estimate is obtained in the form of posterior distributions, which incorporate both scientist's beliefs and the observations, in a well-founded probabilistic framework. In particular, the problem of parameter estimation and model selection can be summarized by the posterior probability of each model and this distribution is meaningful and certainly easier to interpret than say classical P-values.

In Gordon *et al.* [6], analytic solutions to Bayesian estimation problem are only available for a relatively small and restrictive choice of system and measurement models, and the well-known approach is the Kalman filter. The Bayesian formulation means that the Kalman filter is regarded as a way of updating the prior probability distribution when a new observation becomes available to give a revised posterior distribution. The Bayesian approach also enables the analyst to consider the case where several different models are entertained and it is required to choose a single model to represent the process, Chatfield [7]. Kalman [8] estimated coefficient of non-linear differential equations using an optimal sequential estimation techniques often referred to as Kalman filter. Essentially Kalman filtering is a method of signal processing which provides optimal estimates of the current state of a dynamical system.

Kalman's derivation took place within the context of state space models whose core is the recursive least square estimation. Within the state space notation, the Kalman filter derivation rests on the assumption of normality of the initial state vector, and as well as the disturbances of the system. The state of a system is defined to be a minimum set of information from the present and past such that the future behavior of the system can be completely described by the knowledge of the present and the future input. The state space representation is based on the Markov property, which implies that given the present state, the future of a system is independent of its past. Kalman [8] procedure is the most efficient category of prediction models that have an adaptive behavior. In application of Kalman filtering theory, the mathematical formulation of the problem and the computational techniques involved may depend heavily on the computational simplicities of the system model which is used. Kalman filtering is designed to strip unwanted noise out of the stream of data and it addressed the question of getting accurate information out of inaccurate data. A major practical advantage of the Kalman filter is that the calculations are recursive, so that although the current estimates are based on the whole past history of measurements, there is no need for an ever expanding memory, Chatfield [7].

The Kalman filtering approach means that received data can be processed sequentially rather than as a batch. Kalman [7] defined filtering as any mathematical operation which uses past data or measurements on a given dynamical system to make more accurate statement about present, future or past variables in that system. Kalman has based the construction of the filter in probabilistic theory, more specifically on the conditionally Gaussian properties of random variables. The state space model is reduced as an autoregressive moving average (ARMA) process. Akaike [9] was the first to demonstrate that state space models can be reduced to an ARMA (p, q) model. The relationship between the state space model and its reduced forms gives considerable insight into the potential effectiveness of the different ARMA models, Harvey [10]; Chatfield [7]. ARMA models, typically are parsimonious model, Box and Jenkins [11]; Box $et$ $al.$ [12] and is based on the premise that the autocorrelation function (ACF) and the related statistics can be accurately estimated and are stable over time.

## 2 Recursive Bayesian Estimation

For dynamic state estimation, the discrete time approach is wide spread and convenient with the state vector, $x \in \Re^n$ assumed to evolve according to the system model:

$$x_{k+1} = f(x_k, \omega_k) \tag{2.1}$$

Where $f_k : \Re^n \times \Re^m \to \Re^n$ the system transition is function and $\omega \in \mathbb{R}^n$ is a white noise sequence independent of past and current states. The $pdf$ of $\omega_k$ is assumed known and at discrete time, measurements $y_k \in \Re^p$ become available. These measurements (Ages) are related to the state vector via the observation equation

$$y_k = h_k(x_k, v_k) \tag{2.2}$$

Where $h_k : \Re^n \times \Re^r \to \Re^p$ is the measurement function and $v_k \in \Re^r$ is another white noise sequence of known $pdf$ independent of past and present states and the system noise. The initial $pdf$ of the state vector $p(x_1|D_0) \equiv p(x_1)$ is assumed available together with the functional forms $f_i$ and $h_i$ for $i = 1, 2, ..., k$ the available information at time step $k$ is the set of measurements $D_k = \{y_i : i = 1, 2, ..., k\}$. The requirement in the study is to construct the of the current state given all the available information: $p(x_k|D_k)$.

The $pdf$ in principles are obtained recursively in two stages: prediction and update. Suppose that the required $pdf$ $p(x_{k-1}|D_{k-1})$ at time step $k - 1$ is available, then using the system model it

is possible to obtain the prior of the state at time step $k$:

$$p(x_k|D_{k-1}) = \int p(x_k|x_{k-1})p(x_{k-1}|D_{k-1})dx_{k-1}. \tag{2.3}$$

At time step $k$ the measurement $y_k$ become available and is used to update the prior via Bayes' rule:

$$p(x_k|D_k) = \frac{p(y_k|x_k)p(x_k|D_{k-1})}{p(y_k|D_{k-1})} \tag{2.4}$$

where the normalizing denominator is given by

$$p(y_k|D_{k-1}) = \int p(y_k|x_k)p(x_k|D_{k-1})dx_k \tag{2.5}$$

The conditional *pdf* of $y_k$ given $x_k$, $p(y_k|x_k)$ is defined by the measurement model (2.2) and the known statistics of $v_k$ . In the update equation (2.4), the measurement $y_k$ is used to modify the predicted prior from the previous time step to obtain the required posterior of the state. The recurrence relations (2.3) and (2.4) constitute the formal solution to the Bayesian estimation problem. Equation (2.3) and (2.4) follows the conditional independence assumption between observations. The search for improved approximate implementation procedures for general recursive Bayesian filters has been an active area of research for many years. Many alternative approaches have been suggested and these can be broadly split into three group : analytic approximations where the aim is to use some form of distributional approximation to estimate the *pdf* numerical approximations where a grid of nodes is used as the basis for numerical integration strategies and functional approximations; and most recently Monte Carlo methods where random samples are utilized to effectively give filtering by simulation Gordon *et al* [6]. Analytic solution to this problem is feasible through Kalman filter.

# 3    Analytic Approximation using Kalman Filtering Algorithm

The Kalman filter provides an efficient recursive estimator for the unobserved state of a linear dynamic system from a series of noisy measurements. Kalman filters are based on linear dynamical systems discretized in time domain. They are modeled on a Markov chain built on linear operators perturbed by Gaussian noise. For the linear Gaussian estimation problem, the required probability density function *(pdf)* remains Gaussian at every iteration of the filter. The pair of equations in (2.3) and (2.4) constitute the general form of the state space model. The errors in the measurement (or observation) equation in (2.1) and the state (or transition) equation in (2.2) are generally assumed to be serially uncorrelated and also to be uncorrelated with each other at all time periods.

The prediction in (2.3) of the recursive Bayesian estimation can be approximated as

$$p(x_k|D_{k-1}) = \int p(x_k|x_{k-1})p(x_{k-1}|D_{k-1})dx_{k-1} \approx N^{-1}\sum_{i=1}^{N} p(x_k|x_{k-1}) \tag{3.1}$$

In (3.1), the probabilistic model of the state evolution $p(x_k|x_{k-1})$ is a Markov model and is defined by the system equation (2.1) and the know statistics of $\omega_{k-1}$ The n−dimensional hidden state process $x(k+1)$ follows a first order Markovian dynamics, as it only depends on the previous state at time $k$ and is corrupted by a correlated or uncorrelated state noise process $\omega(k)$ . The Kalman filter relations propagate $f_k$ and $h_k$ as linear, while $\omega_k$ and $v_k$ are updated as additive Gaussian noise of known variance.

The Kalman filter recursively evaluates the estimator of the state vector conditional on the past observations up to time $(k-1)$. This means that only the estimated state from the previous time step and the current measurements are needed to compute the estimate for the current state. In contrast to batch estimation techniques, no history of observation and or estimates is required. Thus the state of the filter is represented by two variables: $\hat{X}_{k|k}$, the a− posterior state estimate at time $k$ given observations up to and including at time $k$; and $P_{k|k}$ the a- posteriori error covariance matrix which is a measure of the estimated accuracy of the state estimate. The predicted a-priori state is given as

$$\hat{X}_{k|k-1} = F_k \hat{X}_{k-1|k-1} + B_k u_k \tag{3.2}$$

And the predicted a- priori estimate covariance:

$$P_{k|k-1} = F_k P_{k-1|k-1} F_k' + Q_k \tag{3.3}$$

Equations (3.2) and (3.3) are the prediction equations. Equation (3.3) follows from standard results on variance-covariance matrices for vector random variables, Chatfield [[7]]. When new observation has been observed, the estimator for $X_k$ can be modified to take account of this extra information. At time $(k-1)$, the best forecast of $y_t$ is given as $H_k \hat{X}_{k|k-1}$ so that the prediction error is given by

$$\xi_t = y_k - H_k \hat{X}_{k|k-1} \tag{3.4}$$

$\xi_t$ in (3.4) is called the prediction error. This quantity can be used to update the estimate of $X_k$ and of its variance-covariance matrix and the best way to do this is by means of the following equation

$$\hat{X}_k = \hat{X}_{k|k-1} + K_t \xi_t \tag{3.5}$$

and

$$P_t = P_{k|k-1} - K_k F_k P - k|k-1 \tag{3.6}$$

where

$$K_t = P_{k|k-1} F_k [\, F_k P_{k|k-1} F_k' + \sigma_v^2\,]^{-1} \tag{3.7}$$

$K_k$ in (3.7) is called the Kalman gain matrix and is a vector of size $(m \times 1)$. Equation (3.5) and (3.6) constitute the second updating stage of the Kalman filter and are called the updating equations. The true state is assumed to be an unobserved Markov process and the measurements are the observed states of a hidden Markov model. Markov processes are an important class of models because they are fairly general and good numerical techniques exist for computational statistics about time evolutions of probability distributions of state variables. The transition probability for events is determined by the Markov chain. The transition probability is a conditional probability for the next state given the current state. The analysis has been performed on the data and the Bayesian Dynamical system modeling implemented in an R package. In order to handle this process within the framework of the classical time series analysis, the observed number of patients with hypertension must be transformed by differencing the process in order to get a stationary process. The transform process is then

$$y_k = (1 - B)^d X_k \tag{3.8}$$

Such a model is called an integrated model because the stationary model that is fitted to the difference data has to be summed or integrated to provide a model for the original non stationary data. Describing the *dth* difference of $X_k$ is said to be an $ARIMA(p, d, q)$ process. In practice, the first differencing is often formal to be adequate to make a series stationary. It may turn out that there is more than one plausible model and based on the use of Akaike information criterion ($AIC$), the goodness of fit of different models is to be compared by assuming that the data are normally distributed. The $AIC$ is defined as

$$AIC = -2\,maximized\ log - likelihood + 2n \approx Tln\sigma^2 + 2n + const \tag{3.9}$$

where $T$ is the length of the observed series after any differencing, $n$ is the number of fitted parameters and $\sigma^2$ is the estimated white noise variance. The model with the smallest value of the $AIC$ is judged to be the most appropriate.

# 4    Applications and Results

Data were collected from the University of Lagos Medical Centre on the total number of patients, and the number of patients with high blood pressure (HBP) within different age groups in each month of the year 2015 and 2016 respectively as shown in Tables 1 and 2. The collated data is graphically displayed in Figures 1, 2 and 3 given below. The total number of patients within a particular age group h is represented by Nh while nh represents the number of patients with hypertension within the age group. The following age groups are considered: below 30, 30-40, 41-50, 51-60, and above 60 years respectively.

Table 1: Number of Patients with Hypertension within Different Age Groups in 2015

| Month | Total Number Of Patients | Number Of Patients With Hypertension | Age Group | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Below 30 | | 30-40 | | 41-50 | | 51-60 | | Above 60 | |
| | | | $N_h$ | $n_h$ | $N_h$ | $n_h$ | $N_h$ | $n_h$ | $N_h$ | $n_h$ | $N_h$ | $n_h$ |
| Jan | 322 | 130 | 29 | 8 | 31 | 13 | 120 | 38 | 128 | 62 | 14 | 9 |
| Feb | 435 | 132 | 6 | 2 | 59 | 13 | 163 | 16 | 177 | 83 | 30 | 18 |
| Mar | 617 | 225 | 13 | 1 | 70 | 15 | 249 | 79 | 239 | 99 | 46 | 31 |
| April | 426 | 150 | 3 | 1 | 56 | 9 | 201 | 69 | 138 | 60 | 28 | 11 |
| May | 448 | 164 | 14 | 2 | 98 | 17 | 167 | 63 | 150 | 71 | 19 | 11 |
| Jun | 432 | 111 | 28 | 4 | 287 | 56 | 68 | 23 | 35 | 21 | 14 | 7 |
| July | 75 | 22 | 8 | 0 | 27 | 7 | 22 | 7 | 13 | 4 | 5 | 4 |
| Aug | 125 | 58 | 5 | 2 | 25 | 7 | 58 | 20 | 32 | 25 | 5 | 4 |
| Sept | 220 | 99 | 10 | 5 | 58 | 25 | 97 | 45 | 52 | 23 | 3 | 1 |
| Oct | 103 | 35 | 3 | 0 | 48 | 15 | 20 | 10 | 22 | 5 | 10 | 5 |
| Nov | 291 | 152 | 18 | 7 | 92 | 35 | 118 | 83 | 57 | 25 | 6 | 2 |
| Dec | 347 | 117 | 12 | 3 | 202 | 45 | 71 | 36 | 42 | 18 | 20 | 15 |

Table 2: Number of Patients with Hypertension within Different Age Groups in 2016

| Month | Total Number Of Patients | Number Of Patients With Hypertension | Age Group | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Below 30 | | 30-40 | | 41-50 | | 51 − 60 | | Above 60 | |
| | | | $N_h$ | $n_h$ | $N_h$ | $n_h$ | $N_h$ | $n_h$ | $N_h$ | $n_h$ | $N_h$ | $n_h$ |
| Jan | 65 | 18 | 3 | 0 | 15 | 1 | 27 | 9 | 15 | 5 | 5 | 3 |
| Feb | 56 | 23 | 2 | 0 | 7 | 2 | 18 | 7 | 26 | 11 | 3 | 3 |
| Mar | 60 | 15 | 4 | 3 | 12 | 1 | 17 | 2 | 23 | 7 | 4 | 2 |
| April | 154 | 41 | 5 | 0 | 29 | 2 | 57 | 11 | 51 | 19 | 12 | 9 |
| May | 249 | 77 | 9 | 2 | 62 | 11 | 88 | 29 | 73 | 26 | 17 | 9 |
| Jun | 230 | 64 | 10 | 2 | 59 | 11 | 76 | 22 | 68 | 23 | 17 | 6 |
| July | 259 | 82 | 15 | 1 | 57 | 10 | 85 | 24 | 86 | 40 | 16 | 7 |
| Aug | 226 | 63 | 14 | 2 | 60 | 7 | 74 | 22 | 64 | 24 | 14 | 8 |
| Sept | 244 | 69 | 28 | 3 | 23 | 7 | 96 | 23 | 82 | 27 | 15 | 9 |
| Oct | 197 | 56 | 6 | 0 | 61 | 8 | 60 | 18 | 63 | 26 | 7 | 4 |
| Nov | 148 | 43 | 5 | 1 | 26 | 4 | 54 | 13 | 54 | 19 | 9 | 6 |
| Dec | 234 | 70 | 22 | 0 | 49 | 14 | 80 | 20 | 67 | 29 | 16 | 7 |

The first step in state space modeling is to find an optimal autoregressive (AR) model that fits the data. The selection of a tentative model is frequently accomplished by matching estimated autocorrelations with the theoretical autocorrelation and partial autocorrelation functions. Table 3 is the ACF, PACF and the AIC of the observed number of patients with hypertension and the correlogram is as in Figures 2 and 3. The R package use the Akaike Information Criterion (AIC)
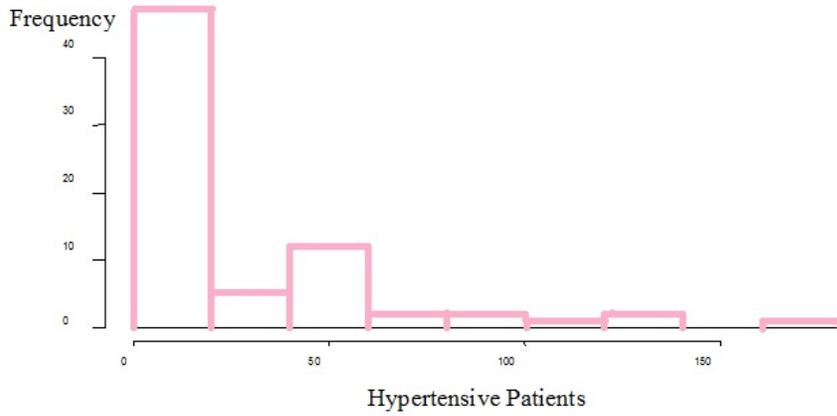
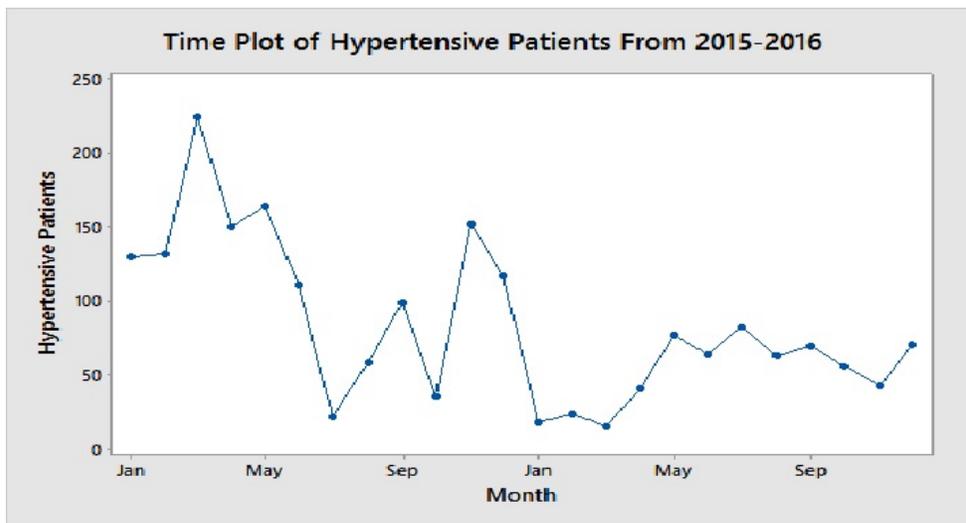Figure 1: Histogram of Hypertensive Patients
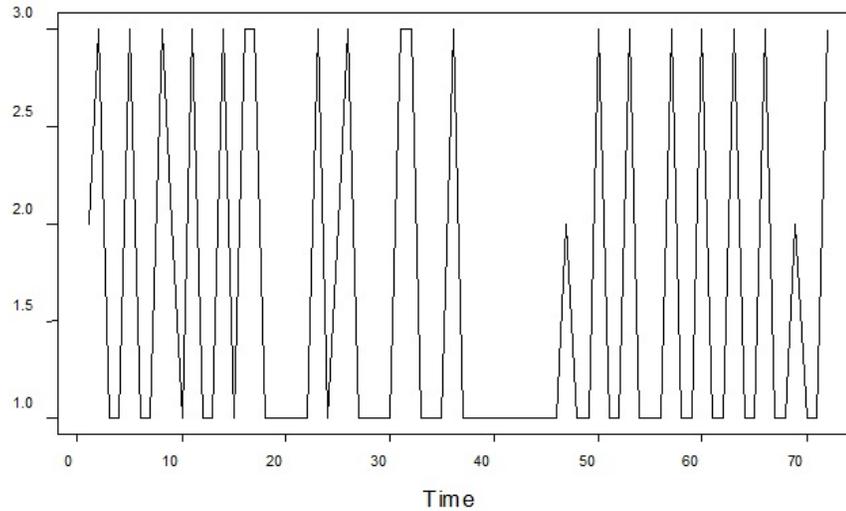


Figure 2: Time Plot of Hypertensive Patients

Figure 3: Stationary Series of Hypertensive Patients

to provide an optimal or best fit for the autoregressive model. The Gauss Markov signal model generated from the observed number of patients with hypertension data using ARMA (p,q) model is

$$\hat{X}_t = 0.478\hat{X}_{t-1} + \omega_t, t \geq 0$$

With mean equal to zero and $\sigma^2_\omega = 0.993$

The Kalman gain $K_k$ as defined in (3.7) is $K_k = 0.0083$. The prediction error variance as defined in (3.4) is $\xi = 0.091$ . The Kalman filter is asymptotically given as

$$\hat{X}_{t|t} = 0.437\hat{X}_{t-1|t-1} + 0.0083y_t$$

Table 3: Sample ACF, PACF and AIC for State Space

| Lag K | ACF | PACF | AIC |
|-------|-------|--------|-------|
| 1 | 0.52 | 0.52 | 0.000 |
| 2 | 0.226 | -0.061 | 1.559 |
| 3 | 0.168 | 0.102 | 2.301 |
| 4 | 0.191 | 0.099 | 3.119 |
| 5 | 0.218 | 0.096 | 4.000 |
| 6 | 0.110 | -0.080 | 5.234 |
| 7 | -0.031 | -0.109 | 5.805 |
| 8 | 0.029 | 0.105 | 6.474 |
| 9 | 0.172 | 0.146 | 5.890 |
| 10 | 0.146 | 0.049 | 7.765 |

Based on the ACF, PACF and AIC of the observed number of patients with hypertension in Table 3, the study recommends an AR (1). The R package use the Akaike's Information Criterion (AIC) to provide an optimal or best fit for the observed number of patients with hypertension. The value of the AIC is minimum at $p = 1$ . The correlogram is on Figures 4 and 5.
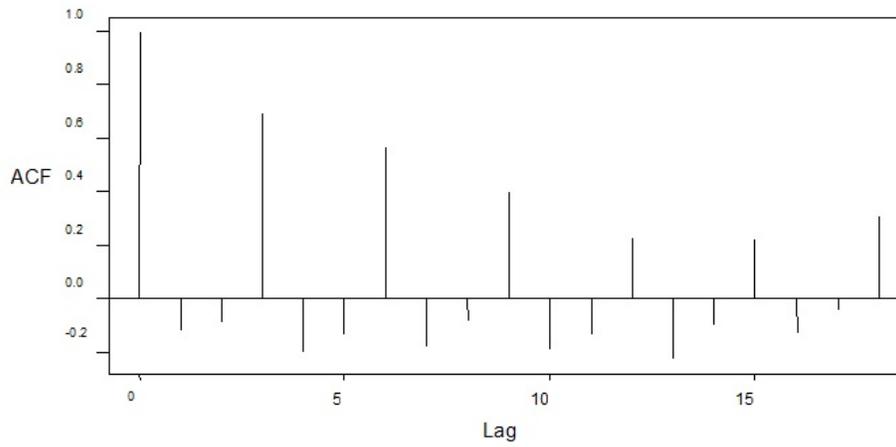
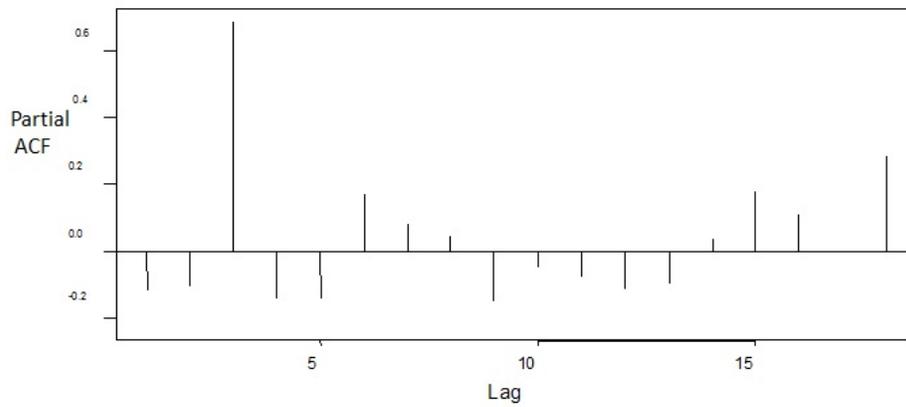Figure 4: ACF of Hypertensive Patients Series



Figure 5: Partial ACF of Hypertensive Patients Series

# 5   Conclusion

Bayesian inference has found application in a wide range of activities, including science, engineering, philosophy, medicine, sport and law; Davison [13]; Jackman [14]; Ogundeji and Okafor [15]. Bayesian analysis is not without problems, however in practice one is forced to establish prior beliefs in the form of prior probability distributions on the model under consideration. For linear Gaussian estimation problem the required probability distribution function remains Gaussian at every iteration of the filter, and the Kalman filter relations propagates and update the mean and covariance of the distribution.

# Competing financial interests

The authors declare no competing financial interests.

# References

# References

[1] Adediran, O. S., Okpara, I. C., Adeniyi, O. S. & Jimoh, A. K. Prevalence of Hypertension in Akwa-Ibom State, South-South Nigeria: Rural Versus Urban Communities Study. *Journal of Medicine and Medical Sciences.* vol. 4, issue 4,pp. 149-154. (2013).

[2] Ekwunife, O. I. & Aguwa, C. N. A Meta analysis of Prevalence Rate of Hypertension In Nigerian Populations. *Journal of Public Health and Epidemiology*, vol. 3, issue 13, pp. 604 - 607.(2011).

[3] Wang, Z., Qingyun, D., Liang S., Nie, K., Lin D., Chen Y. & Li, J. Analysis of the Spatial Variation of Hospitalization Admissions for Hypertension Disease in Shenzhen, China. *Int. J. Environ. Res. Public Health.* vol. 11, pp. 713-733. doi:10.3390/ijerph110100713,(2014).

[4] Adeloyea, D., Basquilla, C., Aderemib, A. V., Jacqueline, Y. Thompsonc, J. Y. & Obid, F. A. An Estimate of the Prevalence of Hypertension in Nigeria: A Systematic Review and Meta-analysis. *Journal of Hypertension.* DOI: 10.1097/HJH.0000000000000413 Source: PubMed,(2014).

[5] Rao, C. R., Kamath, V. G., Shetty A., & Kamath, A. A Quantitative Analysis from Coastal Karnataka, India. *Hindawi Publishing Corporation.* ISRN Preventive Medicine. vol. 2013, Article ID 574973, 6 pages .http://dx.doi.org/10.5402/2013/574973,(2013).

[6] Gordon, N. J., Salmond, D. J. & Smith, A. F. M. A novel approach to nonlinear/ non-Gaussian Bayesian state estimation, *IEE Proceedings on Radar and Signal Processing.* vol. 140, issue 2, pp. 1434-1443, (1993).

[7] Chatfield, C. The Analysis of Times Series An Introduction.6th ed. *Chapman and Hall/CRC*,Boca Raton,(2004).

[8] Kalman, R. E. A new approach to linear filtering and prediction problems, *Journal of Basic Engineering.* vol. 82 pp. 35-45,(1960).

[9] Akaike, H. A New Look at statistical model identification. *I.E.E.E. Transactions on Automatic control*, AU - 19, pp. 716 - 722,(1974).

[10] Harvey, A. C. Forecasting, structural time series models and the Kalman filter. *Cambridge Univ. Press, Cambridge*,(1989).

[11] Box, G. E. P. & Jenkins, G. M . Time series analysis, forecasting and control (rev.ed), *San Franscisco, Holden-Day,*(1976).

[12] Box, G. E. P, Jenkins, G. M, & Remsel, G.C. Time Series Analysis; Forecasting and Control, *Pearson Education*, Delhi, (1994).

[13] Davison, A.C. Statistical Models. *Cambridge University Press*, New York, (2008).

[14] Jackman, S. Bayesian Analysis for the Social Sciences. *Wiley Series in Probability and Statistics*, (2009).

[15] Ogundeji, R. K. & Okafor, R. (2010). Empirical Bayes Approach to Estimate Mean CGPA of University Graduates. *International Journal of Applied Mathematics and Statistics*, vol. 17, issue 10, pp:77-89, (2010).