# State space modelling and forecasting of sugarcane yield in Haryana, India

## Suman[*] and Urmil Verma

[*]Department of Mathematics, Statistics and Physics, CCS Haryana Agricultural University, Hisar - 125004 (Haryana), INDIA
*Corresponding author. E-mail: vermas21@hotmail.com

**Abstract:** Box and Jenkins' Autoregressive Integrated Moving Average (ARIMA) models are widely used for analyzing and forecasting the time-series data. In this approach, the underlying parameters are assumed to be constant however the data in agriculture are generally collected over time and thus have the time-dependency in parameters. Such data can be analyzed using state space (SS) procedures by the application of Kalman filtering technique. The purpose of this article is to illustrate the usefulness of state space models in sugarcane yield forecasting and to provide some empirical evidence for its superiority over the classical time-series analysis. ARIMA and state space models individually could provide the suitable relationship(s) to reliably forecast the sugarcane yield in Karnal, Ambala, Kurukshetra, Yamunanagar and Panipat districts of Haryana (India). However, the state space models with lower error metrics showed the superiority over ARIMA models for this empirical study. The sugarcane yield forecasts based on SS models in the districts under consideration showed good agreement with State Department of Agriculture (DOA) yields by showing 3-6 percent average absolute deviations.

**Keywords:** Autocorrelation function, Kalman filtering technique, State space procedures, Akaike's information criterion, Sugarcane yield forecast

## INTRODUCTION

The importance of agriculture for Indian society can hardly be over emphasized, as its role in economy, employment, food security, self-reliance and general well-being does not need reiteration. India has a very well established system for collection of crop statistics at village level and aggregating it at different administrative levels. The official forecasts (advance estimates) of major cereal and commercial crops are issued by the Directorate of Economics and Statistics, Ministry of Agriculture, New Delhi. However, the final estimates are given a few months after the actual harvest of the crop. Thus, one of the limitations of conventional methods is timeliness and quality of the statistics. Hence, there is a considerable scope of improvement in the conventional system. Timely and in-season crop production forecasting entails judicious planning based on information related to various aspects of agriculture. Thus, information on crop acreage and production are important inputs for strategic planning.

Time series models have advantages in certain situations. They can be used more easily for forecasting purposes because the historical sequences of observations upon study variables are readily available at equally spaced intervals over discrete point of time. These successive observations are statistically dependent and time series modelling is concerned with techniques for the analysis of such dependence. ARIMA forecasts are based only on past values of the variable being forecast. They are not based on any other data series and especially suited to short-term forecasting. The stationarity requirement for the applicability of Box Jenkins' (1976) ARIMA methodology seems to be quite restrictive. The Box-Jenkins procedure for finding a good forecasting model consists of three stages i.e. Identification, Estimation and Diagnostic checking stage(s).

The methods used by the state space procedure also assume the input series to be stationary. Therefore, the first step is to examine the data and test the requirement of differencing. Expositions of the state space approach to multivariate forecasting were observed in the studies of Meinhold and Singpurwala (1983), Kitagawa and Gersh (1984) and Aoki (1987), Jong and Penzer (2004), Brockwell and Davis (2002), Durbin (2002), Bordoloi (2009), Saini and Mittal (2014) etc. At national level, not much work has been done on state space modelling in the field of agriculture. To cite a few more; Stevenson *et al.* (2001), Piepho and Ogutu (2007), Yusof and Kane (2012), Verma *et al.* ( 2015), Yemitan and Shittu (2015), Omekara *et al.* (2016 ) have also given a good account on the application of state space models.

India is one of the largest sugarcane producers in the world, producing around 300 million tonnes of cane

per annum. Sugarcane ranks third in the list of most cultivated crops in India after paddy and wheat (Source: www.mapsofIndia.com/ indiaa-griculture/) (Source: www.mapsofIndia.com/ india-agriculture). Broadly, there are two distinct agro-climatic regions of sugarcane cultivation in India viz., tropical and subtropical. Around 55 per cent of total cane area in the country is in the sub-tropics; Uttar Pradesh, Bihar, Haryana and Punjab comes under this region and Haryana is the largest producer of sugarcane in subtropical region (Source:www.agricoop.nic.in/statistics). Indian planting season of sugarcane in subtropical region falls during September-October to February- March whereas in tropical region, it is January- February to October - November. Keeping in view the above subject matter, the sugarcane yield forecasts of Karnal, Ambala, Kurukshetra, Yamunanagar and Panipat districts have been obtained with the emphasis to see the forecasting performance of the developed models.

## MATERIALS AND METHODS

The study dealt with modeling the time- series sugarcane yield data in Karnal, Kurukshetra, Panipat, Ambala and Yamunanagar districts of Haryana. The sugarcane yield data of State Department of Agriculture for the period 1960-61 to 2014-15 of Karnal and Ambala districts, 1972-73 to 2014-15 of Kurukshetra district and 1989-90 to 2014-15 of Yamunanagar and Panipat districts were compiled from the Statistical Abstracts of Haryana/Punjab (Source: esaharyana.gov.in/ State Statistical Abstract/).

**Box-Jenkins' ARIMA and state space modeling procedures:** The **ARIMA** forecasts are based only on past values of the variable being forecast, however, the data should be available at equally spaced discrete time intervals. The data has to be made stationary before fitting an appropriate ARIMA model. One of the simplest transformations called 'differencing' is applied when the mean of a series changes over time and log transformation is used when the variance of a series changes through time. The two important tools at the identification stage are the estimated autocorrelation function (acf) and partial autocorrelation function (pacf). The estimated acfs i.e. $r_k$ showed the correlation between ordered pairs ( $\overline{Y}_t$ , $\overline{Y}_{t+k}$) separated by various time spans (k = 1,2,3,.....). The estimated pacfs i.e. $\hat{\phi}_{kk}$ showed the correlation between ordered pairs ( $\overline{Y}_t$ , $\overline{Y}_{t+k}$) separated by various time spans (k = 1, 2, 3,…) with the effect of intervening observations ( $\overline{Y}_{t+1}$, $\overline{Y}_{t+2}$, … $\overline{Y}_{t+k-1}$) being accounted for. The functional form of ARIMA (p,d,q) used is expressed as:

$f_p(B) \, \Delta^d \, Y_t = c' + \theta_q(B) \, e_t,$ where $c' = 0$ if $Y_t$ is adjusted for its mean       …. (i)

where $Y$ = Variable under forecasting , $B$ = Lag opera-

tor , $e$ = Error term ($Y - \hat{Y}$ , where $\hat{Y}$ is the estimated value of $Y$), $t$ = time subscript , $f_p(B)$ = non-seasonal AR process, $(1-B)^d$ = non-seasonal difference, $\theta_q(B)$ = non-seasonal MA process, $f$'s and $\theta$'s are the parameters to be estimated (Pankratz, 1991).

At the estimation stage, the precise estimates of a small number of parameters of the model were obtained. Linear least-squares can be used to estimate only pure AR models however non-linear least squares (NLS) method is used for all other models. Finally, the diagnostic tests were performed to check if the random shocks were independent or not.

**The state space model:** represented a univariate time series through auxiliary variables, some of which may not be directly observable. These auxiliary variables, called the state vector summarized all the information from the present and past values of time series relevant to the prediction of future values of the series. The observed time series has been expressed as linear combinations of the state variables.

Let $y_t$ be the $r \times 1$ vector of observed variables after differencing if needed and subtracting the sample mean. Let $z_t$ be the state vector of dimension $s$, $s \geq r$, where the first $r$ components of $z_t$ consist of $y_t$. Various forms of the state space model have been in use but the model fitted with the help of STATESPACE procedure in SAS for this study is based on Akaike (1976). The state space model defined by the state transition equation is

$z_{t+1} = F \, z_t + G \, e_{t+1}$       … (ii)

$z_t$ is a state vector of dimension $s$, whose first $r$ elements are $y_t$ and whose last s-r elements are conditional prediction of future $y_t$. F is an $s \times s$ transition matrix. G is an $s \times r$ input matrix; for model identification, the first $r$ rows and $r$ columns of G are set to an $r \times r$ identity matrix. $e_t$ is a sequence of independent normally distributed random vectors of dimension $r$ with mean 0 and covariance matrix $\Sigma_{ee}$. In addition to the state transition equation, state space models usually include a *measurement equation* or *observation equation* that gives the observed values $y_t$ as a function of the state vector $z_t$.

The measurement equation used by the STATESPACE procedure is

$y_t = H \, z_t$ , H= [$I_r$ 0] and $I_r$ is an $r \times r$ identity matrix       … (iii)

The procedure first fitted a sequence of unrestricted vector autoregressive (VAR) models and computed Akaike's information criterion (AIC) value for each model. The VAR models were estimated using the sample autocovariance matrices and Yule-Walker equations. The order of the VAR model producing the smallest AIC value was chosen as the order (number of lags into the past) to be used in the canonical correlation analysis. The elements of the state vector were determined through a sequence of canonical correla-

tion analysis of the sample autocovariance matrices through the selected order. This analysis computed the sample canonical correlations of the past with an increasing number of steps into the future. Then the state space model was assigned to the data using the Kalman filtering technique.

## RESULTS AND DISCUSSION

**The Box Jenkins' methodology** was applied in obtaining the suitable ARIMA models for district-level sugarcane yield forecasting in Haryana. Autocorrelation functions of sugarcane yield shown in Figure1 indicated that the data series were non-stationary for all the districts under consideration. Differencing of order one was sufficient for making an appropriate stationary series.

The orders of AR and MA components were determined through acfs and pacfs of the stationary series. Marquardt algorithm (1963) was used to minimize the sum of squared residuals. Log Likelihood, AIC (1969), Schwarz's Bayesian Criterion, SBC (1978) and residual variance decided the criteria for the selection/ estimation of AR and MA coefficients in the model. The residual acfs along with the Chi-square test (Ljung and Box, 1978) were used to ascertain the random shocks as white noise.

After experimentation with different lags of moving average and autoregressive processes, ARIMA (0,1,1) for Karnal and Ambala districts and ARIMA (1,1,0) for Kurukshetra, Yamunanagar and Panipat districts were fitted for achieving sugarcane yield forecasts. The fitted ARIMA (0,1,1) and ARIMA (1,1,0) models may be elaborated as below:

$$Y_t = Y_{t-1} - \theta_1 e_{t-1} + e_t \qquad ----- \text{(iv)}$$
$$Y_t = (1+f_1)Y_{t-1} - f_1 Y_{t-2} + e_t \qquad ------\text{(v)}$$

The equations iv & v are the corresponding forecast equations. The presence of lagged values of dependent variable and random shocks in equation-iv indicates the presence of autoregressive and moving average components both. While in equation-v, the presence of lagged values of dependent variable indicates the presence of only autoregressive component. The parameter estimates of fitted ARIMA models are presented in Table 1.

**State space modeling:** The state space model assumes that the time series are stationary. Hence, the data was checked for stationarity. Here, $y_t$, the $r \times 1$ vector of observed variables after differencing and subtracting the sample mean from $Y_t$, can be expressed as follows: The selection of AR orders; five, three, four, five and one seemed reasonable for Karnal, Kurukshetra, Panipat, Ambala and Yamunanagar districts respectively with the use of AIC statistics. Next, the Yule-Walker estimates of the selected AR models were obtained (Table 2). After the autoregressive order selection process of determining the number of lags used in canonical correlation analysis, the state vector was selected.

Information from the canonical correlation and preliminary autoregression analyses were used to form the preliminary parameter estimates of state space models as shown in Table 3.

**The fitted state space models of all the districts can**

| | Karnal | Kurukshetra | Panipat | Ambala | Yamunanagar |
|---|---|---|---|---|---|
| $y_t =$ | $(1-B)Y_t$ - 0.761 | $(1-B)Y_t$ - 0.933 | $(1-B)Y_t$ - 0.443 | $(1-B)Y_t$ - 0.775 | $(1-B)Y_t$ - 0.392 |

**be elaborated as:**

**Karnal**

$$\begin{bmatrix} y_{t+1} \\ y_{t+2/t+1} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -0.125 & -0.113 \end{bmatrix} \begin{bmatrix} y_t \\ y_{t+1/t} \end{bmatrix} + \begin{bmatrix} 1 \\ -0.82 \end{bmatrix} (31.158)$$

……(vi)

**Kurukshetra**

$$\begin{bmatrix} y_{t+1} \\ y_{t+2/t+1} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -0.086 & -0.086 \end{bmatrix} \begin{bmatrix} y_t \\ y_{t+1/t} \end{bmatrix} + \begin{bmatrix} 1 \\ -0.771 \end{bmatrix} (37.39)$$

……(vii)

**Panipat**

$$\begin{bmatrix} y_{t+1} \\ y_{t+2/t+1} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -0.017 & -0.081 \end{bmatrix} \begin{bmatrix} y_t \\ y_{t+1/t} \end{bmatrix} + \begin{bmatrix} 1 \\ -0.686 \end{bmatrix} (31.964)$$

......(viii)

**Ambala**

$$\begin{bmatrix} y_{t+1} \\ y_{t+2/t+1} \\ y_{t+3/t+1} \\ y_{t+4/t+1} \\ y_{t+5/t+1} \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ -0.439 & -0.698 & -0.255 & -0.534 & -0.915 \end{bmatrix} \begin{bmatrix} y_t \\ y_{t+1/t} \\ y_{t+2/t} \\ y_{t+3/t} \\ y_{t+4/t} \end{bmatrix} + \begin{bmatrix} 1 \\ -0.517 \\ 0.035 \\ -0.174 \\ -0.155 \end{bmatrix} (25.512)$$

……(ix)

**Yamunanagar**

$$y_{t+1} = -0.359 \, y_t + 38.684$$

……(x)

The forecasting performance(s) of the contending models were observed in terms of average absolute percent deviations of sugarcane yield forecasts in relation to the observed/DOA yield(s) and RMSEs as well. The sugarcane yield forecasts of the post sample period(s) i.e. 2010-11, 2011-12, 2012-13, 2013-14 and 2014-15 based on ARIMA and SS models were obtained to check the validity of the developed models. The forecast yield(s) along with percent relative deviations of all the districts are presented in Tables 4 &5. However, the root mean square errors (RMSEs) of sugarcane yield forecasts based on both the models are depicted in Table 6.

## Conclusion

ARIMA and state space models individually could provide the suitable relationship(s) to reliably estimate the sugarcane yield of the districts under consideration. However, the state space models with lower error metrics showed the superiority over ARIMA models in capturing the percent relative deviations pertaining to

**Table 1.** Parameter estimates of fitted ARIMA models.

| District/Model | | Estimate | Standard error | p-value |
|---|---|---|---|---|
| Karnal sug-ARIMA (0,1,1) | Constant Difference | 0.67 1 | 0.14 | <0.01 |
| | MA Lag 1 | 0.85 | 0.09 | <0.01 |
| Ambala sug- ARIMA (0,1,1) | Constant Difference | 0.93 1 | 0.20 | <0.01 |
| | MA Lag 1 | 0.74 | 0.10 | <0.01 |
| Kurukshetra sug- ARIMA (1,1,0) | Constant Difference | 0.63 1 | 0.91 | 0.49 |
| | AR Lag 1 | -0.37 | 0.16 | 0.02 |
| Yamunanagar sug- ARIMA (1,1,0) | Constant Difference | 0.79 1 | 1.03 | 0.44 |
| | AR Lag 1 | -0.36 | 0.18 | 0.06 |
| Panipat sug- ARIMA (1,1,0) | Constant Difference | 0.80 1 | 0.80 | 0.32 |
| | AR Lag 1 | -0.48 | 0.17 | 0.01 |

**Table 2.** Yule-Walker estimates of selected AR models.

| Districts | Selected Autoregressive order | | | | |
|---|---|---|---|---|---|
| | Lag=1 | Lag=2 | Lag=3 | Lag=4 | Lag=5 |
| Karnal | -0.806 | -0.709 | -0.535 | -0.306 | -0.284 |
| Kurukshetra | -0.709 | -0.604 | -0.402 | -0.252 | |
| Panipat | -0.601 | -0.408 | -0.336 | -0.484 | |
| Ambala | -0.578 | -0.293 | -0.279 | -0.457 | -0.265 |
| Yamunanagar | -0.359 | | | | |

**Table 3.** Parameter estimates of the state space models.

| Districts | Parameter | Estimate | Standard Error | t-Value |
|---|---|---|---|---|
| Karnal | F(2,1) | -0.125 | 0.180 | -0.69 |
| | F(2,2) | -0.113 | 0.206 | -0.55 |
| | G(2,1) | -0.82 | 0.141 | -5.79 |
| Kurukshetra | F(2,1) | -0.087 | 0.173 | -0.50 |
| | F(2,2) | -0.086 | 0.182 | -0.47 |
| | G(2,1) | -0.771 | 0.153 | -5.04 |
| Panipat | F(2,1) | -0.017 | 0.25 | -0.07 |
| | F(2,2) | -0.081 | 0.292 | -0.28 |
| | G(2,1) | -0.686 | 0.188 | -3.63 |
| Ambala | F(5,1) | -0.439 | 0.164 | -2.68 |
| | F(5,2) | -0.698 | 0.29 | -2.41 |
| | F(5,3) | -0.255 | 0.339 | -0.75 |
| | F(5,4) | -0.534 | 0.339 | -1.58 |
| | F(5,5) | -0.915 | 0.365 | -2.51 |
| | G(2,1) | -0.517 | 0.142 | -3.62 |
| | G(3,1) | 0.035 | 0.161 | 0.22 |
| | G(4,1) | -0.174 | 0.16 | -1.09 |
| | G(5,1) | -0.155 | 0.149 | -1.04 |
| Yamunanagar | F(1,1) | -0.359 | 0.176 | -2.04 |

**Table 4.** Sugarcane yield estimates and their associated percent deviations from the observed yield(s) based on ARIMA models.

Univariate **ARIMA modeling**

| Year | Karnal | | | Ambala | | | Kurukshetra | | | Yamunanagar | | | Panipat | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Obs. yield q/ha | Est. yield (q/ha) | Percent Deviation | Obs. yield q/ha | Est. yield (q/ha) | Percent Deviation | Obs. yield q/ha | Est. yield (q/ha) | Percent Deviation | Obs. yield (q/ha) | Est. yield (q/ha) | Percent Deviation | Obs. yield q/ha | Est. yield (q/ha) | Percent Deviation |
| 2010-11 | 79.77 | 71.52 | 10.34 | 67.22 | 69.75 | -3.76 | 74.16 | 66.09 | 10.88 | 68.29 | 65.16 | 4.58 | 75.99 | 73.9 | 2.75 |
| 2011-12 | 78.38 | 72.2 | 7.88 | 71.58 | 70.7 | 1.23 | 69.93 | 65.64 | 6.13 | 66.02 | 69.34 | -5.03 | 83.72 | 77 | 8.03 |
| 2012-13 | 81.6 | 72.87 | 10.70 | 79.68 | 71.64 | 10.09 | 77.09 | 66.68 | 13.50 | 74.01 | 68.95 | 6.84 | 74.26 | 76.68 | -3.26 |
| 2013-14 | 78.81 | 73.55 | 6.67 | 71.23 | 72.58 | -1.90 | 75.47 | 67.16 | 11.01 | 68.66 | 70.4 | -2.53 | 76.91 | 78.03 | -1.46 |
| 2014-15 | 85.04 | 74.22 | 12.72 | 70.55 | 73.52 | -4.21 | 81.64 | 67.84 | 16.90 | 69.9 | 71.11 | -1.73 | 83.56 | 78.57 | 5.97 |
| Average absolute percent deviation | | | 9.66 | | | 4.24 | | | 11.69 | | | 4.14 | | | 4.29 |

**Table 5.** Sugarcane yield estimates and their associated percent deviations from the observed yield(s) based on State Space models.

Univariate **State Space modeling**

| Year | Karnal | | | Ambala | | | Kurukshetra | | | Yamunanagar | | | Panipat | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Obs. yield q/ha | Est. yield (q/ha) | Percent Deviation | Obs. yield q/ha | Est. yield (q/ha) | Percent Deviation | Obs. yield q/ha | Est. yield (q/ha) | Percent Deviation | Obs. yield q/ha | Est. yield (q/ha) | Percent Deviation | Obs. yield q/ha | Est. yield (q/ha) | Percent Deviation |
| 2010-11 | 79.77 | 72.48 | 9.14 | 67.22 | 66.26 | 1.42 | 74.16 | 72.44 | 2.32 | 68.29 | 68.73 | -0.64 | 75.99 | 74.42 | 2.06 |
| 2011-12 | 78.38 | 74.92 | 4.41 | 71.58 | 70.57 | 1.41 | 69.93 | 74.99 | -7.23 | 66.02 | 65.64 | 0.57 | 83.72 | 79.65 | 4.86 |
| 2012-13 | 81.6 | 76.52 | 6.22 | 79.68 | 73.94 | 7.20 | 77.09 | 73.76 | 4.32 | 74.01 | 67.36 | 8.99 | 74.26 | 81.45 | -9.68 |
| 2013-14 | 78.81 | 78.68 | 0.16 | 71.23 | 68.84 | 3.36 | 75.47 | 75.94 | -0.62 | 68.66 | 70.45 | -2.61 | 76.91 | 80.63 | -4.83 |
| 2014-15 | 85.04 | 79.51 | 6.50 | 70.55 | 69.89 | 0.93 | 81.64 | 75.92 | 7.00 | 69.9 | 69.42 | 0.69 | 83.56 | 82.87 | 0.82 |
| Average absolute percent deviation | | | 5.29 | | | 2.86 | | | 4.30 | | | 2.70 | | | 4.45 |

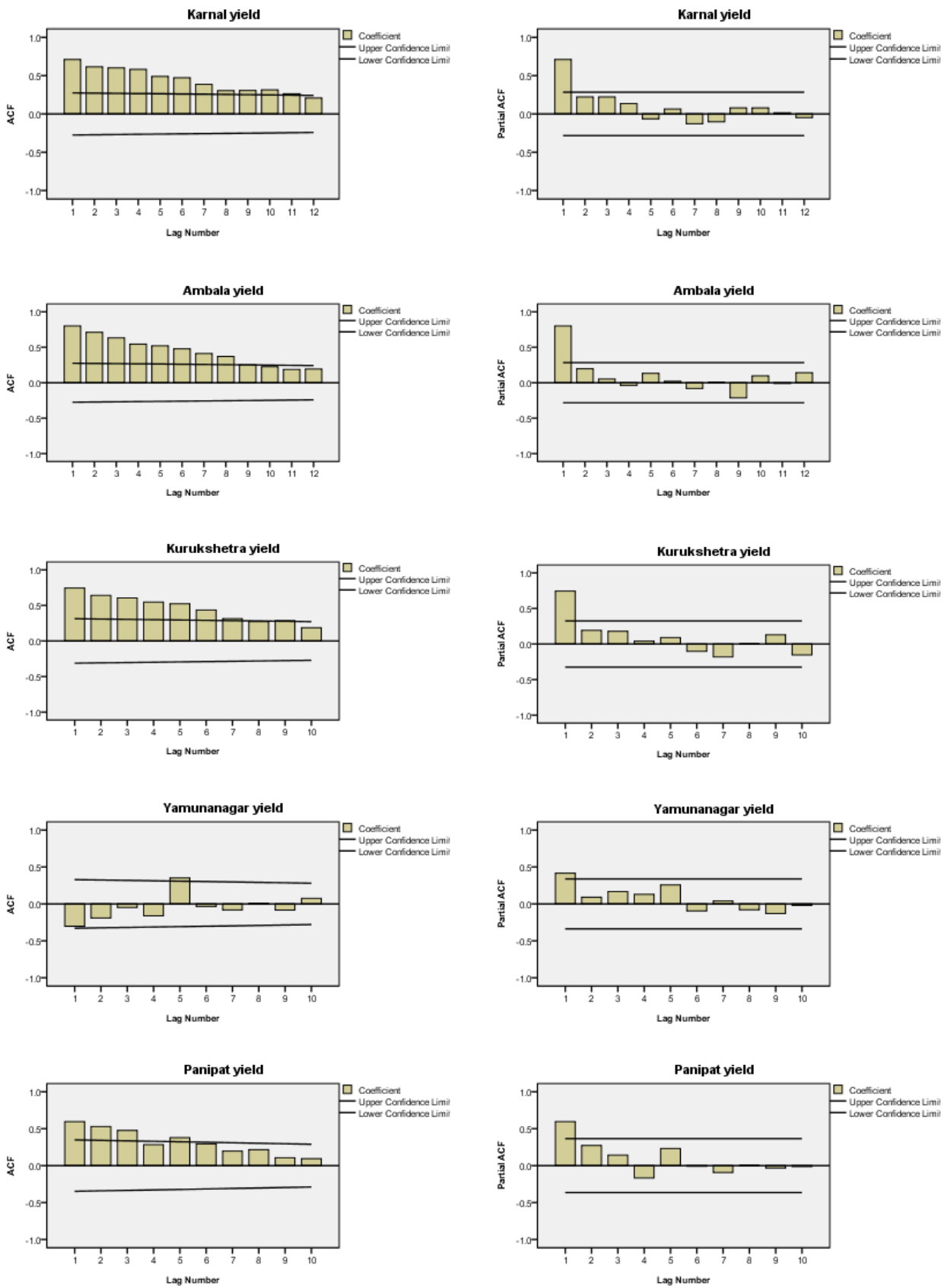Percent Deviation = 100 ' [{observed (obs.) yield – estimated (est.) yield}/ obs. yield]

**Fig. 1.** *Autocorrelation and partial autocorrelation of sugarcane yield for all the districts.*

**Table 6.** Root mean square errors of sugarcane yield forecasts based on alternative models.

| Districts | RMSEs | |
|---|---|---|
| | ARIMA model | SS model |
| Karnal | 8.09 | 4.93 |
| Ambala | 4.06 | 3.69 |
| Kurukshetra | 9.50 | 3.81 |
| Yamunanagar | 3.19 | 3.10 |
| Panipat | 4.04 | 4.12 |

district-level sugarcane yield forecasts in Haryana. The sugarcane yield forecasts based on state space models in the districts under consideration showed good agreement with DOA yield estimates by showing 3-6 percent average absolute deviations. On the basis of this empirical study, it is inferred that the state space modeling may be effectively used pertaining to Indian agriculture data, as it can take into account the time dependency of the underlying parameters which may further enhance the predictive accuracy of the forecast models.

## ACKNOWLEDGEMENTS

## REFERENCES

Akaike, H. (1969). Fitting autoregressive models for prediction. *Ann. Inst. Statist. Math*. 21: 243-247.

Akaike, H. (1976). Canonical correlations analysis of time series and the use of an information criterion in advances and case studies in system identification (R. Mehra and D.G. Lainiotis (Eds.)). Academic Press, New York.

Aoki, M. (1987). State Space modeling of time series. Springer, Berlin.

Bordoloi, S. (2009). Estimation of price level in India through state-space model. *Statistics and Applications*, 7&8(1&2): 17-36.

Box, G.E.P. and Jenkins, G.M. (1976). Time series analysis: Forecasting and control. *Holden Day, San Franscisco.*

Brockwell, P. J., Davis, R. A. (2002). Introduction to time series and forecasting. Springer, *New York*.

Durbin J. (2002). The Foreman Lecture: The State Space approach to time series analysis and its potential for Official Statistics. *Austral. New Zealand J. Statist*. 42: 1 -23.

Jong, P. D. and Penzer, J. (2004). The ARMA model in state space form. *Statistics and Probability letters*, 70(1): 119 -125.

Kitagawa, G. and Gersch, W. (1984). A smoothness priors-state space modeling of time series with trend and seasonality. *J. Amer. Statist. Assoc.* 79 : 378-389.

Ljung, G.M. and Box, G.E.P. (1978). On a measure of lack of fit in time series models. *Biometrika*, 6 : 297-303.

Marquardt, D.W. (1963). An algorithm for least-squares estimation of non-linear parameters. *J. Soc. Ind. Appl. Math*. 2 : 431-441.

Meinhold R.J. and Singpurwalla N.D. (1983). Understanding the kalman filter. *Amer. Statist.* 37 : 123-127.

Omekara, C.O., Okereke, O.E. and Ehighibe, S.E. (2016). Time series analysis of interest rate in Nigeria: A comparison of Arima and state space models. *International Journal of Probability and Statistics*, 5(2) : 33-47.

Pankratz, A. (1991). Forecasting with dynamic regression models. Wiley-Interscience.

Piepho, H.P. and Ogutu, J.O. (2007). Simple state-space models in a mixed model framework. *Amer. Statist.* 61 : 224 -232.

Saini, N and Mittal, A. K. (2014). Forecasting volatility in indian stock market using State Space models. *Journal of Statistical and Econometric Methods*, 3(1): 115-136.

Schwarz, G. (1978). Estimating the dimension of a model. *Ann. Stat.* 62 : 461- 464.

Stevenson, F.C., Knight, J.D., Wendroth, C., Kessel, V.C. and Nielsen, D.R. (2001). A comparision of two methods to predict the landscape-scale variation of crop yield. *Soil and Tillage Research*, 58: 163-181.

Verma, U., Goyal, A. and Goyal, M. (2015). ARIMA versus state space modelling: An application in agriculture. *Adv. Appl. Res.* 7 : 91-95.

Yemitan, R. A. and Shittu, O.I. (2015). Forecasting Inflation in Nigeria by state space modeling. *International Journal of Scientific & Engineering Research*, 6 : 778-786.

Yusof, F. and Kane, I.L. (2012). Modelling monthly rainfall time series using ETS state space and SARIMA models. *International Journal of Current Research*, 4 : 195-200.