

Time series modeling and forecasting of vegetable production in United Arab Emirates

P K Muhammed Jaslam^{1*}, Amogh P Kumar² and Ajay Kumar^{3*}

<https://doi.org/10.5958/2455-7560.2023.00002.X>

¹Department of Forest, Rangeland, and Fire Sciences, University of Idaho 83844 (USA)

²Department of Economics and Finance, University of Canterbury, Christchurch 8140, New Zealand

³Department of Agricultural Statistics, Maharana Pratap Horticultural University 132 116, India

Received: 13 September 2020; Accepted: 19 January 2021

ABSTRACT

Crop forecasting is a formidable challenge. The national and state governments need such predictions before harvesting for various policy decisions relating to storage, distribution, pricing, marketing, import-export etc. In this paper, univariate forecasting models such as random walk, linear trend, quadratic trend, exponential trend, S-curve trend, simple exponential smoothing, Holt's linear exponential smoothing and Autoregressive Integrated Moving Average (ARIMA) models are used to predict vegetable production in the United Arab Emirates. For empirical analysis, a set of 9 different vegetable groups have been considered, contingent upon availability of required data. Annual data from 1974-75 to 2018-19 was used to forecast the next five years since 2019. Suitable models were selected based on the lowest RMSE and minimum of AIC criterion. Model diagnostic checking was done through Runs above and below the median, Runs up and down and Ljung-Box tests on ACF and PACF of residual terms. For onions and green shallots linear trend model was selected as the best fit, whereas simple exponential smoothing model was most suitable in cauliflowers and broccoli, pumpkins, squash and gourds and spinach. The optimum model obtained for forecasting carrots and turnips was Holt's linear exponential smoothing model and ARIMA model was the best fit for the rest of vegetable groups.

KEY WORDS: Forecast models, Vegetable production, AIC, ARIMA, ACF, PACF

Several statistical and econometric forecasting models have been developed in the literature that could be used to forecast various issues, including agricultural production, marketing, demand, trade, etc. (Hanke and Wichern, 2008). Al-Karablieh and Salman (1999) Verma *et al.* (2015), Kumar *et al.* (2019), Naidu *et al.* (2018) etc. are working on various forecasting aspects in agriculture. Fildes and Lusk (1984) advise that forecasters should consider a range of methods and analyze their comparative performance over a random selection of series. Reliable and timely forecasts provide useful and practical advice for effective, foresighted and insightful planning, especially in agriculture, which is full of uncertainties.

Study Area : Food security is at the top of the national agenda in the United Arab Emirates (UAE), but growing crops under the harsh weather conditions of the UAE can be quite a challenge. The landscape of

the United Arab Emirates (UAE) is dominated by low-lying, sandy desert. In the country, about 34 per cent of the area is affected by different levels of salinity, where the growth of healthy plants is almost impossible (Qureshi, 2017).

The UAE imports 80 per cent of its food; this is a significant challenge for the country's food security (Sandhya, 2019). To tackle this hurdle, economical production of food has to be approached at a macro level by examining cross-border efficiencies. Sustainable use of natural resources is a crucial evaluation criterion of modern agricultural production systems. Environmentally-controlled agriculture is a significant source of global agricultural production, especially in the UAE where vegetable consumption rates are going up in addition to the boom of the ornamental plants market (Fadel *et al.*, 2014). For the past four decades, UAE's plant holdings had increased 38-fold from 157 ha in 1971 to 5,935 ha in 2018.

*Corresponding author : ajaystatistics@gmail.com

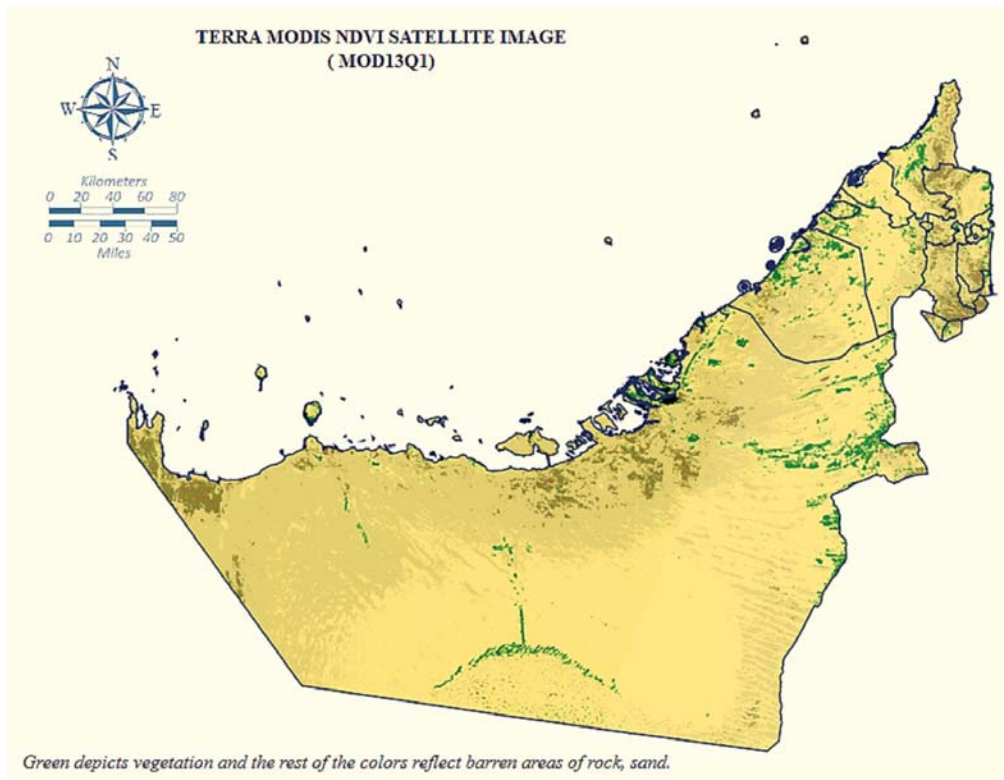


Fig. 1: NDVI image (2019) for the United Arab Emirates

Land classification based on the Normalized Difference Vegetation Index (NDVI) over the UAE is given in Fig. 1. Meanwhile, arable land has expanded by 33-fold from 5530 ha to 1,85,297 ha during the same period. According to the Statistical Book of Abu Dhabi (2019), there are 18,269 greenhouses in Abu Dhabi with 1,415 ha producing vegetables. The smart production of food locally can help alleviate import requirements and as a result, minimize environmental impact in a country that has one of the highest per caput carbon footprints in the world. Duncan (2018) reported high demand for locally-produced fruits and vegetables in UAE supermarkets. Therefore, study was carried out in the United Arab Emirates to find out the trend and forecast the major groups of vegetables using suitable statistical models, which may be useful to the public, researchers, and decision-makers with baseline data.

MATERIALS AND METHODS

The study was carried out by using time series data from 1974-75 to 2018-19, *i.e.* time series data of 45 years. Forecasting of production for next five years of major vegetable groups: (1) cabbages and other brassicas, (2) carrots and turnips, (3) cauliflowers and broccoli, (4) cucumbers and gherkins, (5) eggplant (aubergines), (6) onions and green shallots, (7)

pumpkins, squash and gourds, (8) Spinach and (9) tomatoes. The time-series data were collected from Statistics Division, Food and Agriculture Organization of the United Nations (FAOSTAT) and were analyzed in STATGRAPHICS Centurion 18 software.

Analytical Techniques : The study was tried to fit univariate forecasting models such as random walk, linear trend, quadratic trend, exponential trend, S-curve trend, simple exponential smoothing, Holt's linear exponential smoothing and Autoregressive Integrated Moving Average (ARIMA) models to predict vegetable production. Model diagnostic checking was done through: (1) Runs above and below median (2) Runs up and down and (3) Ljung-Box tests on ACF and PACF of residual terms. Similarly, diagnostic checking can also be done through a minimum of root mean squared error (RMSE) and mean absolute percentage error (MAPE), minimum of Akaike Information Criteria (AIC) and Schwarz Bayesian Criteria (SBC) etc.

Random walk model: It is a non-stationary stochastic time series model also denote as I(1) model. Suppose a_n is a white noise error term with mean 0 and variance σ^2 . Then series Y_n is said to be random walk if

$$Y_n = Y_{n-1} + a_n$$

It means the value of Y (production) at time t is equal to the sum of its value at (n-1) and a random shock.

The above equation can be re-written as

$$Y_n - Y_{n-1} = a_n = \Delta Y_n$$

where, Δ denotes the differencing operator.

Linear trend model

The linear trend is a simple function described as a straight line along with several points of time series value in time series graph and has a typical pattern.

$$Y_n = c + bT_n$$

where c is the constant of production at base period, and b is the coefficient of trend line direction. Method least squares can be applied to find these coefficients.

$$b = \frac{N \sum Y_n T_n - \sum Y_n \sum T_n}{N \sum T_n^2 - (\sum T_n)^2} \text{ and } c = \bar{Y}_n - b\bar{T}_n$$

Non-linear Trend model

In several cases, linear trend was not suitable for time series data. These cases occur when a time series has a different gradient between the beginning phase of the data and the next phase. For these cases, it is better to use a non-linear trend than linear trend. There are several non-linear trends, and in this study, the following models were used:

Quadratic trend $Y_n = c + bT_n + hT_n^2$

Exponential trend $Y_n = cb^{T_n}$

S-curve or Logistic trend $Y_n = \frac{1}{1 + e^{c+bT_n}}$

Non-linear equations can be solved using linearization, Newton Raphson methods etc. see Weisberg (2005).

Exponential smoothing methods

It is a specific kind of moving average technique that is applied to data from time series, used to make a smooth data for projection, or to predict. This method weights preceding observations by diminishing weights exponentially to the prediction of future values.

Simple Exponential Smoothing

It is a process that continually repeats enumeration through the use of the newest data. This approach can be used if trend and seasonal factor do not significantly affect the results. A parameter called the smoothing constant (α) is required to smooth out the data with single exponential smoothing. A convinced weighting is given for each data point, α for the newest data and $(1-\alpha)$ for older data etc. The value of α must be 0 to 1. The following is a smoothed-value equation:

$$S_n = \alpha[Y_n + (1-\alpha)Y_{n-1} + (1-\alpha)^2 Y_{n-2} + \dots]$$

Forecasting value with single exponential smoothing can be done by substituting this equation:

$$\hat{Y}_{n+1} = \alpha Y_n + (1-\alpha)\hat{Y}_n$$

The initial value S_0 can be calculated from the average of several observations. The first several observations can be chosen to determine S_0 .

Double exponential smoothing (Holts)

Holts Method uses different parameters than the one used in the original series. Exponential smoothing prediction can be achieved by using two smoothing constants (with values between 0 and 1) and the following three equations:

$$S_n = \alpha Y_n + (1-\alpha)(S_{n-1} + T_{n-1})$$

$$T_n = \gamma(S_n - S_{n-1}) + (1-\gamma)T_{n-1}$$

$$\hat{Y}_{n+m} = S_n + T_n m$$

The 1st Equation calculates smoothing value S_n from the trend of the previous period T_{n-1} added by the last smoothing value S_{n-1} . Equation 2nd calculates trend value T_n from S_n , S_{n-1} and T_{n-1} . Finally, from equation 3 forward prediction is obtained from trend T_n , multiplied with the amount of next period forecasted m , and added to basic value S_n .

The initial value, *i.e.* S_0 & T_0 , can be estimated with the least-squares method is used. The estimation value for S_0 is the intercept value of linear estimation, while T_0 is the slope value.

Autoregressive Integrated Moving Average methodology (ARIMA)

Univariate Box-Jenkins ARIMA forecasts are based only on past values of the variable being forecast. They are not based on any other data series, and uniquely suited to short-term forecasting. The Box-Jenkins procedure for finding a good forecasting model consists of the following three stages. At the identification stage, two graphical devices estimated ACF and estimated PACF are used to measure the statistical relationships within a data series in a somewhat crude way, but helps in giving a feel for the pattern in the available data. These functions act as a guide for choosing one or more ARIMA models that seem to be appropriate. Whatever model is selected from the identification stage, is merely a tentative candidate for the final model. At the estimation stage, one gets precise estimates of the coefficients of the model chosen at the identification stage based on the available data. At the diagnostic checking stage, the residuals are used to test hypothesis about the independence of the random shocks and to help determine if an estimated model is statistically adequate.

This model is generalized model of the non-stationary ARMA model denoted by ARMA (p, q) can be written as:

$$Y_n = \phi_1 Y_{n-1} + \phi_2 Y_{n-2} + \dots + \phi_p Y_{n-p} + e_n - \theta_1 e_{n-1} - \theta_2 e_{n-2} - \dots - \theta_q e_{n-q}$$

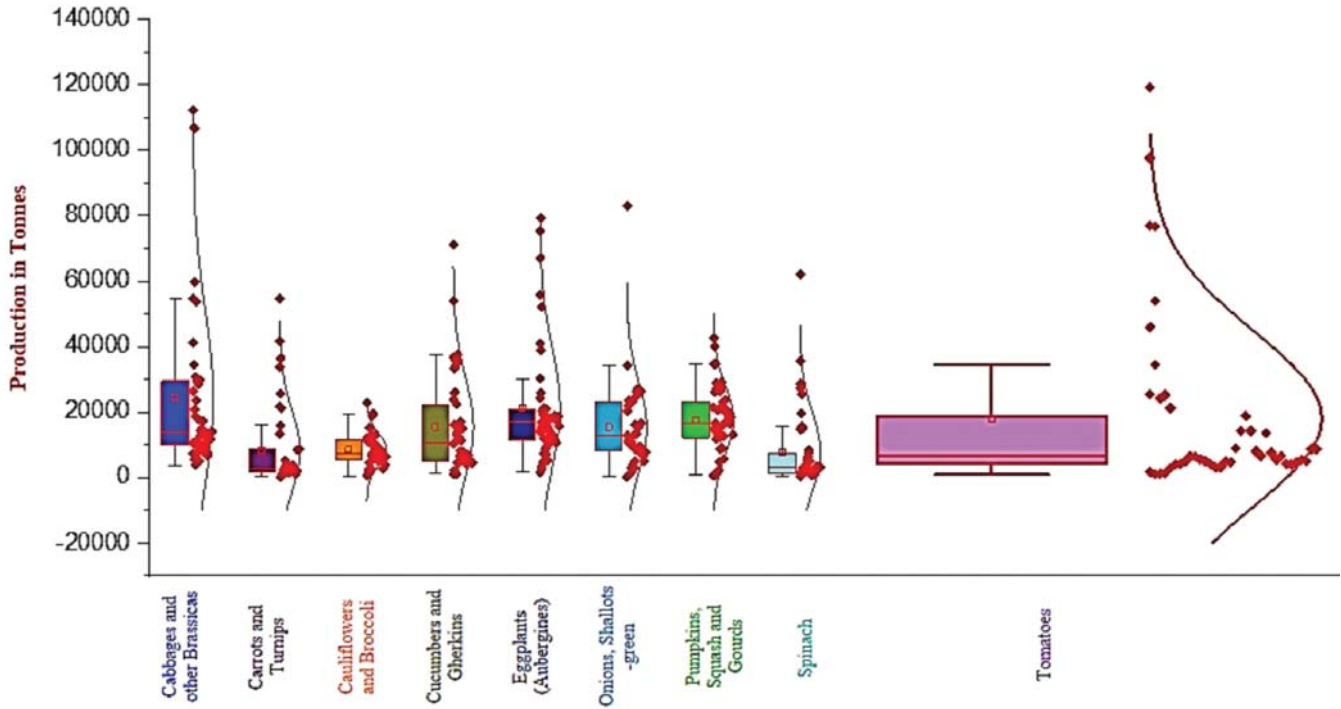


Fig. 2: Box-normal plot of vegetable production in the UAE (1974-75 to 2018-19)

where, Y_n is the original series, for every n , we assume that is independent of $Y_{n-1} + Y_{n-2} + \dots + Y_{n-p}$. A time-series $\{Y_n\}$ is said to follow an integrated autoregressive moving average (ARIMA) model if the d^{th} difference $W_n \nabla^d Y_n$ is a stationary ARMA process. If W_n follows an ARMA (p, q) model, we say that $\{Y_n\}$ is an ARIMA (p, p, q) process. Fortunately, for practical purposes, we can usually take $d = 1$ or at most 2.

Diagnostics checking

Models that are estimated are acceptable only when the residuals are random. For this purpose, several alternative models that may be appropriate were to be fitted. The ACF and PACF of the residuals of these models are then estimated. If the plot of these ACF and PACF exhibit a non-significant pattern, then the corresponding model is valid and can be considered for forecasting. Three standard tests to test the randomness of residuals based on ACF and PACF are: (1) Runs above and below median (2) Runs up and down and (3) Ljung-Box tests. See Box and Jenkins (2008).

To measure the adequacy of the fitted model, RMSE and AIC values are utilized, and it can be computed as follows:

$$RMSE = \sqrt{\frac{1}{N} \sum e_n^2} \quad AIC = 2 \ln(RMSE) + \frac{2k}{N}$$

where, k is the number of estimated model parameters.

RESULTS AND DISCUSSION

Descriptive statistics

The data indicated that there is no stability in the production of all vegetable crops/(variables) over the period. The production of certain variables indicates increasing trends, although non-linear in many cases, with some showing declining trends and others having mixed data (Fig. 4). In time-series language we can say most of the variables are non-stationary in nature and this is reflected in descriptive statistics (Table 1). The descriptive statistics demonstrates the range, minimum, maximum and average values along with other statistical properties. Tomatoes have highest production significantly. It may be attributed to favorable government policies to encourage the production of tomatoes, high acceptance of modern agricultural practices such as hydroponics, greenhouse cultivation etc. Vegetable production shows a subtle pattern as the standard deviation (Std. Dev.) value is too high for all variables. The coefficient of variation (CV) is a useful statistic for comparing the degree of variation from one data series to another, even if the means are drastically different from each other. It is calculated as the ratio of the standard deviation to mean. The coefficient of variation is lesser for Pumpkins, Squash & Gourds, *i.e.* 58 per cent and higher for Tomatoes, *i.e.* 152 per cent, which means that the production of these series is less and more dispersed from the mean values.

Table 1. Descriptive statistics for vegetable production in the UAE (1974-75 to 2018-19)

Variable	Range	Minimum	Maximum	Mean	Std. Dev.	CV	Skewness	Kurtosis
Cabbages and other Brassicas	108383	3836	112219	24541.84	26455.492	107.798	2.309	5.056
Carrots and turnips	54471	129	54600	8136.71	12314.763	151.348	2.294	4.931
Cauliflowers and broccoli	22355	452	22807	8548.98	4841.855	56.637	.825	.948
Cucumbers and gherkins	70218	1132	71350	15419.89	15078.885	97.789	1.770	3.511
Eggplants (aubergines)	77758	1638	79396	21350.87	18220.429	85.338	1.919	3.315
Onions and green shallots	83021	90	83111	15412.71	13559.401	87.975	2.895	13.380
Pumpkins, squash and gourds	42037	589	42626	17544.82	10179.228	58.018	.349	.029
Spinach	61571	406	61977	7931.13	11867.568	149.633	2.808	9.410
Tomatoes	1152720	4700	1157420	170434.82	260048.89	152.580	2.444	5.660

Table 2. Comparison of different time series models based on selection criteria

Variable	Selection Criteria		Forecast Models										
	Random walk	Linear trend	Quadratic trend	Exponential trend	S-curve trend	Simple exponential smoothing	Holt's linear exp.	ARIMA					
Cabbages and other Brassicas	RMSE	13.949	26.728	22.523	27.989	27.419	13.949	14.252	12.472				
	AIC	5.271	6.660	6.362	6.753	6.711	5.315	5.403	5.225				
Carrots and turnips	RMSE	6.326	8.813	5.900	7.844	12.682	5.626	4.940	4.962				
	AIC	3.689	4.441	3.683	4.208	5.169	3.499	3.284	3.293				
Cauliflowers and broccoli	RMSE	3.508	4.848	3.324	5.299	4.405	3.383	3.422	3.412				
	AIC	2.510	3.246	2.536	3.424	3.054	2.482	2.549	2.499				
Cucumbers and gherkins	RMSE	9.136	11.023	10.564	10.781	14.778	9.063	9.115	8.272				
	AIC	4.424	4.889	4.848	4.844	5.475	4.453	4.509	4.404				
Eggplants	RMSE	9.408	10.382	15.562	19.577	17.939	9.408	11.989	7.769				
	AIC	4.483	4.725	5.623	6.038	5.863	4.528	5.057	4.278				
Onions and green shallots	RMSE	12.553	10.760	10.818	12.814	12.164	11.492	11.160	11.210				
	AIC	5.060	4.840	4.896	5.190	5.086	4.928	4.914	4.922				
Pumpkins, squash and gourds	RMSE	5.931	8.532	6.452	11.519	8.172	5.716	5.857	5.758				
	AIC	3.560	4.377	3.862	4.977	4.290	3.531	3.624	3.546				
Spinach	RMSE	12.148	12.001	10.594	12.770	12.528	10.563	10.691	10.685				
	AIC	4.994	5.059	4.854	5.183	5.145	4.759	4.828	4.782				
Tomatoes	RMSE	131.635	260.619	224.928	279.754	272.257	131.637	168.109	125.626				
	AIC	9.760	11.215	10.965	11.357	11.302	9.805	10.338	9.756				

Table 3. Model summary and forecast values of vegetable production in the UAE

Selected model summary					Forecast values (tonnes)				
					2019	2020	2021	2022	2023
Cabbages and other Brassicas									
ARIMA(2,1,2)									
Parameter	Estimate	SE	t	P-value					
AR(1)	0.101	0.141	0.719	0.476	13071	9787	10059	12106	12146
AR(2)	-0.615	0.145	-4.248	0.000					
MA(1)	-0.098	0.070	-1.385	0.174					
MA(2)	-1.004	0.035	-28.823	0.000					
Carrots and turnips									
Holt's linear exp. smoothing					52420	56869	61318	65766	70215
alpha = 0.1276 and beta = 0.9999									
Cauliflowers and broccoli									
Simple exponential smoothing					6312	6312	6312	6312	6312
alpha = 0.7194									
Cucumbers and gherkins									
ARIMA(2,1,2)									
Parameter	Estimate	SE	T	P-value					
AR(1)	0.359	0.111	3.226	0.003	57833	51949	62656	72077	65303
AR(2)	-0.948	0.098	-9.715	0.000					
MA(1)	0.669	0.092	7.275	0.000					
MA(2)	-0.920	0.075	-12.340	0.000					
Eggplants									
ARIMA(2,1,2)									
Parameter	Estimate	SE	T	P-value					
AR(1)	0.476	0.138	3.455	0.001	27079	29406	26807	24182	24481
AR(2)	-0.596	0.141	-4.223	0.000					
MA(1)	0.103	0.040	2.581	0.014					
MA(2)	-0.961	0.038	-25.267	0.000					
Onions and green shallots									
Linear trend = 0.686041 + 0.64029 t									
Parameter	Estimate	SE	t	P-value					
Constant	0.686	3.262	0.210	0.834	30139	30780	31420	32060	32701
Slope	0.640	0.124	5.184	0.000					
Pumpkins, Squash and Gourds									
Simple exponential smoothing					20712	20712	20712	20712	20712
alpha = 0.7448									
Spinach									
Simple exponential smoothing					3142	3142	3142	3142	3142
with alpha = 0.3617									
Tomatoes									
ARIMA(1,0,1)									
Parameter	Estimate	SE	t	P-value					
AR(1)	0.864	0.085	10.200	0.000	66944	57841	49976	43180	37308
MA(1)	-0.299	0.156	-1.916	0.062					

Fig. 2 is the box plot of variables that is the structured way to show the distribution of data based on a five-number summary (minimum, first quartile (Q1), median, third quartile (Q3), and maximum) which provides idea on the variability or dispersion of data.

Identification and estimation of model

The results of fitting different models to the data are compared (Table 2). The model with the lowest value of RMSE and AIC was selected and used to

generate the forecast values. The summary of the chosen model is given (Table 3).

It is noteworthy that variable 6 (onions & green shallots) trails in an incremental linear fashion and linear trend model has been selected as the best fit. This model assumes that the best forecast for future best-fit forecasting model data is given by the linear regression line fit to all previous data. Simple exponential smoothing with alpha values 0.7194, 0.7448

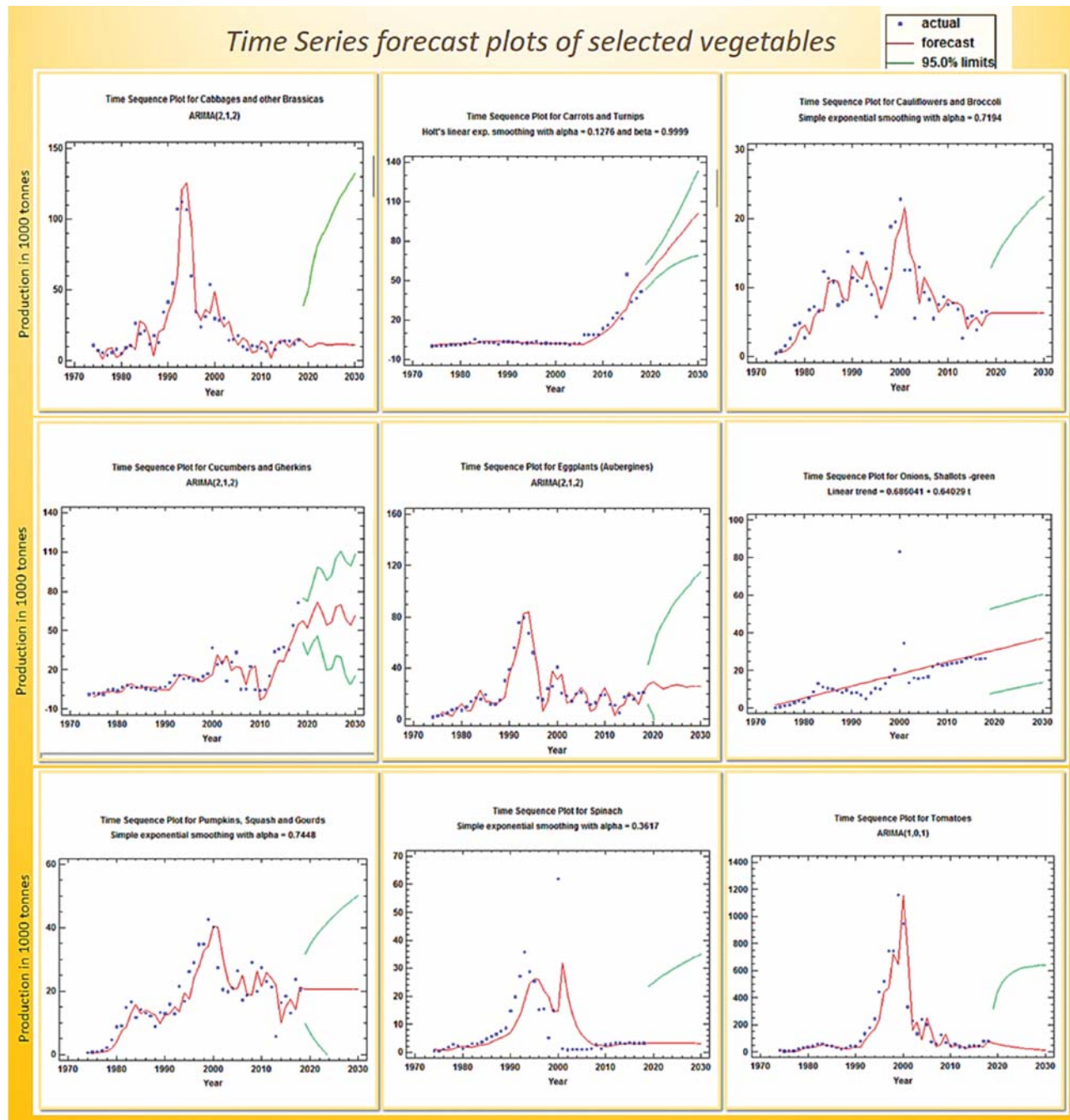


Fig. 3: Time series forecast plots for vegetable production data in the UAE

and 0.3617 fitted best for 3rd variable (cauliflowers and broccoli), 7th variable (pumpkins, squash and gourds) and 8th variable (spinach) respectively. This model assumes future data forecast, is given an exponentially weighted average of all previous data value. Whereas, the forecast model selected for 2nd variable (carrots and turnips) was Holt's linear exponential smoothing with $\alpha = 0.1276$ and $\beta = 0.9999$ and this model assumes that the best forecast for future data is given by a linear trend estimated by exponentially. In the case of the first variable, i.e. cabbages & other brassicas, model ARIMA (2, 1, 2) has been selected with the standard deviation of the input white noise equals 12.6221. Similarly ARIMA (2, 1, 2) also found best fit for variable 4 and 5 (Cucumbers, gherkins and eggplants) with the standard deviation of the input white noise equals to 8.29 and 7.80 respectively. Finally, for the final variable (tomatoes), ARIMA (1, 0, 1) was found to be the best fit with a standard deviation of the input white noise equals to 125.626. The ARIMA model assumes that the best forecast for future data is given by a parametric model relating the most recent data values and previous noise.

Diagnostic checking

The plots of residual normal probability analysis, ACF and PACF are given in the Annexure.

Three tests mentioned in the methods and materials have been run to determine whether or not the residuals form a random sequence of numbers. A series of random numbers is often called white noise since it contains equal contributions at many frequencies. The pooled results indicate that in case of variable 2, 6 and 8, the residuals are not wholly random and that the selected model does not capture all of the structure in the data. In contrast, for rest of the variable, the tests suggest we cannot reject the hypothesis that the series is random at the 95.0% or higher confidence level. Hence in this, forecast value is spoken as a point value without declaring any confidence interval.

Forecasting

The five-year point forecast value obtained by estimating the selected model are reported in Table 3. It is interesting to note that in the case of variables 3, 7 & 8, the chosen models provided a single value forecast for all the forecast years and the forecast value depends solely on the recent production values. Various forecast values with changing trends are seen for the remaining variables and noticeable in figure 3.

CONCLUSIONS

Time series models were found as adequate for forecasting major vegetables of the UAE. Other researchers, people in business, policy-makers, food

producers and many more in the supply chain can use these selected models for information, resource planning and decision-making on vegetable production. The time series modelling for each significant food crop production was appropriate. Based on the forecast trend report, it can be concluded that an increase in production can be expected for carrots, turnips, onions and green shallots. In contrast, almost stagnant output of cabbages and other brassicas, cauliflowers and broccoli, eggplant, pumpkins, squash, gourds and spinach can be expected. It is observed that the production of cucumbers and gherkins is in an undulating pattern. Tomato production shows a declining trend. Also, at the same time, the Box-Jenkins ARIMA model gives a good representation of short-time forecasting. Thus, methodology will encourage other researchers working in the area of vegetable production to develop more efficient and better-grounded forecasting models and techniques.

REFERENCES

- Box G E P and Jenkins G M. 1976. *Time series analysis: forecasting and control*, Holden-Day, University of Michigan.
- Duncan G. 2018. High demand for locally-produced fruit and vegetables, UAE supermarkets say. The National news report. Retrieved on March 04, 2020 from: <https://www.thenational.ae/uae/high-demand-for-locally-produced-fruit-and-vegetables-uae-supermarkets-say-1.743339>.
- Fadel M A, Al Mekhmary M and Mousa M. 2014. Water and energy use efficiencies of organic tomatoes production in a typical greenhouse under UAE weather conditions. *Acta Hort.* **1054**: 81-88.
- Fildes R and Lusk E J. 1984. The choice of a forecasting model. *Omega* **12**(5): 427-35.
- Hanke J E and Wichern D W. 2008. *Business forecasting*, Pearson Education, New Delhi.
- Karablieh E and Salman A. 1999. Forecasting models for barley production in Jordan. *Emir. J. Food Agri.* **11**(1): 59-81.
- Kumar, Deepankar A, Jaslam P K M and Kumar A. 2019. Wheat yield forecasting in Haryana: A time series approach. *Bull. Envir. Pharmaco. Life Sci.* **8**(3): 63-69.
- Naidu G M, Reddy B R and Murthy B R. 2018. Time series forecasting using ARIMA and neural network approaches. *International J. Agric. Stat. Sci.* **14**(1): 275-78.
- Sandhya, D. Sustainable farming on the rise in UAE. August 31, 2019. Khaleej Times news report. Retrieved on March 02, 2020 from: <https://www.khaleejtimes.com/news/general/sustainable-farming-on-the-rise-in-uae>.
- Qureshi A S. 2017. Sustainable use of marginal lands to improve food security in the United Arab Emirates. *J. Exp. Biol. Agri. Sci.* **5**(Spl-1- SAFSAW): 41-49.
- Verma, U and A. Goyal. 2015. Linear mixed modeling for mustard yield prediction in Haryana State (India). *J. Math. Stat. Sci.* pp. 96-105.
- Weisberg, S. 2015. *Applied Linear Regression*, John Wiley & Sons, Inc.