



Modelling and forecasting sorghum (*Sorghum bicolor*) production in India using hierarchical time-series models

SOUMEN PAL¹ and RANJIT KUMAR PAUL²

ICAR–Indian Agricultural Statistics Research Institute, Library Avenue, Pusa, New Delhi 110 012

Received: 1 April 2015; Accepted: 25 December 2015

ABSTRACT

Hierarchical time-series comprises of several dataset maintaining certain hierarchical relationship among them. There are certain specialized strategies, viz. top-down, bottom-up, middle-out and optimal approaches which take care of predicting future values for such multi-level data. For forecasting of individual series at different levels of hierarchy, a method of aggregation or disaggregation is followed. In the present study, these methodologies are investigated thoroughly. Further, state-wise seasonal sorghum [*Sorghum bicolor* (L.) moench] production data of India is analyzed by employing the hierarchical forecasting approaches. A comparative study on performance of different methods is carried out from the viewpoint of multi-step-ahead forecasts on the basis of Mean absolute error (MAE) and Root mean square error (RMSE). The findings show that the middle-out technique outperforms other approaches and traditional method of forecasting as well. This fact has been confirmed statistically by using pair-wise *t*-test. Finally, using the middle-out approach, forecasts of sorghum production for 2015 till 2017 have been carried out at all hierarchical levels. For statistical analysis, R software package has been employed.

Key words: Bottom-up, Forecasting, Hierarchical time-series, Middle-out, Sorghum Production, Top-down

Time series forecasting is an important statistical analysis technique used as a basis for manual and automatic planning in many application domains (Gooijer and Hyndman 2006). Forecasts are calculated using mathematical models that capture a parameterized relationship between past and future values to express behavior and characteristics of a historic time series. The parameters of these forecast models are estimated on a training data set to fit the specifics of the time series by minimizing the forecast error. Time-series data collected in many situations are hierarchical in structure. These dataset generally contain information in clusters which can be combined into another series of interest. Here, the time series are aggregated along the hierarchy based on dimensional attributes such as location (Hyndman *et al.* 2011). Thus, a hierarchical time-series is a collection of several time-series data that are correlated in a hierarchical manner. By contrast, a collection of time-series that are aggregated in a number of non-hierarchical ways, are called a grouped time-series.

In other words, the order by which the series can be grouped is not unique. Forecasting in these environments is especially complex since, it is necessary to involve data and

entities across hierarchical levels and to ensure forecasting consistency among them. This has been the subject of increasing attention recently (Moon *et al.* 2012; Athanasopoulos *et al.* 2009) and has application in diverse fields.

In agriculture, forecasting of crop production is of utmost importance from the view point of policy making. Sorghum [*Sorghum bicolor* (L.) Moench] is the fifth most important cereal crop and is the dietary staple of more than 500 million people in 30 countries. India shares around 20% of the world's sorghum area and is the fourth largest producer of this cereal crop. The cultivation area of sorghum in India was more than 16 million ha in 1981, but has gradually decreased to 6.3 million ha in 2012. Production of this cereal has also faced a decline from 12 million tonnes to 6 million tonnes during this period. Unfortunately, yield of sorghum, on the contrary, has climbed up from 7.3 tonnes/ha to only 9.5 tonnes/ha in the same tenure. However, this crop has much growing potentiality as it is one among the few resilient crops that can adapt well to future climate change conditions, particularly the increasing drought, soil salinity and high temperatures. Thus, to the policy makers, forecasting of sorghum production in this country is certainly a matter of key concern in formulating strategies to enhance the situation. In India, this crop grows in the rainy (*khari*) as well as in the post rainy (*rabi*) season. Forecasting of seasonal sorghum production can incorporate added visibility

¹Scientist (e mail: soumen.4345@gmail.com), Division of Computer Applications; ²Scientist (e mail: ranjitstat@gmail.com), Division of Statistical Genetics,

Under assumption that the base (independent) forecasts are unbiased, it can be written that

$$E[\hat{Y}_n(h)] = E[Y_n(h)]. \tag{3}$$

The necessary condition for the revised hierarchical forecasts to be unbiased is:

$$E[\tilde{Y}_n(h)] = E[Y_n(h)] = SE[Y_{H,n}(h)]. \tag{4}$$

If $\beta_n(h) = E[Y_{H,n+h} | Y_1, \dots, Y_n]$ represents the mean of the forecast values of the bottom level H , then

$$E[\tilde{Y}_n(h)] = SPE[\hat{Y}_n(h)] = SPS\beta_n(h). \tag{5}$$

Therefore, the unbiasedness of the revised forecast will hold if and only if the following condition is satisfied:

$$SPS = S. \tag{6}$$

Bottom-up approach

A frequently applied method for hierarchical forecasting is the bottom-up approach (Dangerfield and Morris 1992, Zellner and Tobias 2000, Espasa *et al.* 2002). This can be represented by the general form of eq. (2), due to Athanasopoulos *et al.* (2009):

$$P = [0_{i \times j} \times I - I_H \mathbf{1}_H]. \tag{7}$$

where $0_{i \times j}$ is the $i \times j$ null matrix. Here, P matrix extracts only bottom-level forecasts from $\hat{Y}_n(h)$. These bottom-level forecasts are then aggregated by the summation matrix S to produce the revised forecasts for the whole hierarchy. By using this approach, the property of unbiasedness in (6) is satisfied. In case of bottom-up approach, the revised forecasts for the bottom level series are equal to the base forecasts.

Top-down approach

Another commonly applied method in hierarchical forecasting is the top-down approach (Narasimhan *et al.* 1994, Fliedner 1999, Dekker *et al.* 2004, Zotteri *et al.* 2005, Widiarta *et al.* 2007). This approach involves first generating base forecasts for the *Total* series on the top of the hierarchy and then disaggregating these downwards based on the proportions of the data. Once the bottom level forecasts have been generated, the summing matrix S can be used to generate forecasts for the rest of the series in the hierarchy. It is evident that for top-down approaches the top level revised forecasts are equal to the top level base forecasts. In terms of the general form of eq. (2), due to Athanasopoulos *et al.* (2009), it can be written that

$$P = [p \mathbf{0}_{1 \times (l-1)}]. \tag{8}$$

where $p = [p_1, p_2, \dots, p_{l-1}]'$ are a set of proportions for the bottom level series. So, the role of P here is to distribute the top level forecasts for forecasting the bottom level series.

Middle-out approach

In this approach, forecasts are produced at an intermediate level of hierarchy, and then disaggregated to obtain forecasts at lower levels and aggregated for higher level forecasts. Thus, the middle-out approach is a

combination of bottom-up and top-down approaches. In practice, production houses apply this method to study demand forecasting (Lo *et al.* 2008). At first, base forecasts are produced for all the series of the selected middle level. Then, for the series above the middle level, revised forecasts are constructed using the bottom-up approach by aggregating the middle-level base forecasts upwards. For the series below the middle level, revised forecasts are generated using a top-down approach by disaggregating the middle level base forecasts downwards.

Optimal combination approach

The optimal combination approach due to Hyndman *et al.* (2011), involves first generating independent base forecast for each series in the hierarchy and, provided the base forecasts are unbiased, produces unbiased revised forecasts which are consistent across the levels of the hierarchy. The h -step-ahead base forecasts can be written by the linear regression model as

$$\hat{Y}_n(h) = S\beta_n + \epsilon_h \tag{9}$$

where $\beta_n = E(\hat{Y}_{H,n}(h) | Y_1, \dots, Y_n)$ is the unknown mean of the base forecasts of the bottom level H and ϵ_h and has zero mean and covariance matrix $Var[\epsilon_h] = \Sigma_h$. Provided the Σ_h is known, the generalised least squares estimation procedure can be used to obtain the minimum variance unbiased estimate of β_n . This can be written as

$$\beta_n(h) = (S'\Sigma_h^{-1}S)^{-1} S'\Sigma_h^{-1}\hat{Y}_n(h) \tag{10}$$

where Σ_h^{\dagger} is the Moore-Penrose generalized inverse of Σ_h . This leads to the following revised forecasts

$$\tilde{Y}_n(h) = S\hat{\beta}_n(h) = SP\hat{Y}_n(h) \tag{11}$$

where $P = (S'\Sigma_h^{\dagger}S)^{-1} S'\Sigma_h^{\dagger}$. This fulfils the unbiasedness property of (6).

RESULTS AND DISCUSSION

For hierarchical time-series data under this study, modeling and forecasting have been done using *hts* package (Hyndman *et al.* 2013) in R software. Fig 2 and 3 illustrate the time-series data of sorghum production across all hierarchical levels. The bottom-most level, i.e. level 2 exhibits time-series of seasonal sorghum production for each of the states, however, for clarity, (seasonal) production of some of the major growing states are only illustrated in Fig 3. Level 1 shows production of sorghum in *kharif* and *rabi* seasons separately in India. Time-series of total sorghum production in India is displayed in level 0. For each of the graphical representation, Y-axis represents production in '000 tonnes and x-axis indicates the time period (year).

An optimal Autoregressive integrated moving average (ARIMA) model (Box *et al.* 2007) using the automatic algorithm of Hyndman and Khandakar (2008) has been employed for each series. The parameters of the model are re-estimated using a moving window beginning with the model fitted using the first 44 (1963-2006) observations. Forecasts from the fitted model are produced for up to 6

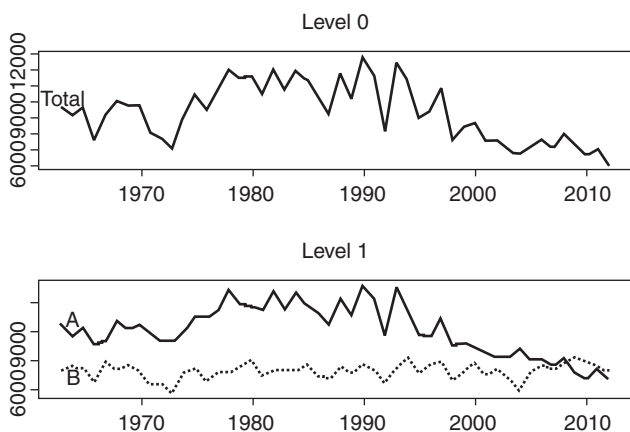


Fig 2 Hierarchical time-series of sorghum production at level 0 and 1

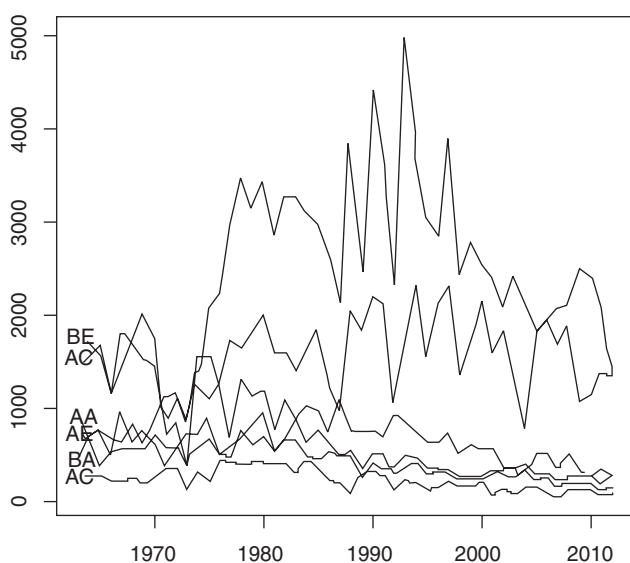


Fig 3 Hierarchical time-series of sorghum production at level 2 for selective states

steps ahead. This process is iterated, increasing the sample size by one observation until 2011 and thus produces 6 one-step-ahead, 5 two-step-ahead, and up to 1 six-step-ahead forecast. This rolling window approach is applied to evaluate the out-of-sample forecast performance of the hierarchical methods under consideration. Independent forecasting using ARIMA has also been included for each time-series separately for comparison purpose.

Using eq. (12) and (13), MAE and RMSE are obtained by employing different approaches for each forecast horizon (h) and presented in Table 3 and 4. In Table 3, for a particular h, under any approach, each MAE value is the average of total MAE of all the series. The same is applied for Table 4 also. The level for middle-out method has been set here as 1.

$$MAE = 1/h \sum_{i=1}^h |y_{t+i} - \hat{y}_{t+i}| \tag{12}$$

$$RMSE = [1/h \sum_{i=1}^h \{(y_{t+i} - \hat{y}_{t+i})^2\}]^{1/2} \tag{13}$$

For each method, the final columns of Table 3 and 4 labeled “Average” show average MAE and RMSE,

respectively, across all the forecast horizons. The bold entries identify the minimum among all the average values. It is evident from the above two tables that middle-out method produces best out-of-sample forecasts followed by top-down approach.

Pairwise t-tests have also been performed between each pair of methods to examine whether differences in the MAE (or RMSE) values are statistically significant. In case of MAE, pairwise t-values between each pair of last 3 methods, viz. middle-out, optimal and independent are statistically significant. Thus, employing the same ARIMA methodology, middle-out hierarchical forecasting approach can significantly perform better than the optimal and independent methods. However, for RMSE, instead of the minimum average value produced in case of middle-out approach, none of the t-values are statistically significant. As the middle-out method outperforms other methods for this particular dataset, the final forecasting of sorghum production for the year 2015 till 2017 has been produced by using this approach. Table 5 depicts the forecasted values for all the series at each level. It has been observed that forecasts obtained at bottom level for *kharif* sorghum production in most of the states and India (level 1) as well

Table 2 Hierarchical structure

		Top level	
1	Total	India	
		Level 1: Season	
2	A	Kharif	
3	B	Rabi	
		Level 2: State/Union territory	
4	AA	Kharif-Andhra Pradesh	
5	AB	Kharif-Bihar*	
6	AC	Kharif-Gujrat	
7	AD	Kharif-Haryana	
8	AE	Kharif-Karnataka	
9	AF	Kharif-Madhya Pradesh**	
10	AG	Kharif-Maharashtra	
11	AH	Kharif-Orrisa	
12	AI	Kharif-Rajasthan	
13	AJ	Kharif-Tamil Nadu	
14	AK	Kharif-Uttar Pradesh	
15	AL	Kharif-Delhi	
16	AM	Kharif-Others***	
17	BA	Rabi-Andhra Pradesh	
18	BB	Rabi-Gujrat	
19	BC	Rabi-Karnataka	
20	BD	Rabi-Madhya Pradesh	
21	BE	Rabi-Maharashtra	
22	BF	Rabi-Tamil Nadu	

* (including Jharkhand), ** (including Chhattisgarh), *** (Comprising of Dadra and Nagar Haveli, Jammu and Kashmir, Kerala, Nagaland, Puducherry, Punjab and West Bengal whose contributions are either nil in some of the years or insignificant compared to total sorghum production of India).

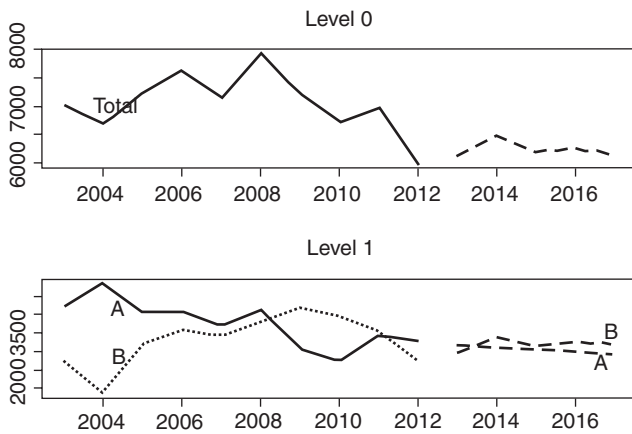


Fig 4 Hierarchical forecasting of sorghum production at level 0 and 1

Table 3 MAE for each forecast horizon

Methods	MAE						
	Forecast horizon (h)						Average
	1	2	3	4	5	6	
Bottom-up	186.74	144.96	136.48	150.63	154.82	157.85	155.25
Top-down	151.36	96.74	133.38	178.00	152.27	153.77	144.25
Middle-out	113.29	105.23	110.16	156.83	151.92	155.11	132.09
Optimal	167.80	124.03	131.87	169.95	157.26	158.10	151.50
Independent	160.27	124.88	125.90	162.67	151.94	154.25	146.65

Table 4 RMSE for each forecast horizon

Methods	MAE						
	Forecast horizon (h)						Average
	1	2	3	4	5	6	
Bottom-up	186.74	166.83	175.04	174.31	174.37	179.70	176.17
Top-down	151.36	121.86	158.09	194.73	172.45	176.25	162.46
Middle-out	113.29	114.15	126.36	175.09	170.28	177.46	146.10
Optimal	167.80	145.94	167.40	189.35	177.08	181.40	171.49
Independent	160.27	144.12	155.18	182.67	177.08	178.16	166.25

are declining gradually as the year progresses. However, this is not applicable in case of state-wise production in *rabi* season as the forecasted values are fluctuating for the coming years. A similar situation has also been observed for overall sorghum production of India in *rabi* season at first level of hierarchy and total national production at topmost level as well.

Fig 4 and 5 illustrate forecasting of sorghum production across all hierarchical levels for each series by using middle-out approach.

The bottom-most level, i.e. level 2 exhibits, through dotted lines, forecast value of seasonal sorghum production for each of the states, however, for clarity of graphical representation, only few of those are labeled. Level 1 shows forecasted production of sorghum in *kharif* and *rabi* seasons separately for India. Predicted future values of total sorghum

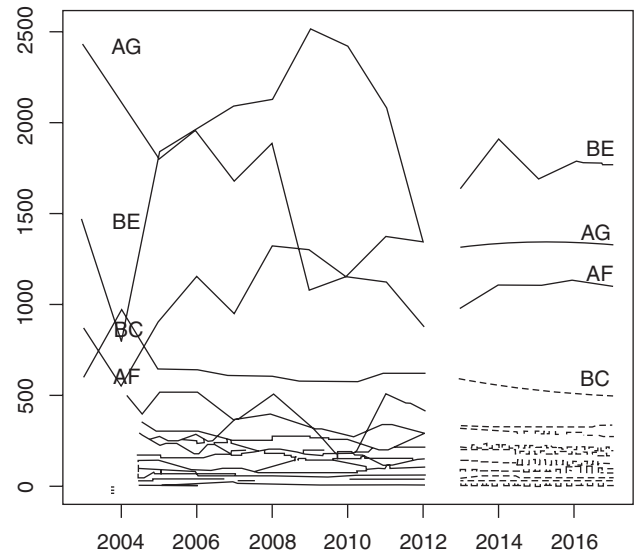


Fig 5 Hierarchical forecasting of sorghum production at level 2

Table 5 Forecasting of sorghum production at all levels for 2015-2017

Level		2015	2016	2017
Top level				
1	Total	6188.44	6257.35	6124.05
Level 1				
2	A	3051.93	2991.52	2931.12
3	B	3136.51	3265.82	3192.93
Level 2				
4	AA	119.50	108.39	97.28
5	AB	2.94	2.88	2.87
6	AC	80.28	76.38	72.49
7	AD	30.33	29.98	29.63
8	AE	293.89	284.73	275.60
9	AF	532.25	512.66	493.13
10	AG	1338.82	1334.01	1329.40
11	AH	5.53	5.52	5.52
12	AI	326.72	326.90	327.13
13	AJ	132.90	127.80	122.40
14	AK	182.14	175.42	168.72
15	AL	3.33	3.31	3.30
16	AM	3.30	3.53	3.65
17	BA	215.04	210.49	194.35
18	BB	51.41	50.34	46.50
19	BC	1100.41	1131.02	1098.46
20	BD	2.72	2.45	2.05
21	BE	1681.28	1780.05	1766.38
22	BF	85.65	91.46	85.18

production in India are displayed in level 0. For each of the graphical representation, y-axis represents production in '000 tonne.

In this paper, a new forecasting strategy has been discussed for hierarchical time-series. State-wise seasonal sorghum production data from 1963 till 2012 form the

lowest level of the hierarchy under study. The hierarchical forecasting approaches, viz. top-down, bottom-up, middle-out and optimal are used to examine the out-of sample forecasting performance for individual method. A traditional forecasting approach has also been employed for each series independently. It has been observed that the middle-out method outperforms the conventional approach of forecasting for the particular dataset and the result is confirmed using the pairwise *t*-test. Finally, the forecasting is made for sorghum production from 2015 till 2017 at all levels of hierarchy i.e. state-wise seasonal, seasonal nation wise and for the whole country. A decreasing future trend of sorghum production in *kharif* season has been observed for most of the growing states. This is affecting the production of sorghum at national level as well. This result confirms the scenario of decreasing area and production of this crop at a much faster rate than the increasing productivity in Indian context. This may be a matter of concern for the policy makers of our country. Finally, the approaches advocated here are very general and can be used for forecasting of hierarchical time series data of any crop as well.

REFERENCES

- Athanasopoulos G, Ahmed R A and Hyndman R J. 2009. Hierarchical forecasts for Australian domestic tourism. *International Journal of Forecasting* **25**(1): 146–66.
- Box G E P, Jenkins G M and Reinsel G C. 2007. *Time-Series Analysis: Forecasting and Control*, 3rd edition. Pearson Education, India.
- Dangerfield B J and Morris J S. 1992. Top-down or bottom-up: Aggregate versus disaggregate extrapolations. *International Journal of Forecasting* **8**(2): 233–41.
- Dekker M, Donselaar K V and Ouwehand P. 2004. How to use aggregation and combined forecasting to improve seasonal demand forecasts. *International Journal of Production Economics* **90**(2): 151–67.
- Espasa A, Senra E and Albacete R. 2002. Forecasting inflation in the European Monetary Union: A disaggregated approach by countries and by sectors. *European Journal of Finance* **8**(4): 402–21.
- Fliedner G. 1999. An investigation of aggregate variable time-series forecast strategies with specific subaggregate time-series statistical correlation. *Computers and Operations Research* **26**(10-11): 1 133–49.
- Gross C and Sohl J. 1990. Disaggregation methods to expedite product line forecasting. *Journal of Forecasting* **9**(3): 233–54.
- Gooijer J G D and Hyndman R J. 2006. 25 years of time series forecasting. *International Journal of Forecasting* **22**(3): 443–73.
- Hyndman R J, Ahmed R A, Athanasopoulos G and Shang H L. 2011. Optimal combination forecasts for hierarchical time-series. *Computational Statistics and Data Analysis* **55**(9): 2 579–89.
- Hyndman R J, Ahmed R A and Shang H L. 2013. hts: Hierarchical and grouped time-series. R package version 4.3, URL <http://CRAN.R-project.org/package=hts>.
- Hyndman R J and Khandakar Y. 2008. Automatic time-series forecasting: The forecast package for R. *Journal of Statistical Software* **27**(3): 1–22.
- Lo S, Wang F and Lin J T. 2008. Forecasting for the LCD monitor market. *Journal of Forecasting* **27**(4): 341–56.
- Moon S, Hicks C and Simpson A. 2012. The development of a hierarchical forecasting method for predicting spare parts demand in the South Korean Navy-A case study. *International Journal of Production Economics* **140**(2): 794–802.
- Narasimhan S L, McLeavey D W and Billington P. 1994. *Production Planning and Inventory Control*, 2ndEdn. Allyn & Bacon, USA.
- Widiarta H, Viswanathan S and Piplani R. 2007. On the effectiveness of top-down approach for forecasting autoregressive demands. *Naval Research Logistics* **54**(2): 176–88.
- Zellner A and Tobias J. 2000. A note on aggregation, disaggregation and forecasting performance. *Journal of Forecasting* **19**(5): 457–65.
- Zotteri G, Kalchschmidt M and Caniato F. 2005. The impact of aggregation level on forecasting performance. *International Journal of Production Economics* **93-94**: 479–91.