

Forecasting grouped time series demand in supply chains

Dejan Mircetic^a, Bahman Rostami-Tabar^b *, Svetlana Nikolicic^a, Marinko Maslaric^a,

^a*University of Novi Sad, Trg Dostiteja Obradovica 6, 2100 Novi Sad, Republic of Serbia*

^b*Cardiff Business School, Cardiff University, Aberconway Building, Colum Drive, Cardiff CF10 3EU, UK*

ABSTRACT

Purpose

Decisions are made at different levels in a supply chain hierarchy. Different decisions also require forecasts at various levels of supply chains. Grouped time series forecasting approaches indeed can provide forecasts to inform decisions at different part of supply chains. Existing approaches often involve bottom-up, top-down and optimal combination approaches. The purpose of this paper is to evaluate the effectiveness of grouped time series approaches in supply chain forecasting and to investigate conditions under which each approach provides forecasts that are more accurate using historical time series demand of a supply network from a major European brewery.

Design/methodology/approach

We evaluate the effectiveness of grouped time series forecasting approaches in supply chains using a simulation and an empirical study. We also develop a regression model to investigate the effect of time series characteristics on the forecast performance of each approach.

Findings

We find that optimal combination and bottom-up methodologies provide forecasts that are more accurate and robust across various levels of supply chain. Moreover, the results show that forecast combination of approaches are as accurate as the optimal approach. We also observe that the presence of entropy and linearity in the bottom series improve the forecast accuracy of approaches.

Practical implications

Our paper should be of interest to practitioners who need to generate forecasts at various levels of supply chains. It will help them to understand the superiority conditions of each approach and their combinations. Therefore, It will assist them in selecting an appropriate methodology for their forecasting task.

Originality/value

This paper contributes to the supply chain forecasting by evaluating the effectiveness of grouped time series forecasting approaches and their combinations on the forecast accuracy using a supply chain network of a brewery company. It offers insights into the choice of grouped time series approaches in supply chains and provides empirical evidences on their superiority conditions.

Keywords: Supply chain forecasting, Forecast combination, Grouped time series forecasting, Time series characteristics.

1. INTRODUCTION

Demand forecasting is the starting point for most planning and control organizational activities (Rostami-Tabar et al., 2015). It is vital to supply chains (SCs), as it provides the basic inputs for the planning and control of all functional areas, including logistics, marketing, production, and finance (Ballou, 2004). Demand forecasting performance is subject to the uncertainty underlying the time series demand a SC is dealing with (Rostami-Tabar, 2013). Therefore, capturing, managing and characterizing uncertainty in SCs represents one of the main problems confronting managers while planning and synchronizing operations in SCs. The demand uncertainty is among the most important challenges facing modern SCs (Syntetos et al., 2016; Chen and Blue, 2010; Mircetic et al., 2017; Teunter et al., 2018; Babai et al., 2012; Trapero et al., 2019; Kalchschmidt et al., 2006), and it poses considerable difficulties in terms of the SC planning and control (Syntetos et al., 2016). Hence, the purpose of demand forecasting in SCs is to inform SC planning decisions by providing an accurate estimation of the future demand in a given situation.

Demand forecasting for SC often concerns many items. SC forecasters may extrapolate the time series for each Stock Keeping Unit (SKU) individually. However, most of the SC time series have natural groupings of SKUs; that is, the SKUs may be aggregated to get higher levels of forecasts across different dimensions such as geographical areas, customer types, supplier types and product families (Chen and Boylan, 2007). Therefore, various levels of forecasts are required for different parts of the SC. The level at which forecasting is performed then it will depend on the function the forecasts are fed into (Rostami-Tabar et al., 2015). For instance, a retailer may use the point-of-sale data to produce forecasts at the store level. However, a manufacturer may use the forecasts of aggregated demand series for the production planning (Chopra and Meindl, 2007). For the transportation manager in charge of the distribution planning, crucial information may include spatial fragmentation of demand, shipments size (replenishment orders) for each distribution channel, a timing of the shipments and the type of product in shipments. Inventory manager might be interested in forecasts related to the type of materials needed at the SKU level, how much will be needed, and when it will arise (Caplice and Sheffi, 2006). Therefore, SC managers may disaggregate

the total demand on the dimensions which are important for a particular party in the chain. This is the moment where the forecasting a single SKU is directed toward the hierarchical forecasting (HF) or grouped forecasting (GF)^a.

A considerable part of the forecasting literature has been dedicated to methods for single time series, but in reality, there are often many related time series that can be organized hierarchically or in groups (Rostami-Tabar et al., 2015). Hierarchical time series can be represented as a hierarchically organized multiple time series that may be aggregated at several different levels in groups according to the different features (Hyndman et al., 2011). Grouped time series are hierarchical time series that do not impose a unique hierarchical structure in the sense that the order by which the series can be grouped is not unique (Hyndman and Athanasopoulos, 2018).

HF naturally reflects important SC characteristics and offers an ample scope for the introduction of innovative forecasting methodologies (Syntetos et al., 2016), improvement of the forecast accuracy and planning, reduction of the overall forecast burden and delivery of the high service level (Strijbosch et al., 2008; Caplice and Sheffi, 2006). Existing approaches for forecasting hierarchical and grouped times series may involve bottom-up (BU), top-down (TD) and optimal combination (OC) approach. Therefore, in this paper we may use GF or HF approaches to refer to BU, TD and OC approaches. In the TD approach, the univariate forecast is generated at the top level of the forecasting structure and then disaggregated to the bottom level series. Oppositely, the BU generates the multiple univariate forecasts in the bottom level of the forecasting structure and then aggregates these forecasts to the upper levels in the hierarchy. Hyndman et al. (2011) propose the OC approach as a new methodology for HF. OC is using all the information available in the hierarchy by forecasting all of the series independently and then uses a regression model to combine and reconcile created forecasts.

When forecasting demand for a hierarchical/grouped SC network, practitioners need to determine: *i)* the univariate forecasting model to use when generating the base forecasts^b *ii)* the appropriate forecasting structure and *iii)* an approach which provides the most accurate forecasts. The latter has attracted the attention of many researchers as well as practitioners over the last few decades (Rostami-Tabar, 2013). Although this has been studied for decades, however, there is no agreement on which HF approach provides more accurate forecasts. Moreover, there is a lack of studies in the literature linking time series characteristics to the accuracy of HF models especially using real

^a GF can be considered as a special case of HF. Depending on the demand structure of the SC, HF or GF methodology might be used.

^b Base forecasts are independent forecasts created at different levels of the forecasting structure by some of the univariate forecasting models.

datasets of a SC structure. To the best of our knowledge, this is the first study that uses a dataset of a multi-echelon SC to investigate not only the effectiveness of HF approaches and their links to the time series characteristics but also the forecast performance of HF combinations. Kahn (1998) was the first to suggest that it is time to combine the existing methodologies so that we can enjoy the good features of both methods, but no specific idea was provided in that discussion.

In this paper, we evaluate the performance of different approaches in a SC context. To do so, we conduct a simulation study and an empirical investigation using real data from a SC distribution network of a major European brewery company. The study aims to: *i)* evaluate the performance of BU, TD and OC approaches; *ii)* examine the accuracy of forecast combination of different approaches; *iii)* propose a model to analyse the effect of time series characteristics on the performance of BU, TD and OC; and *iv)* demonstrate the application of grouped demand forecasting in the SC. For the simulation study, we generate time series at the bottom level using *arima.sim* function in *R* software. The hierarchical structure in the simulation study is a two level hierarchy with eight series in a bottom level and 13 series in total. In the empirical study, the grouped structure of the beer distribution network contains eleven levels with 56 series in the bottom level and 169 series in total. The empirical study aims to produce forecast required by manufacturing, marketing, finances and logistics. For generating the base forecasts, exponential smoothing state space (ETS) models are used in all levels as a univariate forecasting method.

Our contribution to the literature is threefold: *i)* we demonstrate the application of grouped demand forecasting in SC and compare the effectiveness of approaches on a multi-echelon SC from a major European brewery company; *ii)* we examine the forecast performance of HF combination and *iii)* we develop a model to analyse the effect of time series characteristics on the forecast performance of different approaches.

The remainder of the paper is structured as follows: the theoretical background of the HF/GF is introduced in the next section. Section 3 provides the forecasting approaches. Section 4 and Section 5 present the simulation and empirical evaluation, subsequently. Section 6 introduces the effect of time series characteristics on the forecasting performance of HF models. We discuss the findings in Section 7 and conclude the paper with future research and final remarks.

2. HIERARCHICAL AND GROUPED TIME SERIES FORECASTING

Compared with traditional forecasting of univariate time series, forecasting hierarchical or group time series is a more challenging and demanding task for the forecasters. One of the main reasons is because hierarchical or grouped data structures impose additional aggregation constraint, which

needs to be taken into account during the forecasting process. This constraint is related to generating the forecasts which need to be consistent through all levels in the hierarchy or grouped structure^c. That is, an objective is to generate the final forecast that will add up in a way that is consistent with the aggregation structure of the collection of time series (Hyndman and Athanasopoulos, 2018). Therefore, forecasters have to provide forecasts that are accurate and in the same time consistent through all levels of hierarchy or grouped structure. Therefore, the HF/GF could be seen as the principle on how the base forecasts are aggregated, disaggregated, reconciled or combined during the process of generating the final forecasts for each series in the forecasting structure.

There are three main methodologies which can be used when dealing with forecasting hierarchical and grouped time series: BU, TD and OC. The main criterion for selecting among different methodologies is their forecast accuracy. This is essential as an effective planning and operation logistics system require the use of accurate, disaggregated demand forecasts (Caplice and Sheffi, 2006). Errors in forecasting may cause significant misallocation of the resources in inventory, facilities, transportation, sourcing, pricing, and even in information management (Chopra and Meindl, 2007). Forecasting accuracy is directly connected to inventory management, lower errors result in reduced stock keeping without compromising the service level (Trapero et al., 2012). Moreover, inaccurate forecasts will inevitably lead to inefficient, high cost operations and/or poor levels of customer service (Forslund and Jonsson, 2007). Therefore, one of the most important action we may take to improve the efficiency and effectiveness of the logistics process is to improve the quality of the demand forecasts (Caplice and Sheffi, 2006). Starting from the 1950s, there have been extensive discussions in the literature about the merits of TD and BU models. Studies that favour the BU approach are predominantly in the field of an economy (Dunn et al., 1971; Dunn et al., 1976; Kinney, 1971; Edwards and Orcutt, 1969; Collins, 1976). Others argue that TD can produce more accurate aggregate forecasts at top levels (Grunfeld and Griliches, 1960; Aigner and Goldfeld, 1973; Barnea and Lakonishok, 1980). Generally, the proponents of a TD approach argue that the lower-level data is often more error-prone and more volatile (Vogel, 2013) and suggest that the TD approach is superior because of its lower cost and greater accuracy during times of a reasonably stable demand (Weatherford et al., 2001). On the other hand, researchers suggest that BU should be used the distinction between demand patterns for individual items is important (Dunn et al., 1976; Weatherford et al., 2001). Schwarzkopf et al. (1988) argue against using the TD approach for forecasting the bottom level series in a hierarchy. They also challenge the premise that aggregating series reduce the variability in the top level by developing equations which demonstrate that the

^c In literature authors usually refer to this constraint as “aggregate consistency” or “coherent forecasts”.

variability will increase in cases of positive correlation between bottom level series. We found similar conclusions in Gordon et al. (1997); Dangerfield and Morris (1992). While the empirical results tend to point towards the superiority of the BU approach, there is no general consensus on whether a TD or BU approach performs better (Vogel, 2013). In a recent study, Rostami-Tabar et al. (2015) provide the superiority conditions for BU and TD in the one level hierarchy with two nodes in sub-aggregate level, where series follow a non-stationary integrated moving average process of order one. The application of the HF/GF especially in SCs requires the need for accurate forecasts at all levels and not only in the aggregate top level (Fliedner, 1999; Vogel, 2013). In more recent researches, where the emphasis is directed towards the forecast accuracy on all levels in the hierarchy, BU shows better overall performance (Hyndman et al., 2011; Athanasopoulos et al., 2009; Seongmin et al., 2012). There is still a “dead heat race” between the accuracy of TD and BU in the top level of the hierarchy; however, when the entire hierarchy is considered, BU significantly outperforms the TD approach. At the same time, the OC represents a new promising methodology, which has shown excellent results and outperformed others in forecasting tourism, mortality, prison population and labour market data (Hyndman et al., 2011; Shang and Hyndman, 2017; Hyndman et al., 2016; Hyndman and Athanasopoulos, 2018). However, to the best of our knowledge, it has not yet been tested on the SC data. Therefore, there is a need to evaluate their effectiveness on supply chains. Besides choosing among different methodologies, there is also an additional dilemma about choosing the right forecasting model, which comes from the diversity of models in TD and OC methodologies. There are several variations of models in TD and OC methodologies which all have different forecasting performances.

By considering the fact that there is no consensus on which approach provides the most accurate forecasts, we fill the gap in the literature by comparing the forecast accuracy of the different approaches using empirical data of a supply chain and a simulation study. Additionally, we examine the performance of combining different models by creating a unique forecast and compare it against individual approaches. Finally, we develop a model to analyse the impact of time series characteristic on the superiority condition of existing approaches.

3. FORECASTING APPROACHES FOR GROUPED TIME SERIES

Common approaches to forecast group time series often include bottom-up, top-down and the optimal combination approach. Each of them has its own unique principle as well as advantages and disadvantages, which will be further explained in the following subsections. In addition to these approaches, we introduce two combination schemes for combining the forecasts of different approaches in the subsection 3.4.

3.1. Bottom up (BU) methodology

BU approach first generates the base forecasts in the bottom level of the forecasting structure, using a univariate forecasting model. All other forecasts in the structure are generated through aggregating of the base forecast to the higher levels, in a manner which is consistent with the observed data structure. For that purpose summing matrix \mathbf{S} can be used to represent the matrix that dictates how the aggregation of higher level series is calculated from the bottom level series. Therefore, the final forecasts in the BU approach can be expressed as following:

$$\tilde{\mathbf{y}}_h = \mathbf{S} \cdot \hat{\mathbf{y}}_{B,h}. \quad (1)$$

Where $\tilde{\mathbf{y}}_h$ represents the vector of all final forecasts in a given structure for the h - step ahead periods and $\hat{\mathbf{y}}_{B,h}$ represents the vector of all the bottom level forecasts, generated for h - steps ahead.

Since the BU creates the base forecasts at the bottom level, it uses a significant amount of information available in the data. This could result in a better capturing of the individual dynamics of the series in the bottom level. On the other hand, series in the bottom level may be noisy and hard to forecast which may lead to inaccurate forecasts, especially in the top level of the forecasting structure.

3.2. Top down(TD) methodology

TD consists of generating the forecast at the top level of the structure and then disaggregate it to the bottom level in the structure. For disaggregating the top level forecasts, TD methodology uses the disaggregation proportions (p_j). Hence, the forecasting principle of TD can be presented as:

$$\tilde{\mathbf{y}}_h = \mathbf{S} \cdot \hat{\mathbf{y}}_h \cdot \mathbf{p}. \quad (2)$$

Where $\hat{\mathbf{y}}_h$ represents the top level base forecast generated for the h - step ahead periods and $\mathbf{p} = [p_j]$ is a vector containing all disaggregation proportions corresponding to the series in the bottom level. Where $j = 1, \dots, n$; and n is the number of bottom level series in the forecasting structure.

Generally, there is a lot of criticism in the literature regarding the performance of TD methodology in the lower levels of the forecasting structures. The poor performance of the TD approach in the lower levels lies in the disaggregation proportions. There are several variations of the TD approach based on how the disaggregating proportions are determined. These variations could be classified into two groups: approaches that use historical proportions and those that use future forecasts to determine disaggregation proportions.

3.2.1. Top down approaches based on the historical proportions

In the literature, there are three TD approaches based on historical proportions to determine disaggregation weights. Gross and Sohl (1990) examined twenty-one different proportional disaggregation schemes, which include simple averages of the sales proportions, lagged proportions and combined lagged proportions. They suggest two disaggregation proportions as best for disaggregating the top level forecasts: *i)* average historical proportions (TD1) and *ii)* proportions of the historical averages (TD2). The majority of practitioners are still using disaggregation proportions, suggested by Gross and Sohl (1990).

For the TD1, the disaggregation proportions are determined in the following way:

$$p_j = \frac{1}{T} \sum_{t=1}^T \frac{y_{j,t}}{y_t}. \quad (3)$$

Disaggregation proportions of TD1 represent the mean value of the proportions between the series in the bottom level ($y_{j,t}$) and the top level series (y_t), observed in the historical period $t = 1, \dots, T$. Similarly, disaggregation proportions for TD2 reflect the relationship between the average historical values of the same series and they are determined as follows:

$$p_j = \frac{\sum_{t=1}^T \frac{y_{j,t}}{T}}{\sum_{t=1}^T \frac{y_t}{T}}. \quad (4)$$

Based on the TD1 and TD2 models, Chen et al. (2008) attempted to improve the accuracy of the TD methodology. Accordingly, they proposed minimizing the sum of squared demand errors and determining the disaggregating proportions as a result of that process. In their approach, disaggregating proportions are determined as follows:

$$p_j = \frac{\sum_{t=1}^T y_{j,t} y_t}{\sum_{t=1}^T y_t^2}. \quad (5)$$

We will refer to this approach as TD3 in the following Sections.

3.2.2. Top down approaches based on future forecasts

Bearing in the mind that disaggregation proportions can change over time that could significantly deteriorate the forecast accuracy in the bottom level, it is crucial to capture the dynamic nature of disaggregation proportions. Therefore, there is another direction for obtaining disaggregating proportions, which consists of using the future forecasts of the series in the forecasting structure.

Fliedner (2001) was among the first to propose such a TD model and suggested using final forecasts of the BU model for that purpose. The author proposed calculating the ratio of the direct child

forecast divided by the sum of the direct child forecasts comprising their families. The parent forecast is multiplied by this ratio. For more details refer to Appendix in (Fliedner, 2001). We will refer to this approach as TD4 in the following Sections.

Top down forecasted proportions (TDFP) is another TD approach for generating the disaggregating proportions by using future forecasts (Athanasopoulos et al., 2009). For that purpose, the TDFP is using future forecasts of the top and bottom level series, which as a result significantly improved the accuracy of the TD methodology. Boylan (2010) note that although this has not been tested on SC data, the use of forecasted proportions rather than historical proportions appears to be promising. The principle of determining TDFP forecasted proportions is the following:

$$p_j = \prod_{l=0}^{K-1} \frac{\hat{y}_{j,h}^{(l)}}{\hat{s}_{j,h}^{(l+1)}}. \quad (6)$$

Where $\hat{y}_{j,h}^{(l)}$ is h - step ahead base forecast of the node that is l levels above j , and $\hat{s}_{j,h}^{(l)}$ refers to the sum of the h - step ahead base forecasts below the node which is l levels above the node j and directly connected to that node (Hyndman and Athanasopoulos, 2014).

3.3. Optimal combination (OC) methodology

The OC approach uses all the information that is available in the series by generating the univariate forecasts for all of the series in the forecasting structure. Since the independent univariate forecasts do not meet the condition of “aggregate consistency”, OC is performing the reconciliation of the forecasts. The aim of reconciliation is to produce the final forecasts which are mutually coherent and at the same time close to the initial independent base forecasts. The generic formula for producing all final h – step ahead forecasts ($\tilde{\mathbf{y}}_h$) in the OC approach is the following:

$$\tilde{\mathbf{y}}_h = \mathbf{S}(\mathbf{S}'\mathbf{W}_h^{-1}\mathbf{S})^{-1}\mathbf{S}'\mathbf{W}_h^{-1}\hat{\mathbf{y}}_h. \quad (7)$$

Where \mathbf{W}_h represents the variance-covariance matrix of the base forecast errors.

There are four variations of the OC approach, depending on how the estimation of the \mathbf{W}_h matrix is performed. These estimators are: *i)* ordinary least square, *ii)* weighted least squares, *iii)* structural scaling and *iv)* the minimum trace. In this paper, we used the minimum trace estimator since it provided the most accurate forecasts in the simulation and empirical study^d.

For more details regarding the OC approach and its different estimators, see (Hyndman et al., 2007; Hyndman et al., 2016; Hyndman and Athanasopoulos, 2018).

^d Due to the space restrictions, we here only present the results of the OC with minimum trace estimator. Results of other estimators are presented in the following online Shiny platform and more details about their performance could be obtained by request from the authors via email.

3.4. Combination approaches

In this paper, we also used the two combination approaches based on existing approaches: *i)* COMB – the combination of models with no weights which is shown in the Equation 8 and *ii)* COMB_w – the wighted combination of models, shown in the Equation 9 (Ballou, 2004).

$$COMB = \frac{F_{m1} + F_{m2} + \dots + F_{mn}}{n}; \quad (8)$$

$$COMB_w = w_1 \cdot F_{m1} + w_2 \cdot F_{m2} + \dots + w_n \cdot F_{mn}; \quad (9)$$

where: $F_{m1}, F_{m2}, \dots, F_{mn}$ – different forecasting models; $w_i = \frac{\frac{Total\ error}{Forecasting\ error_i}}{\sum_{j \in U} \frac{Forecasting\ error_j}{Total\ error}}$; for $i = 1, \dots, n$

and $Total\ error = \sum_{j \in U} Forecasting\ error_j$, $U \in \{F_{m1}, F_{m2}, \dots, F_{mn}\}$, $Forecasting\ error$ – some of the measurement of forecasting performance.

We use the combination approaches in the simulation and the empirical study. Different comination of BU, TD and its varities and OC approaches are used depending on their individual performances in simulation and empirical study. We also demonstrate the performance of the combination approaches via two forecast accuracy measures.

4. NUMERICAL SIMULATION

In Section 4, we perform a simulation study to evaluate: *i)* the relative performance of the TD, the BU and the OC approaches; and *ii)* the performance of the combination approaches.

4.1. Experiment design

For evaluating the performance of the different approaches, a simulation study is performed. The simulation hierarchy consists of two levels, where the top aggregated series (Total) is subdivided into four series at level 1 (A, B, C and D) and each of series is further disaggregated into two additional series at the level 2 (AA, AB, BA, BB, CA, CC, CB, DA and DD). Therefore, there are eight time series in the bottom and 13 series in total (Fig. 1). The *Autoregressive Integrated Moving Average* (ARIMA) process is used for generating the simulated series at the bottom level of Fig. 1, since the ARIMA framework of the analysis, has been the most useful for research in the SC forecasting (Syntetos et al., 2016; Rostami-Tabar et al., 2015).

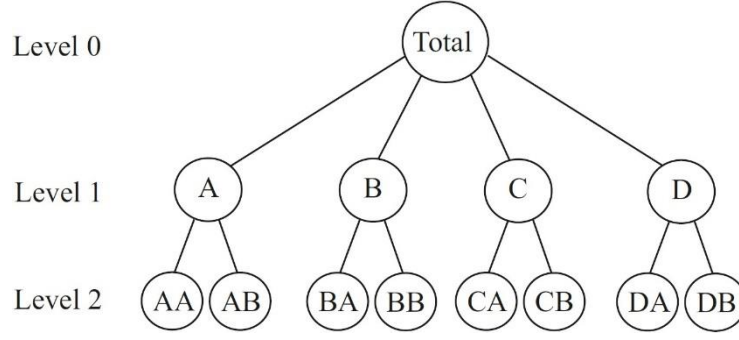


Fig. 1. The hierarchical structure of the simulation study.

During the simulation, orders of the ARIMA process (d - differencing, p - autoregression and q - moving average) were chosen randomly and restricted to values of 0, 1 and 2. Moving average parameter (θ) and autoregressive parameter (ϕ) in ARIMA process are also chosen randomly from the interval $[-0.99, 0.99]$. Therefore, we generate bottom level series ($\mathbf{y}_{B,t}$) corresponding to eight nodes at level 2 of the Fig. 1. We then obtain all other series (\mathbf{y}_t) by aggregating the bottom level series. The process of obtaining all the series in the hierarchy could be represented as:

$$\begin{bmatrix} y_t \\ y_{A,t} \\ y_{B,t} \\ y_{C,t} \\ y_{D,t} \\ y_{AA,t} \\ y_{AB,t} \\ y_{BA,t} \\ y_{BB,t} \\ y_{CA,t} \\ y_{CB,t} \\ y_{DA,t} \\ y_{DB,t} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} y_{AA,t} \\ y_{AB,t} \\ y_{BA,t} \\ y_{BB,t} \\ y_{CA,t} \\ y_{CB,t} \\ y_{DA,t} \\ y_{DB,t} \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}}_S \cdot \underbrace{\begin{bmatrix} y_{AA,t} \\ y_{AB,t} \\ y_{BA,t} \\ y_{BB,t} \\ y_{CA,t} \\ y_{CB,t} \\ y_{DA,t} \\ y_{DB,t} \end{bmatrix}}_{\mathbf{y}_{B,t}}$$

Or in compact form:

$$\mathbf{y}_t = S \cdot \mathbf{y}_{B,t}. \tag{10}$$

Each generated series has 56 observations and all are restricted to be positive. Those series that were generated with negative observations are scaled up until all observations become positive. The simulation is repeated for 500 times, producing 500 different scenarios of the series in the bottom level. In the literature, there were only two similar simulations studies to look upon (Hyndman et al.,

2011; Hyndman et al., 2007). They used 600 and 1000 simulations, respectively. In this paper, we used 500 simulations since after 300-350 simulations, forecasting errors become stable.

For evaluating the forecasting performance of each model, we divide each simulated series into in-sample/training and out-of-sample/test sets. Training data are initially set with 12 observations, and the test set with 8 observations. We use ETS forecasting model to produce out of sample base forecasts. Forecasting horizon is set to 8-step-ahead and 1 to 8 steps ahead forecast is produced. After that, the out of sample error is determined for every time series from Fig. 1. We use the rolling forecasting procedure for the evaluation. The process is consisted from iteratively adding one observation to the training set and generating 1 to 8 steps ahead forecasts. The procedure is constantly repeated until training data reached 48 observations. This process yields 32 different error sets for each node in one simulation scenario. After that, the process is repeated with another simulation scenario. Finally, we use the Root Mean Square Error (RMSE) to summarise and report the accuracy by finding the average value of RMSE across all different error sets and the simulation scenarios. Additionally, we also provide a Mean Absolute Percentage Error (MAPE) in Table A.1.

4.2. Numerical results

In subsection 4.2, we present the result of simulation investigation on the performance of the HF approaches and HF combinations.

4.2.1. Forecasting performance of the hierarchical forecasting approaches

Table 1 presents the performance of HF models for the hierarchical structure illustrated in Fig. 1.

Table 1. RMSE of different models based on the 500 simulation scenarios^e.

Level	Node	RMSE						
		BU	TD1	TD2	TD3	TD4	OC	TDFP
Level 0	Total	55.54	62.12	62.12	62.12	62.12	54.83	62.12
Level 1	A	20.62	180.79	152.71	354.97	23.10	20.44	24.10
	B	24.44	101.22	82.20	217.69	26.24	23.84	26.07
	C	13.37	121.97	97.52	268.99	17.45	13.58	17.78
	D	25.64	179.45	152.62	328.92	29.91	25.47	30.00
Level 2	AA	11.27	115.55	102.27	188.94	13.09	11.28	21.45
	AB	12.78	78.51	63.88	183.73	14.56	12.73	21.73
	BA	13.80	47.96	37.63	104.33	14.99	13.60	23.11
	BB	13.65	69.23	59.37	133.55	14.95	13.47	23.46

^e Best results are bolded.

CA	7.09	35.06	29.56	74.30	8.58	7.47	10.38
CB	8.98	100.63	76.12	225.78	11.65	9.15	13.47
DA	11.20	97.37	81.35	181.23	13.09	11.23	13.96
DB	17.10	92.41	78.44	165.58	19.54	17.16	20.72
Average	18.11	98.64	82.75	191.55	20.71	18.01	23.72

Results show that overall, the OC demonstrates the most accurate forecasts in the hierarchy. It performs better than the BU, although the following tests failed to identify any statistically significant difference between these models. OC and BU are closely followed by TD4. TDFP performed well in the upper levels of the hierarchy but failed to produce such accurate forecasts in the bottom level. Results also show the underperformance of TD1, TD2 and TD3 approaches comparing to others in all levels of the hierarchy. The reason could be in the fact that simulated series have varying levels of correlation and participation of the bottom level series in the top aggregate series, which we found to be important in the question of accuracy of forecasting approaches. Likewise, a great number of the bottom level series has a vigorous and dynamic trend, that TD1, TD2 and TD3 are not able to appropriately capture and incorporate in the future forecasts. Therefore, forecasts of these models prove to be unreliable and inaccurate. Moreover, these approaches might be drastically underperformed in some of the simulation scenarios (more details in Section 6) which consequently deteriorate their performance. In Contrast, OC, BU, TD4 and TDFP approaches produce more robust and stable forecasts.

This is demonstrated in Fig. 2 where BU, TD4, OC and TDFP have a narrow interquartile range, in the box plots which represent their forecasting performance. This suggests that observed models have consistent forecasting errors across all levels and series in the hierarchy. Conversely, TD1, TD2 and TD3 have much larger interquartile range indicating higher dispersion of the forecasting errors, through different levels and series in the hierarchy. Therefore, the outperformance of the TD1, TD2 and TD3 is clear. Differences presented in Fig. 2 and Table 1 are tested on statistical significance using the Kruskal–Wallis test. The p -value for the χ^2 in the Kruskal–Wallis test is below the 0.05. This suggests that there is a significant difference in forecasts between models. Dunn’s post-hoc test is used to identify the pairs of significantly different forecasts. The test revealed that TD1, TD2 and TD3 generated statistically indistinguishable forecasts. At the same time, the test identified a statistical difference among forecasts of TD1, TD2 and TD3 and others. Supplementary, the test failed to identify any statistically significant difference among the forecasts of the BU, OC, TD4 and TDFP.

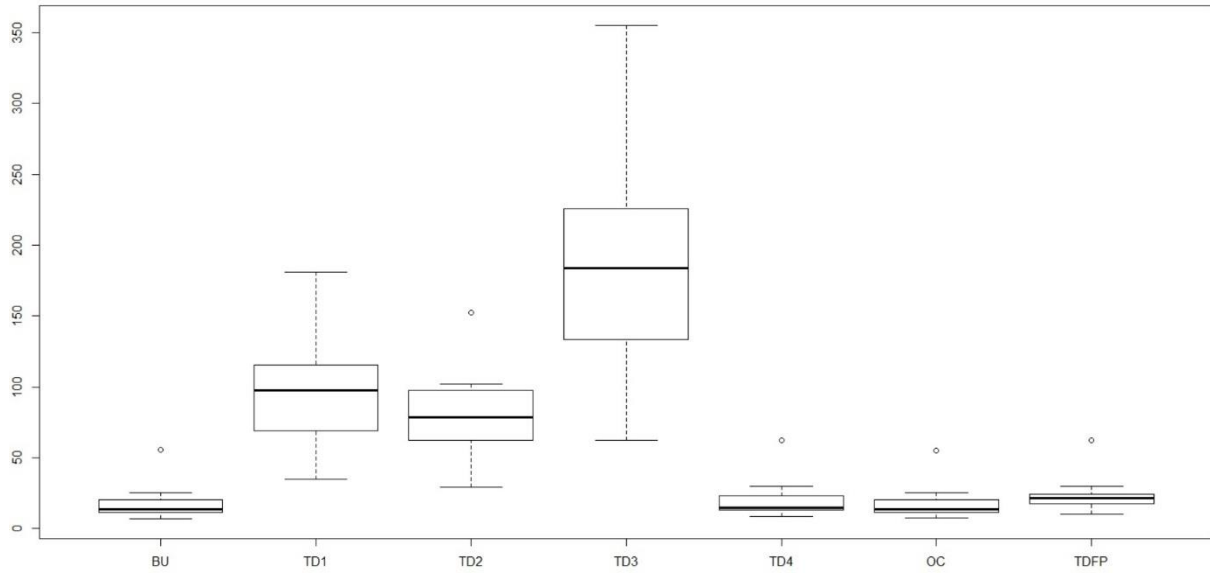


Fig. 2. Box plots for RMSEs of different models tested on simulated data.

Fig. 3 displays the performance of different HF models through the hierarchy levels, and demonstrates diverging performances by moving from top to the bottom level series. The figure presents that TD1, TD2 and TD3 underperformed compared to other competing models. Moreover, it is noticeable that TD1, TD2 and TD3 perform better only at the highest level of the hierarchy and that their performance deteriorates in all other levels. However, other models show a consistent performance regardless of the hierarchical level. We develop a shiny application that allows readers to perform the comparison between models and their performances. It is available in the following link: https://forecastingsupplychain.shinyapps.io/simulation_study/.

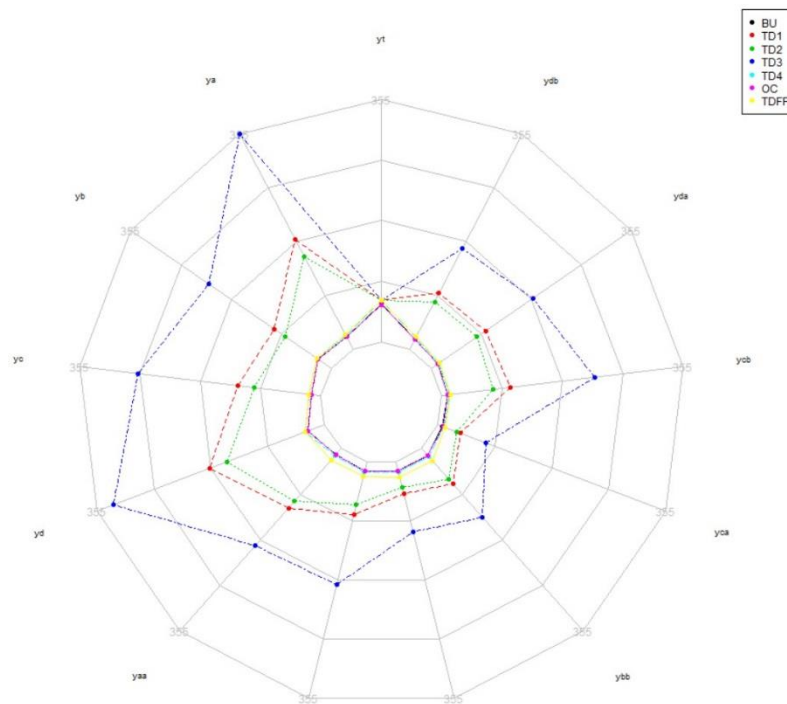


Fig. 3. Performance of different HF models through the hierarchy levels.

We also present the performance of different HF models using MAPE in Table A.1. Results show that BU and OC are the best performing models and that differences in their performance are not statistically significant. Therefore, the insight drawn from Table A.1 is aligned with Table 1.

The overall performance of BU model was surprisingly good since it managed to be competitive with the OC model. Similar results could be found in the (Athanasopoulos et al., 2009; Hyndman et al., 2011). It is possible that BU could demonstrate much poorer results in the hierarchies which have more levels and a higher degree of disaggregation of the data. Forecasting highly disaggregated data (possibly intermittent) could be a challenging task for BU since some of the series features could be undetected and therefore not incorporated in the future forecasts. As a consequence, generating all other forecasts in the hierarchy by aggregating these granular forecasts, could possibly produce misleading and not accurate overall forecasts.

4.2.2. Performance of the combination approach

Bearing in the mind that several HF models may produce different outcomes (Table 1), combining approaches may be the potential solution to improve the forecast accuracy. This recommendation is usually given for individual (univariate) forecasting. We examine the combination of HF approaches to see whether combining forecasts of HF models produces more accurate forecasts than using a single HF model.

We examined various combinations of HF models in this research. However, we only present the most accurate combination approach. The HF approaches present in the forecast combination approaches are OC, BU, TD4 and TDFP approaches. These models were top-performing models in Table 1. The combination approaches are described in subsection 3.4. RMSEs of the HF combination models used in the simulation study are presented in the last two columns of Table 2.

Table 2. RMSEs of combined forecasts based on the 500 simulation scenarios^f.

Level	Node	RMSE								
		BU	TD1	TD2	TD3	TD4	OC	TDFP	COMB	COMBw
Level 0	Total	55.54	62.12	62.12	62.12	62.12	54.83	62.12	55.91	51.12
Level 1	A	20.62	180.79	152.71	354.97	23.10	20.44	24.10	20.22	19.99
	B	24.44	101.22	82.20	217.69	26.24	23.84	26.07	24.46	24.11
	C	13.37	121.97	97.52	268.99	17.45	13.58	17.78	14.35	13.71

^f Best results are bolded.

	D	25.64	179.45	152.62	328.92	29.91	25.47	30.00	26.54	25.60
Level 2	AA	11.27	115.55	102.27	188.94	13.09	11.28	21.45	11.34	11.39
	AB	12.78	78.51	63.88	183.73	14.56	12.73	21.73	12.45	12.37
	BA	13.80	47.96	37.63	104.33	14.99	13.60	23.11	13.97	13.85
	BB	13.65	69.23	59.37	133.55	14.95	13.47	23.46	13.65	13.54
	CA	7.09	35.06	29.56	74.30	8.58	7.47	10.38	7.54	7.30
	CB	8.98	100.63	76.12	225.78	11.65	9.15	13.47	9.74	9.33
	DA	11.20	97.37	81.35	181.23	13.09	11.23	13.96	11.72	11.30
	DB	17.10	92.41	78.44	165.58	19.54	17.16	20.72	17.75	17.25
	Average	18.11	98.64	82.75	191.55	20.71	18.01	23.72	18.43	17.75

Results show that the COMBw approach provides the most accurate forecasts on average. It is calculated based on the simple average of RMSE across all levels. Moreover, the COMBw approach significantly outperforms the all other HF approaches in the top level and generated consistently good forecasts in the remaining nodes of the hierarchy. COMBw was closely followed by OC, BU and COMB models. The accuracy of TD1, TD2, TD3, TD4, and TDFP was far behind the accuracy of the COMB and COMBw models (Fig. 4).

Fig. 4 displays the performance of COMB and COMBw models (grey boxplots) compared to other HF models. It indicates that combining the forecasts of OC, BU, TD4 and TDFP reduces forecast errors compared with any individual HF model. Moreover, it indicates that combinations have a thin interquartile range, implying that COMB and COMBw generated reliable forecasts, which resulted in stable forecasting errors through the entire hierarchy. Therefore, COMB and COMBw are robust and consistent for generating forecasts across various hierarchies.

OC model generates the most accurate forecasts in the simulation study (Table 1). Therefore, we use it as a benchmark to compare the forecast accuracy improvement of COMB and COMBw (and all other HF models). This has been shown in Table 3.

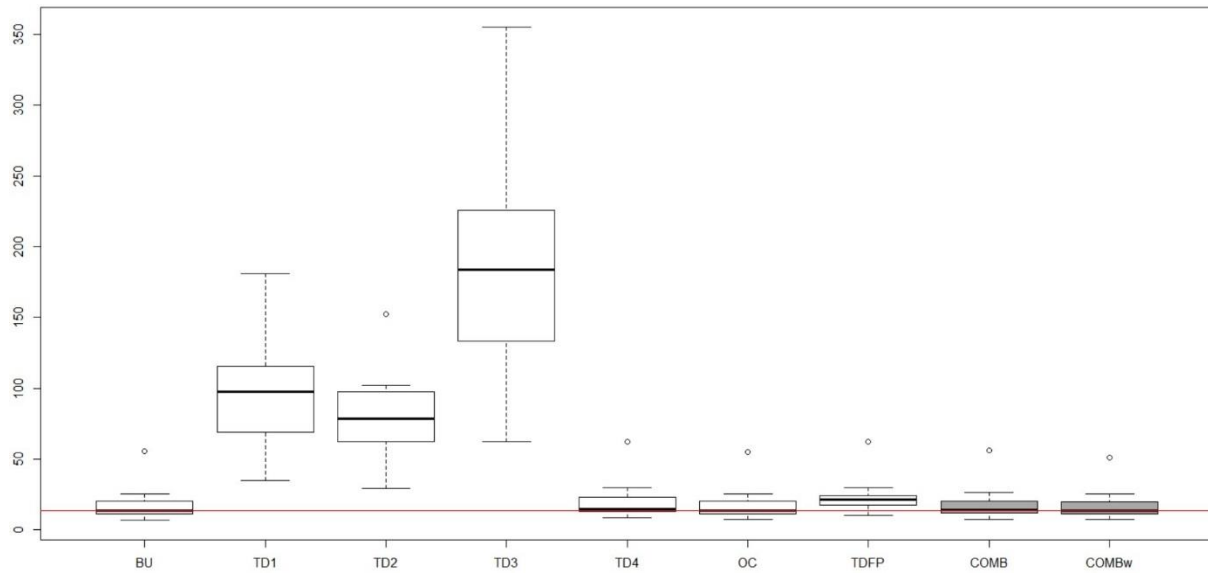


Fig. 4. RMSE forecasting error of COMB and COMBw models in the simulation study^g.

Table 3. Percentage increase or decrease of RMSE forecasting error of all HF models compared to the OC model^h.

Level	Node	%								
		BU	TD1	TD2	TD3	TD4	OC	TDFP	COMB	COMBw
Level 0	Total	1.29	13.29	13.29	13.29	13.29	0	13.29	1.97	-6.76
Level 1	A	0.89	784.62	647.21	1636.91	13.02	0	15.75	-1.07	-2.19
	B	2.56	324.66	244.87	813.29	10.09	0	17.93	2.61	1.15
	C	-1.51	798.42	618.32	1881.4	28.53	0	44.21	5.72	0.96
	D	0.67	604.65	499.27	1191.55	17.43	0	23.25	4.22	0.54
Level 2	AA	-0.05	924.63	806.9	1575.4	16.12	0	34.18	0.58	0.98
	AB	0.45	516.87	401.91	1343.68	14.38	0	7.73	-2.16	-2.82
	BA	1.49	252.65	176.65	667.08	10.21	0	19.52	2.73	1.82
	BB	1.29	413.89	340.7	891.31	10.96	0	20.58	1.34	0.48
	CA	-5.09	369.1	295.44	893.96	14.75	0	31.71	0.81	-2.31
	CB	-1.85	1000.06	732.19	2368.24	27.32	0	48.34	6.45	1.99
	DA	-0.26	766.97	624.36	1513.66	16.56	0	32.62	4.37	0.57
	DB	-0.32	438.56	357.1	864.97	13.88	0	20.3	3.41	0.54
Average		-0.03	554.49	442.93	1204.21	15.88	0	40.42	2.38	-0.38

^g The red line in a figure represents the median value of the RMSE forecasting error of the COMBw model.

^h Best results are bolded.

Table 3 represents the percentage increase or decrease of forecasting error of different HF models, compared to the forecasts of the OC model in every node. Negative values mean forecast improvement comparing to the OC approach. Results demonstrate that by average COMBw decreased forecasting error in the amount of -0.38%, compared to the forecasts of the OC model. All other models generate forecasts that are less accurate than forecasts of OC in each node of the hierarchy, except BU. Which observed by this principle of direct percentage comparison has more accurate forecasts by 0.03%.

Results of Table 2, Table 3 and Fig. 4 indicate that combining the forecasts of OC, BU, TD4 and TDFP provides more accurate forecasts than any of individual HF models. Therefore, results of a simulation study show that HF combinations improved forecast accuracy.

5. EMPIRICAL EVALUATION

In Section 5, we assess the empirical validity of the main findings of this research using real time series of a SC distribution network from a European brewery company. There is a lack of studies evaluating the performance of the BU, TD and OC in the SCs. There are only a few examples linking forecasting to various parts of SCs (Mircetic et al., 2017; Seongmin et al., 2012; Rostami-Tabar et al., 2015; Mircetic, 2018; Villegas and Pedregal, 2018; Pennings and van Dalen, 2017). To the best of our knowledge, this is the first study that examines a comprehensive grouped demand forecasting in a SC network. The empirical study is performed to evaluate the effectiveness of different approaches in a real SC network. Additionally, we also examine GF combinations in SC which has never been investigated.

In subsection 5.1, we first provide details of the real SC distribution network and the empirical data available for the purposes of our investigation along with the experimental structure employed in our work. We then present the actual empirical results in subsection 5.2.

5.1. Supply chain distribution network

Fig. 5 illustrates a distribution structure of brewery company operating in the South-East of Europe. The scale economies in the transport of freight, combined with the market requirement to provide fast and reliable delivery times, drive most large firms to operate multi-echelon distribution inventories (Caplice and Sheffi, 2006). For the same reasons, the observed brewery company has a multi-echelon distribution structure, and its distribution network spreads over several distribution centers (DC) located across various geographical regions. Different DCs are designated to serve only particular market regions. The distribution starts from the central warehouse, which is directly connected to the manufacturing plant. The plant produces more than 200 different beer product

families. The annual output from the central warehouse varies, and it is usually between 250 000-300 000 pallets. Highest demand peaks occur in the spring/summer months (May, June and July) with the demand picking to the 14 000 pallets of different brewery products.

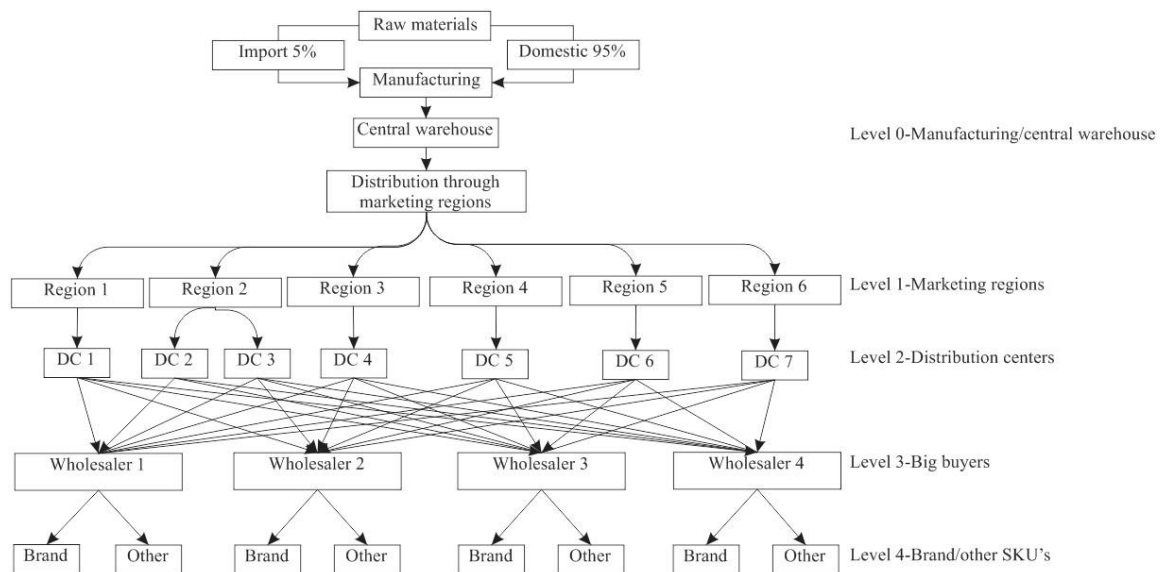


Fig. 5. Multi-echelon brewery distribution network.

The brewery industry is specific in the sense that all of the consumption of the products is accomplished via bars, restaurants and retail shops. There are no direct deliveries and internet sales of products. In order to provide products to a wide consumer network, observed distribution chain is divided into the eight marketing regions. These marketing regions are supplied through seven DCs. Each region has one designated DC, except region 2, which is served through two DCs. Further distribution of goods is carried out through wholesalers. There are four big wholesalers which are dominating in the observed market. Some of those wholesalers are big retail chains, while others act as agents between manufacturers and small retailers, bars and restaurants. Nevertheless, each wholesaler is acquiring brewery products from the DCs and makes further placement of goods on the market. Goods that are provided to the wholesalers are classified as a brand and other products. Brand products represent the most important products for the company since there are providing the majority of revenue in the market. It is a top selling beer which comes in different packaging types. In observed brewery distribution chain, there is no further feedback from the wholesalers regarding the point of sale data, therefore the visibility of the customer data is limited.

Main reasons for using the GF in SCs is to simplify forecasting process, obtain more accurate forecasts, harmonise forecasts from different levels and to provide all information needed for different SC parties. Therefore, in this empirical study, the forecasting structure is designed to generate forecasts that support the planning and execution of the processes in different parts of SC.

Special attention is addressed to the alignment of the time component as well, and not just on the cross-sectional alignment of the grouped structure.

5.2. Grouped structure for forecasting the brewery demand

The demand dataset available for the purpose of our research includes the period from 2012 to 2015; and consists of 56 weekly demand observations from a brewery company. The unit of observation is a pallet. The demand in a given multi-echelon brewery distribution chain has the grouped structure. The structure is provided in Table 3, where each row denotes the level of disaggregation.

Table 3. Grouped structure of brewery demand.

Disaggregation level	Level	Labels	Number of series
0	Total	Total	1
1	Regions (marketing regions)	$R_1, R_2 \dots R_6$	6
2	Distribution centers	$DC_1, DC_2 \dots DC_7$	7
3	Wholesalers	$W_1, W_2 \dots W_4$	4
4	Product types	B and O	2
5	Regions x Distribution centers	$R_1DC_1, R_2DC_2, \dots R_6DC_6$	7
6	Regions x Wholesalers	$R_1W_1, R_1W_2 \dots R_6W_4$	24
7	Regions x Product types	$R_1B, R_1O \dots R_6O$	12
8	Distribution centers x Wholesalers	$DC_1W_1, DC_1W_2 \dots DC_7W_4$	28
9	Distribution centers x Product types	$DC_1B, DC_1O \dots DC_7O$	14
10	Wholesalers x Product types	$W_1B, W_1O \dots W_4O$	8
11	Distribution centers x Wholesalers x Product types	$R_1DC_1W_1B, R_1DC_1W_1O \dots R_6DC_7W_4O$	56
Total number of series			169

At the top level, the total aggregate demand for brewery products is presented. Demand is further divided by marketing regions, DCs, wholesalers, product types and their accompanying interactions. This division provides information related to the manufacturing with total demand, marketing by region demand and product types of demanded products (brand or other products), a financial sector with the large buyers demand and logistics with a spatial fragmentation of demand. The total node represents the central warehouse, while nodes from R_1 to R_6 represent the marketing regions. Nodes at level 2 (from $DC_1, DC_2 \dots DC_7$) represent the DCs. In level 3, four wholesalers are represented with nodes from W_1 to W_4 . In level 4, B and O nodes represent the product types. Further levels represent the interactions between observed disaggregating features. Levels 5, 6, 7 represent the demand disaggregation of different marketing regions by DCs, wholesalers and product types (nodes from R_1DC_1 to R_6O). In levels 8 and 9 demand of DC is further subdivided by the wholesalers and product types (nodes from DC_1W_1 to DC_7O). Nodes in level 10 (from W_1B to W_4O), represent the demand of each wholesaler subdivided by product types. The most disaggregated data arise when

we consider the two product types that are supplied through seven different DCs to four different wholesalers, giving a total of $2 \times 7 \times 4 = 56$ bottom level series in the observed grouped structure. These series represented by the nodes from $R_1DC_1W_1B$ to $R_6DC_7W_4O$. Time plots of time series for the first four levels are presented in Fig. 6. Results show that series are non-stationary, with week trend and pronounced seasonality.

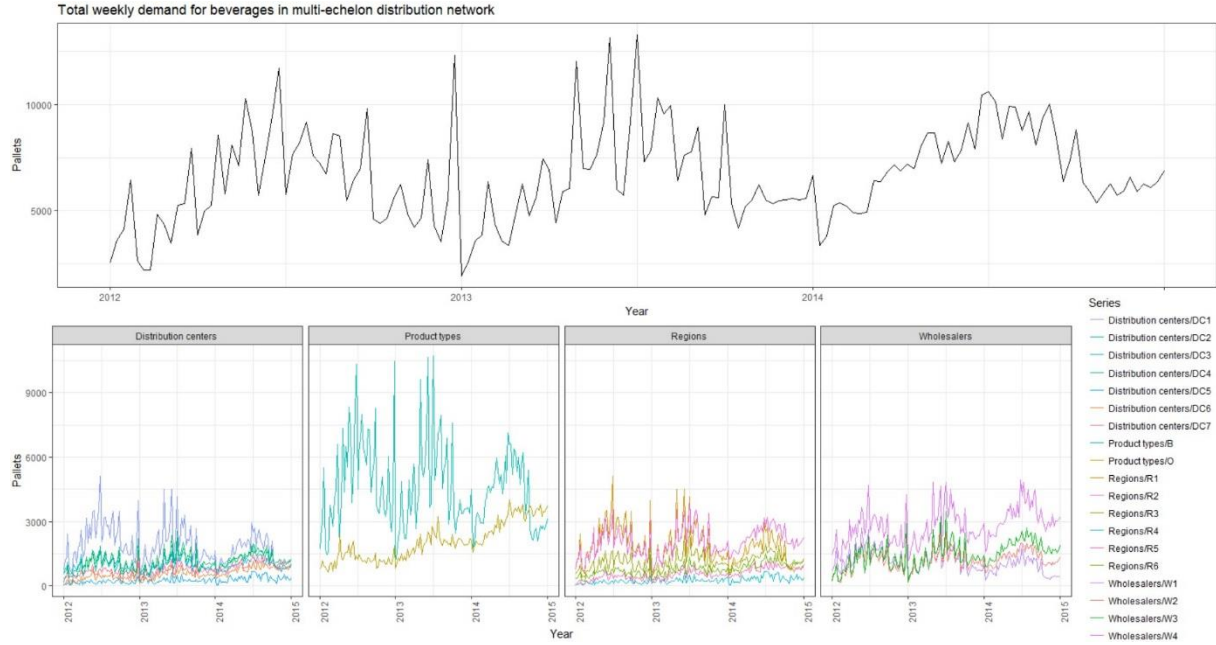


Fig. 6. Total weekly demand, disaggregated by marketing regions, DCs, wholesalers and product types.

5.3. Data and empirical design of experiment

Forecast horizon is equal to the lead time of the decisions driven by the forecast. Since the replenishment orders are required every week and manufacturing needs annual forecasts for creating the production and procurement plan, the demand is forecasted on the weekly level for one year ahead. All other sectors require forecasts between these two periods, so they can be easily determined by looking at the forecasts for the period of their interest (monthly, quarterly, semi-annually and annually).

For evaluating the forecasting performance of HF/GF approaches using real time series, we divide each series at each level into training/in-sample and test/out-of-sample sets. Training data is set to 104 weekly observations and includes the period from 2012 to 2014. The test data is set to 52 weekly observations and it represents the period from 2014 to 2015. As in the simulation study in the empirical study, we also use ETS forecasting model to produce out of sample base forecasts. Forecasting horizon is set to 52-steps-ahead (one year ahead), and 1 to 52 steps ahead forecasts are

generated: After that, the out of sample error is determined for every series in the grouped structure from Table 3. To report the forecast accuracy we use RMSE and MAPE errors.

For forecasting the grouped brewery demand structure, we only keep BU and OC approaches. We exclude TD1, TD2 and TD3 as they drastically underperformed in the simulation study. Moreover, TD4 and TDFP are not used in the empirical study. TD4 has high implementation complexity in forecasting structures with more than two levels. TDFP does not produce aggregately consistent forecasts when forecasting grouped time series. Shang and Hyndman (2017) also suggest that BU and OC are the only approaches that are currently suitable for forecasting the grouped demand structures.

5.4. Empirical results

In this section, we present the results of the empirical investigation. In the subsection 5.4.1, we evaluate the effectiveness of GF approaches and then we examine whether GF forecast combinations improve the forecast accuracy in real data from a brewery SC.

5.4.1. The effectiveness of the grouped forecasting approaches

Results of the empirical study confirm the simulation results. Fig. 7 shows that BU and OC models produce similar forecasts. We observe that the OC model performs slightly better than others followed by BU (please refer to Table A.2 in Appendix for details). However, its performance is not significantly distinguishable, as the Kruskal–Wallis test failed to identify important discrepancy among forecasts of BU and OC.

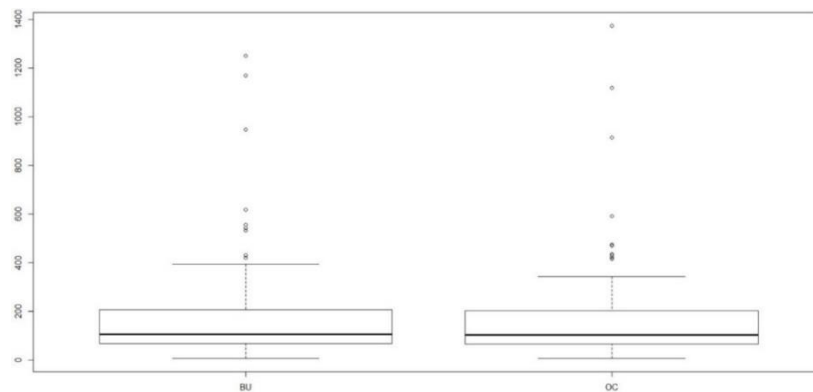


Fig. 7. Box plots for RMSEs of BU and OC models tested on the multi-echelon brewery distribution chain.

Therefore, OC and BU have similar forecasts through all nodes in the forecasting structure as shown in the Figs. 7 and 8. We develop a shiny platform that allows users to compare and visualise the performance of these approaches using real data from the multi-echelon brewery distribution chain, which is available here: https://forecastingsupplychain.shinyapps.io/empirical_case_study/.

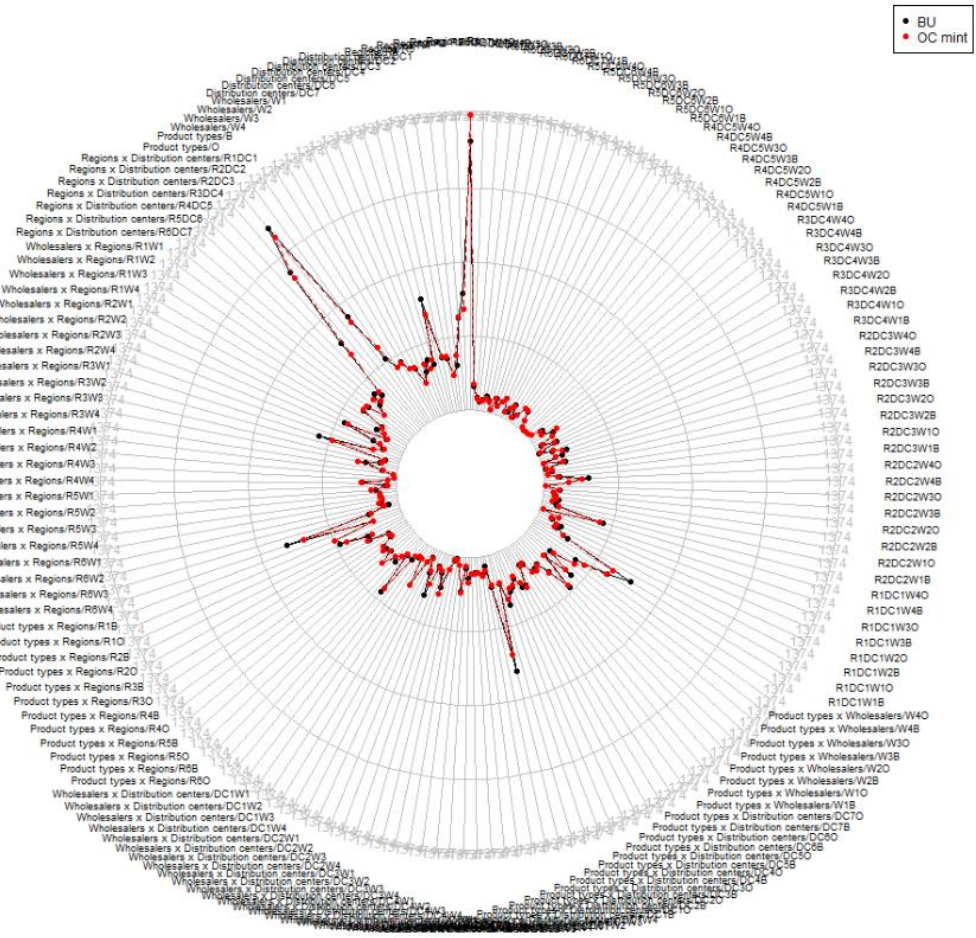


Fig. 8. Performance of OC and BU through the grouped levels while forecasting the demand in the brewery distribution chain.

Fig. 8 shows the performance of OC and BU approaches across all levels in the forecasting structure presented in Table 3. It reveals that both approaches generate similar forecasts in all levels of the forecasting structure.

5.4.2. Performance of the combination approach

In the subsection 5.4.2, we compare the forecast performance of the combination models against the existing approaches using the real dataset from a brewery distribution network. In Fig. 9, the performance of COMB and COMBw is represented next to the performance of OC and BU models.

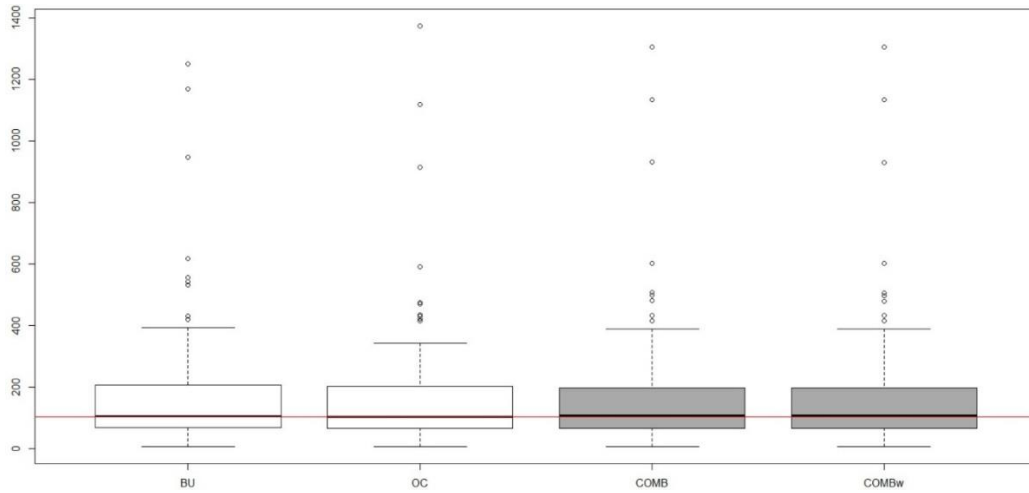


Fig. 9. RMSE forecasting error of COMB and COMBw models in the empirical studyⁱ.

Fig. 9 shows that combining the forecasts of OC and BU through two different combination approaches (Eq. 8 and Eq. 9), produce more accurate forecasts than the forecasts of BU (please refer to Tables A.2 and A.3 for details). COMB was the even the most accurate model according to the MAPE forecasting error (Table A.3).

Since OC model produces the most accurate forecasts in the empirical study (please see Table A.2), we also compare the forecasts of BU, COMB and COMBw with those generated by the OC model (please refer Table A.4 for details). Table A.4 in the Appendix presents the percentage increase or decrease of forecasting errors of different models, compared to the OC forecasts in the whole grouped structure. Results reveal that COMBw and COMB had the closest forecasts to the OC model, generating almost identical forecasts. On average, forecasting errors of COMBw and COMB are only 0.33% and 0.34% higher than the forecasting errors of the OC model. BU forecasts were prominently more inaccurate than forecasts of OC, and on average had forecasting errors of 2.81% higher than the forecasting errors of the OC model. Therefore, forecasts generated by the GF combinations demonstrated very good results and we recommend using GF combinations in contrast to using individual GF approaches to generate grouped time series forecasts in SCs.

6. THE EFFECT OF TIME SERIES CHARACTERISTICS ON THE PERFORMANCE OF FORECASTING APPROACHES

In Section 4, we discuss that the performances of different models are increasingly diverging by moving from the top aggregate level to the bottom level in the hierarchy. At the top level, all HF models showed similar forecasts. In both studies (simulation and empirical), differences between their performances become more apparent in the lower levels of the hierarchy, especially in the

ⁱ The red line in a figure represents the median value of the RMSE forecasting error of the COMBw model.

bottom level. Therefore, we investigate whether there is any connection between the performance of different HF models and characteristics of the series in the bottom level.

To the best of our knowledge, this is the first study that develops a model to evaluate the effect of time series characteristics on the forecasting performance of different HF models. To that end, we develop an additive multiple linear regression model. Multiple linear regression represents a model for forecasting cross-sectional data. It assumes that there is a linear relationship between input features $X=(X_1, X_2, \dots, X_p)$, and the observed variable Y (Eq. 11).

$$Y = f(X) + \varepsilon. \quad (11)$$

where ε is a random error term, which is independent from X and has mean zero. It is the irreducible part of forecasting error of the model, therefore we are only interested in estimating the relationship $f(X)$ shown in the Eq. 12.

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon. \quad (12)$$

where X_j represents the j predictor, and β_j represents the average effect of X_j predictor while holding all other predictors fixed.

The algorithm provides insights into the interaction among characteristics of time series and the accuracy of different HF models. The idea is to measure and extract different characteristics of time series (that comprise the bottom level of the hierarchy), which will then be used as independent features (X) in the multiple linear regressions. RMSEs of different HF models are set as dependent variables (Y). The main goal is to identify the influential (i.e. statistically significant) time series features, rather than creating the most accurate statistical learning algorithm on a given set of data. For every HF model, a separate regression model was created. Therefore, we create seven regression models.

For evaluating the statistical learning algorithm, we form a database which contains the results of the simulation study. For each node in the bottom levels of the simulation study, 19 different time series characteristics and RMSE forecast errors of HF models, are extracted, scaled and recorded in the database. Therefore, the database contains 4000 different entries for time series measures and the HF forecasting errors. We use 70% of the data for training and 30% for the test. Therefore, 19 different time series characteristics are used as independent variables in regression models which are presented in the first column of Table 4 from number 1 to 19. First, 16-time series characteristics are described in detail by Hyndman et al. (2015) and Wang et al. (2009). These characteristics might provide insights into why some HF models perform better than others on the same data. We use given characteristics and added additional three characteristics: *i*) correlation between observed

bottom level series and the top aggregate series (*correlation (bts-top)*); ii) the participation of observed series from the bottom level to the top aggregate series (*aggregate share*); and iii) *coefficient of variation*. Additional characteristics are calculated by the following equations:

$$\text{correlation}(bts - top) = r_j = \frac{\sum_{t=1}^T (y_{j,t} - \bar{y}_{j,t})(y_t - \bar{y}_t)}{\sqrt{\sum_{t=1}^T (y_{j,t} - \bar{y}_{j,t})^2} \sqrt{\sum_{t=1}^T (y_t - \bar{y}_t)^2}}; \quad (13)$$

$$\text{aggregate share} = AS_j = \frac{\sum_{t=1}^T y_{j,t}}{\sum_{t=1}^T y_t}; \quad (14)$$

$$\text{Coefficient of variation} = \frac{\delta}{\mu} = \frac{\sqrt{\frac{1}{T-1} \sum_{t=1}^T (y_{j,t} - \frac{1}{T} \sum_{t=1}^T y_{j,t})^2}}{\frac{1}{T} \sum_{t=1}^T y_{j,t}}. \quad (15)$$

Where $\bar{y}_{j,t}$ represents the mean value of observed j bottom level series ($y_{j,t}$) and \bar{y}_t is the mean value for the top level series (y_t), observed in the historical period $t = 1, \dots, T$ and $j = 1, \dots, n$.

The *Coefficient of variation* measures the volatility of the time series, while *correlation (bts-top)* is measuring the strength and direction of the linear relationship amongst bottom level series and the top aggregate series. *Aggregate share* measures the participation of the observed series from the bottom level in top aggregate series. It provides the information about how “big” or “small” is observed series in the given hierarchy.

Table 4. Summary statistics for the observed time series characteristics.

	Time series characteristics	Mean	Standard deviation	Median	Min	Max	Range	Skew	Kurtosis
1	<i>Lumpiness</i> (Variance of annual variances of remainder)	0.32	0.55	0.01	0.00	5.06	5.06	2.17	5.22
2	<i>Entropy</i> (Spectral entropy)	0.72	0.17	0.68	0.53	1.00	0.47	0.30	-1.53
3	<i>ACF1</i> (First order of autocorrelation)	0.57	0.51	0.88	-0.97	0.99	1.96	-1.14	-0.03
4	<i>Lshift</i> (Level shift)	1.02	0.42	1.00	0.28	2.61	2.33	0.27	-0.51
5	<i>Vchange</i> (Variance change)	0.66	0.69	0.43	0.00	3.65	3.64	1.20	1.03
6	<i>Cpoints</i> (The number of crossing points)	11.60	12.19	6.00	1.00	55.00	54.00	1.00	-0.14
7	<i>Fspots</i> (Flat spots)	7.19	4.84	5.00	1.00	37.00	36.00	1.01	0.58
8	<i>Trend</i> (Strength of trend)	0.71	0.37	0.94	0.00	1.00	1.00	-0.89	-0.86
9	<i>Linearity</i> (Strength of linearity)	1.48	4.30	1.16	-7.49	7.45	14.94	-0.45	-0.57
10	<i>Curvature</i> (Strength of curvature)	0.73	2.20	0.35	-6.17	6.97	13.14	0.64	0.45
11	<i>Spikiness</i> (Strength of spikiness)	0.00	0.00	0.00	0.00	0.01	0.01	2.45	6.52
12	<i>Season</i> (Strength of seasonality)	0.05	0.08	0.03	0.00	0.90	0.90	4.57	30.76
13	<i>Peak</i> (Strength of peaks)	0.35	0.55	0.12	0.00	4.69	4.69	2.56	8.07
14	<i>Trough</i> (Strength of trough)	-0.36	0.49	-0.17	-3.46	0.00	3.45	-2.02	4.84
15	<i>KLscore</i> (Kullback-Leibler score)	2.49	2.46	1.94	0.03	17.30	17.27	1.27	1.74

16	<i>Change.idx</i> (Index of the maximum KL score)	28.28	15.55	32.00	4.00	52.00	48.00	-0.17	-1.44
17	<i>Correlation (bts-top)</i>	0.21	0.52	0.13	-0.99	1.00	1.99	-0.16	-0.67
18	<i>Aggregate share</i>	0.13	0.19	0.04	0.00	0.99	0.99	2.12	4.05
19	<i>Coefficient of variation</i>	3.70	18.12	0.30	0.00	599.86	599.86	17.89	456.80

Table 4 provides the summary statistics for the time series characteristics that are used as input in the statistical learning algorithms. Fig. A.1 in Appendix provides the distributions for the time series characteristics. Distributions have different shapes, but the majority of the time series characteristics have right-skewed distributions.

Table 5. The effect of time series characteristics on the performance of HF approaches.

	Time series characteristics	BU	TD1	TD2	TD3	TD4	OC	TDFP
1	<i>Intercept^j</i>	11.96	80.37	66.09	160.31	13.80	12.00	19.07
2	<i>Lumpiness</i> (Variance of annual variances of remainder)							
3	<i>Entropy</i> (Spectral entropy)	-2.07	-26.64	-20.51	-56.19	-1.86	-1.96	
4	<i>ACF1</i> (First order of autocorrelation)					2.40	2.14	
5	<i>Lshift</i> (Level shift)	0.75	-23.79	-21.33	-44.25	0.74	0.74	
6	<i>Vchange</i> (Variance change)							
7	<i>Cpoints</i> (The number of crossing points)		-22.02	-19.75		1.73	1.50	
8	<i>Fspots</i> (Flat spots)	2.98	25.64	28.33	42.69	3.77	3.04	4.35
9	<i>Trend</i> (Strength of trend)							
10	<i>Linearity</i> (Strength of linearity)	-2.41	-16.93	-12.77	-29.95	-2.11	-2.39	-6.66
11	<i>Curvature</i> (Strength of curvature)	1.78			20.37	1.81	1.74	
12	<i>Spikiness</i> (Strength of spikiness)		-10.64	-9.56				
13	<i>Season</i> (Strength of seasonality)		-6.93	-6.99		-0.54	-0.53	
14	<i>Peak</i> (Strength of peaks)							
15	<i>Trough</i> (Strength of trough)			-5.49				
16	<i>KLscore</i> (Kullback-Leibler score)							
17	<i>Change.idx</i> (Index of the maximum KL score)	1.07			-14.08	0.84	1.03	
18	<i>Correlation (bts-top)</i>		-38.55	-32.91	-74.67			
19	<i>Aggregate share</i>	4.57	31.55	35.95	45.40	6.84	4.42	3.66
20	<i>Coefficient of variation</i>	7.83	94.92	42.80	227.73	7.58	7.86	8.36
	Adjusted R ²	0.6739	0.4861	0.5376	0.4003	0.6487	0.6692	0.056

In Table 5, coefficients are presented only when the effect of each time series characteristics is statistically significant ($p\text{-value} < 0.05$)^k. Each column represents a separate regression model created for the different HF model. For different HF models, different characteristics found to be significant. Table 5 shows that *Fspots*, *linearity*, *aggregate share* and *coefficient of variation* are among time

^j This is not a time series characteristic. It is the intercept of the regression model.

^k We restrict extrapolation of the findings regarding the influence of time series characteristics on the forecasting performance, only on those time series which have similar or closely related summary statistics to the data provided in Table 5 and Fig. A.1.

series characteristics that impact the performance of all HF models in the bottom level. Positive coefficients indicate that they contribute to the increase of the forecasting error, while negative indicates the effect of decreasing the forecasting error. Therefore HF models have a tendency of producing more inaccurate forecasts while forecasting the time series with a higher values of *Fspots*, *coefficient of variation* and *aggregate share*. In contrast, higher values of *linearity*, increase forecast accuracy of all HF models. *Entropy* has the same effect, and it reduces the forecasting error, except it is not significant for TDFP model. *ACF1* and *curvature* decrease the forecast accuracy. However, their influence is much smaller and it is not significant for all HF models. Moreover, *Spikiness*, *season*, *trough* and *correlation (bts-top)* characteristics might also improve forecast accuracy. *Lshift*, *Cpoints* and *Change.idx* might also deteriorate forecast accuracy for all HF models, except for TD1, TD2 and TD3, for which they demonstrate an improvement.

The last row in Table 6 represents the adjusted R^2 , i.e. the percentage of the RMSEs variability of HF error models explained by given additive linear regression models. The adjusted R^2 ranges from 5.6% of the explained error variability for the TDFP model, to the high 67.39% of the explained error variability for BU.

7. CONCLUSION AND FURTHER RESEARCH

Various levels of forecasts are required in SC to support decisions in different departments such as logistics, marketing, manufacturing and finance. The current practice in SC is to generate separate forecasts using univariate forecasting methodologies to support different decisions. Univariate forecasting methodologies can only provide accurate forecasts for the unit for which forecasting is performed such as particular level, sector or echelon in the SC. They are not able to provide coherent forecasts across all levels or echelons of the SC. Consequently, these forecasts might be more damaging for the efficiency of the SC, than having perhaps less accurate HF/GF forecasts that are coherent across different parties of the SC. We argue that SC and HF/GF are naturally matched. In this paper, we demonstrate the application of GF methodology in a multiple-echelon distribution network of a major European brewery industry. Special emphasis is given to the design of the forecasting structure to ensure that generated forecasts are aligned with the need of all parties involved in delivering final products in the brewery distribution network. The forecasting structure consists of eleven levels and 169 nodes in total. It provides forecasts for the planning and the execution of processes in key parts of SC: manufacturing, marketing, finance, and logistics.

In this paper, we also considered the fact that there is no agreement on which HF approach provides the most accurate forecast. Therefore, we evaluate the effectiveness of BU, TD and the OC approaches in the simulation study, and examine whether a combination of these models improves

the forecast accuracy (COMB and COMBw). Moreover, we investigate the impact of time series characteristics on the effectiveness of each HF approach.

The main findings of this paper can be summarized as follows:

- First, OC (with minimum trace estimator) and BU outperform all TD approaches on average and across all levels. Therefore, we recommend practitioners to use these approaches when generating grouped demand forecasts across various levels of supply chain. Moreover, we notice that OC and BU approaches demonstrate robustness and consistency in producing stable and accurate forecasts, regardless of the hierarchy level and time series characteristics. This is a very important result for practitioners as BU shows to be very competitive given its simplicity.
- Second, in comparing various variations of TD approaches, we observe that TD4 and TDFP outperform other TD methodologies considered in this study. These approaches also seem to be more robust to the level of hierarchies.
- Third, our results show that forecast combination of the existing BU, TD and OC models improve the forecast accuracy. We propose two simple combination approaches based on existing models. They are at least as accurate as forecasts generated by BU and OC approaches. Therefore, when dealing with hierarchical and grouped demand structures in a SC, we recommend practitioners to use a combination of models instead of using individual approaches.
- Fourth, we examine the effect of time series characteristics on the forecasting performance of different approaches at the bottom level of the hierarchy. We show that higher values of the *Fspots*, *aggregate share* and the *coefficient of variation* may have a negative impact and deteriorate the forecast accuracy performance of HF models. In contrast, *entropy* and *linearity* may have a positive impact and improve the forecast accuracy of HF models. Other time series characteristics have a much smaller effect on all HF models.
- Finally, we demonstrate the application of grouped forecasting in SC using a multi-echelon distribution network of a major European brewery company. We empirically present the holistic approach for designing the forecasting structure while forecasting the demand in the SC. The structure serves as the information platform to support the planning and execution of processes in manufacturing, marketing, finances and logistics.

As far as the next steps of research are concerned, further work into the following areas would appear to be merited:

- The interface between temporal and hierarchical or grouped aggregation has received minimal attention in both academia and industry and this is an issue that we plan to investigate in the next

steps of our research. Creating a HF/GF methodology to produce the forecasts in the hierarchies or grouped structures with horizontal connections between nodes within the same level is an interesting avenue for further research. Current models can only produce forecasts in the structures, which have vertical connections between nodes in different levels. This kind of hierarchies or groups with horizontal and vertical connections are useful in the SC.

- The impact of time series characteristics from all levels in the forecasting structure and creating an algorithm to link the time series characteristics of the entire forecasting structure to the accuracy of each HF model is another interesting avenue for the future research.
- TD approach might be improved by considering a new disaggregation method. Standard TD approaches consider the disaggregating proportions as a static value, however observing the disaggregating proportions as functions in time, might improve the forecast accuracy of the TD approach. We would like to investigate this in a SC structure in the future works. This approach may use the forecasted proportions principle and still be used not just in hierarchical but also in grouped data structures.

REFERENCES

- Aigner, D.J. and Goldfeld, S.M. (1973), "Simulation and aggregation: a reconsideration", *The Review of Economics and Statistics*, pp.114-118.
- Athanasopoulos, G., Ahmed, R.A. and Hyndman, R.J. (2009), "Hierarchical forecasts for Australian domestic tourism", *International Journal of Forecasting*, Vol. 25 No. 1, pp.146-166.
- Babai, Z., Ali, M. and Nikolopoulos, K. (2012), "Impact of temporal aggregation on stock control performance of intermittent demand estimators: Empirical analysis", *Omega*, Vol. 40 No. 6, pp.713-721.
- Ballou, R.H. (2004), *Business Logistics/Supply Chain Management-Planning, Organizing, and Controlling the Supply Chain*, Pearson/Prentice Hall, New York.
- Barnea, A. and Lakonishok, J. (1980), "An analysis of the usefulness of disaggregated accounting data for forecasts of corporate performance", *Decision Sciences*, Vol. 11 No. 1, pp.17-26.
- Boylan, J. (2010), "Choosing levels of aggregation for supply chain forecasts", *Foresight: The International Journal of Applied Forecasting*, No. 18, pp.9-13.
- Caplice, C. and Sheffi, Y. (2006), "ESD.260J Logistics Systems", working paper, Massachusetts Institute of Technology, United States.
- Chen, A. and Blue, J. (2010), "Performance analysis of demand planning approaches for aggregating, forecasting and disaggregating interrelated demands", *International Journal of Production Economics*, Vol. 128 No. 2, pp.586-602.
- Chen, A., Yang, K. and Hsia, Z. (2008), "Weighted least-square estimation of demand product mix and its applications to semiconductor demand", *International Journal of Production Research*, Vol. 46 No. 16, pp.4445-4462.
- Chen, H. and Boylan, J.E. (2007), "Use of individual and group seasonal indices in subaggregate demand forecasting", *Journal of the Operational Research Society*, Vol. 58 No. 12, pp.1660-1671.
- Chopra, S. and Meindl, P. (2007), *Supply chain management-Strategy, Planning, and Operation*, Pearson Prentice Hall, New Jersey.
- Collins, D.W. (1976), "Predicting earnings with sub-entity data: Some further evidence", *Journal of Accounting Research*, pp.163-177.

- Dangerfield, B.J. and Morris, J.S. (1992), "Top-down or bottom-up: Aggregate versus disaggregate extrapolations", *International Journal of Forecasting*, Vol. 8 No. 2, pp.233-241.
- Dunn, D.M., Williams, W.H. and DeChaine, T. (1976), "Aggregate versus subaggregate models in local area forecasting", *Journal of the American Statistical Association*, Vol. 71 No. 353, pp.68-71.
- Dunn, D.M., Williams, W.H. and Spivey, W.A. (1971), "Analysis and prediction of telephone demand in local geographical areas", *The Bell Journal of Economics and Management Science*, pp.561-576.
- Edwards, J.B. and Orcutt, G.H. (1969), "Should aggregation prior to estimation be the rule?", *The Review of Economics and Statistics*, pp.409-420.
- Flidner, G. (1999), "An investigation of aggregate variable time series forecast strategies with specific subaggregate time series statistical correlation", *Computers & Operations Research*, Vol. 26 No. 10, pp.1133-1149.
- Flidner, G. (2001), "Hierarchical forecasting: issues and use guidelines", *Industrial Management & Data Systems*, Vol. 101 No. 1, pp.5-12.
- Forslund, H. and Jonsson, P. (2007), "The impact of forecast information quality on supply chain performance", *International Journal of Operations & Production Management*, Vol. 27 No. 1, pp.90-107.
- Gordon, T.P., Morris, J.S. and Dangerfield, B.J. (1997), "Top-down or bottom-up: Which is the best approach to forecasting?", *The Journal of Business Forecasting*, Vol. 16 No. 3, pp.13-16.
- Gross, C.W. and Sohl, J.E. (1990), "Disaggregation methods to expedite product line forecasting", *Journal of Forecasting*, Vol. 9 No. 3, pp.233-254.
- Grunfeld, Y. and Griliches, Z. (1960), "Is aggregation necessarily bad?", *The Review of Economics and Statistics*, pp.1-13.
- Hyndman, R.J., Ahmed, R.A. and Athanasopoulos, G. (2007), "Optimal combination forecasts for hierarchical time series", working paper, MONASH University, 17 July 2007.
- 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*, Vol. 55 No. 9, pp.2579–2589.
- Hyndman, R.J. and Athanasopoulos, G. (2014), *Forecasting: principles and practice*, OTexts,
- Hyndman, R.J. and Athanasopoulos, G. (2018), *Forecasting: principles and practice*, OTexts,
- Hyndman, R.J., Lee, A.J. and Wang, E. (2016), "Fast computation of reconciled forecasts for hierarchical and grouped time series", *Computational Statistics & Data Analysis*, Vol. 97, pp.16-32.
- Hyndman, R.J., Wang, E. and Laptev, N. (2015), "Large-scale unusual time series detection", in *Data Mining Workshop (ICDMW), 2015 IEEE International Conference - Yunnan, China*, IEEE, pp.1616-1619.
- Kahn, K.B. (1998), "Revisiting top-down versus bottom-up forecasting", *The Journal of Business Forecasting*, Vol. 17 No. 2, pp.14-19.
- Kalchschmidt, M., Verganti, R. and Zotteri, G. (2006), "Forecasting demand from heterogeneous customers", *International Journal of Operations & Production Management*, Vol. 26 No. 6, pp.619-638.
- Kinney, W.R. (1971), "Predicting earnings: entity versus subentity data", *Journal of Accounting Research*, pp.127-136.
- Mircetic, D. (2018), *Boosting the performance of top down methodology for forecasting in supply chains via a new approach for determining disaggregating proportions*, PhD University of Novi Sad, Serbia.
- Mircetic, D., Nikolicic, S., Stojanovic, D. and Maslaric, M. (2017), "Modified top down approach for hierarchical forecasting in a beverage supply chain", *Transportation research procedia*, Vol. 22, pp.193-202.
- Pennings, C.L. and van Dalen, J. (2017), "Integrated hierarchical forecasting", *European Journal of Operational Research*, Vol. 263 No. 2, pp.412-418.
- Rostami-Tabar, B. (2013), *ARIMA demand forecasting by aggregation*, PhD Université Sciences et Technologies-Bordeaux I.

- Rostami-Tabar, B., Babai, M.Z., Ducq, Y. and Syntetos, A. (2015), "Non-stationary demand forecasting by cross-sectional aggregation", *International Journal of Production Economics*, Vol. 170, Part A, pp.297-309.
- Schwarzkopf, A.B., Tersine, R.J. and Morris, J.S. (1988), "Top-down versus bottom-up forecasting strategies", *The International Journal Of Production Research*, Vol. 26 No. 11, pp.1833-1843.
- Seongmin, M., 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*, Vol. 140 No. 2, pp.794-802.
- Shang, H.L. and Hyndman, R.J. (2017), "Grouped functional time series forecasting: An application to age-specific mortality rates", *Journal of Computational and Graphical Statistics*, Vol. 26 No. 2, pp.330-343.
- Strijbosch, L.W.G., Heuts, R.M.J. and Moors, J.J.A. (2008), "Hierarchical estimation as a basis for hierarchical forecasting", *IMA Journal of Management Mathematics*, Vol. 19 No. 2, pp.193-205.
- Syntetos, A., Babai, Z., Boylan, J., Kolassa, S., et al. (2016), "Supply chain forecasting: Theory, practice, their gap and the future", *European Journal of Operational Research*, Vol. 252 No. 1, pp.1-26.
- Teunter, R.H., Babai, M.Z., Bokhorst, J.A. and Syntetos, A.A. (2018), "Revisiting the value of information sharing in two-stage supply chains", *European Journal of Operational Research*, Vol. 270 No. 3, pp.1044-1052.
- Trapero, J.R., Cardos, M. and Kourentzes, N. (2019), "Empirical safety stock estimation based on kernel and GARCH models", *Omega*, Vol. 84, pp.199-211.
- Trapero, J.R., Kourentzes, N. and Fildes, R. (2012), "Impact of information exchange on supplier forecasting performance", *Omega*, Vol. 40 No. 6, pp.738-747.
- Villegas, M.A. and Pedregal, D.J. (2018), "Supply chain decision support systems based on a novel hierarchical forecasting approach", *Decision Support Systems*, Vol. 114, pp.29-36.
- Vogel, S. (2013), *Demand fulfillment in multi-stage customer hierarchies*, Springer Science & Business Media.
- Wang, X., Smith-Miles, K. and Hyndman, R. (2009), "Rule induction for forecasting method selection: Meta-learning the characteristics of univariate time series", *Neurocomputing*, Vol. 72 No. 10-12, pp.2581-2594.
- Weatherford, L.R., Kimes, S.E. and Scott, D.A. (2001), "Forecasting for hotel revenue management: Testing aggregation against disaggregation", *Cornell hotel and restaurant administration quarterly*, Vol. 42 No. 4, pp.53-64.