

Hierarchical Time Series Forecasting in Emergency Medical Services

Bahman Rostami-Tabar^{a,*}, Rob J. Hyndman^b

^aCardiff University, Cardiff Business School, United Kingdom, CF10 3EU

^bMonash University, Department of Econometrics and Business Statistics, Australia, VIC 3800

Abstract

Accurate forecasts of ambulance demand are crucial inputs when planning and deploying staff and fleet. Such demand forecasts are required at national, regional and sub-regional levels, and must take account of the nature of incidents and their priorities. These forecasts are often generated independently by different teams within the organization. As a result, forecasts at different levels may be inconsistent, resulting in conflicting decisions and a lack of coherent coordination in the service. To address this issue, we exploit the hierarchical and grouped structure of the demand time series, and apply forecast reconciliation methods to generate both point and probabilistic forecasts that are coherent and use all the available data at all levels of disaggregation. The methods are applied to daily incident data from an ambulance service in Great Britain, from October 2015 to July 2019, disaggregated by nature of incident, priority, managing health board, and control area. We use an ensemble of forecasting models, and show that the resulting forecasts are better than any individual forecasting model. We validate the forecasting approach using time-series cross-validation.

Keywords: healthcare, emergency services, forecast reconciliation, ambulance demand, regression

1. Introduction

A failure to match available resources to demand in Emergency Medical Services (EMS) results in patient flow problems, with serious consequences for patients, staff and the entire care system (Ekström et al., 2015; Rostami-Tabar and Ziel, 2022). Demand forecasting in EMS helps service planners to avoid the mismatch, potentially providing massive savings in costs and lives, and leading to better patient outcomes. Accurate daily demand forecasting enables planners and decision makers to manage resources to meet anticipated patients, reconfigure units, and redeploy staff and vehicles as necessary.

Demand forecasts at EMS are typically required at multiple levels of an organization to inform various planning and decision-making processes (Hulshof et al., 2012). There are some planning process at the national level (strategic and long-term) such as workforce resource planning and budgeting; sub-national, regional or healthcare level (tactical and medium-term) such as temporary capacity expansions, resource sharing and staff-shift scheduling; and hospital or station level (operational and short-term) such as planning rosters for staff and ambulance deployment. Demand forecasts might also be required at different levels for a specific area of interest such as the nature of demand or the priority level. Moreover, the time series data in EMS has an inherent hierarchical and grouped structure to support such forecasting requirements. Demand for emergency medical services at the country level can be disaggregated in a geographical

*Corresponding author

Email addresses: rostami-tabarb@cardiff.ac.uk (Bahman Rostami-Tabar), Rob.Hyndman@monash.edu (Rob J. Hyndman)

hierarchy into sub-national, regions, health boards, and stations/hospitals, or divided into groups such as the nature of incidents or demand priority. Therefore, using forecasting methodologies that account for hierarchical and/or grouped structures of time series in EMS is a natural fit.

However, despite a large number of studies dedicated to forecasting for EMS (Shi et al., 2022; Gul and Celik, 2020; Ibrahim et al., 2016; Wargon et al., 2009), the hierarchical data structure has been largely ignored, and the main focus has been on producing independent forecasts at a single level. Generating independent forecasts can result in a lack of consistency and coordination, and therefore leads to less effective planning and decision making. With hierarchical forecasting, plans at any level are based on coherent forecasts and therefore can be aligned. Implementing and sustaining improvements in EMS require alignments and coordination between different stakeholders, without which teams operate in isolation leading to conflicts, duplication work, rework, or work that runs counter to the overall goal to improve the quality of delivery service. Hierarchical forecasting framework can be used as a tool to improve coordination between teams across the care services at the national, sub-national, regional and local levels. The hierarchical forecasting approaches not only create consistent forecasts, but are usually also more accurate than the independent (base) forecasts (Hyndman et al., 2011). To our knowledge, there has been no previous research involving hierarchical and grouped forecasting in the entire field of forecasting for healthcare management.

In this paper, we address this gap by investigating the application of hierarchical forecasting approaches in the EMS using daily time series of verified incidents from 2015 to 2020 in a major ambulance service in Great Britain. The data has hierarchical and grouped structures, with hierarchies at the national, control (i.e. sub-national), health board (i.e. regional) levels, as well as groups by priority and nature of incidents. We produce consistent point forecasts and forecast distributions for all levels, which is critical for an effective planning and associated risk management. We compare the point and probabilistic forecast accuracy of the independent forecasts, bottom-up and optimal reconciliation approaches. We first generate independent/base forecasts using Exponential Smoothing State Space (ETS), Generalized Linear Model (GLM), Poisson regression, a simple empirical distribution and an ensemble method, followed by applying bottom-up and optimal reconciliation approaches. Forecast performance is assessed by the Mean Squared Scaled Error (MSSE) for point forecasts and Continuous Ranked Probability Scores (CRPS) for the probabilistic forecasts. This paper complies with reproducibility principles (Stodden and Miguez, 2013; Boylan et al., 2015), and can be applied to any healthcare service (e.g., emergency department, primary or social care) subject to the time series having a hierarchical and/or grouped structure.

The remainder of this article is structured as follows: In Section 2, we provide a brief review of the literature and discuss its limitation to position our work; in Section 3, we present the experiment design describing the data set, forecasting methods and forecast evaluation metrics. In Section 4, we discuss the hierarchical time series forecasting approaches to generate both point and probabilistic forecasts. In Section 5, we present and discuss our results; in Section 6, we summarize our findings and present ideas for future research.

2. Research background

Emergency medical services (EMS) are a critical component in the delivery of urgent medical care to communities. An effective service delivery requires accurate resource planning that generally relies on demand forecasts at operational, tactical and strategic levels.

There is a substantial number of studies on the application of time series forecasting in the Emergency Medical Services. For example, Ibrahim et al. (2016) provide an extensive review of the models used in forecasting call volume arrivals. Another important area is related to forecasting ambulance demand. Although the definition of demand might not be always clearly stated, this is typically referring to a situation where a physical resource has been deployed to respond to an incident. This might be also called *attended incidents*. Another demand related variable is verified incidents; these are all incidents that require an action: either by sending a physical vehicle, responding via the Clinical Support Desk, requesting an external provider

to respond to it, or forwarding it to other channels such as police, firefighters or general practitioners. Our study is aligned with this stream of the literature. Another similar area that has received considerable attention is Emergency Department forecasting; we refer interested readers to [Shi et al. \(2022\)](#), [Gul and Celik \(2020\)](#) and [Wargon et al. \(2009\)](#) for extensive reviews of the relevant literature. Although crucial to EMS performance, [Aringhieri et al. \(2017\)](#) state that demand forecasting has received limited research attention in the EMS context. In this section, we provide a brief review of studies on forecasting ambulance demand in EMS.

There are generally two main streams of research related to forecasting ambulance demand in EMS: (i) the first stream focuses on the application of time series methods and regression approaches to forecasting aggregate ambulance demand ([Vile et al., 2012](#); [Sasaki et al., 2010](#)); and (ii) the second stream considers forecasting EMS demand in finer temporal and geographical granularities by employing temporal-spatial prediction methods ([Zhou and Matteson, 2016](#); [Zhou, 2016](#)). The focus of our study is related to the first stream of research.

[Sasaki et al. \(2010\)](#) develop a multivariable regression model to estimate future EMS demands. In addition to the historical demand, the population census for different age groups and counts of the number of companies employing more than five people are included in the regression. The census variables describe groups who are more likely to need an ambulance. A stepwise ordinary least squares regression analysis is used for estimating the parameter and generating forecast. The only performance measure reported in this study is R^2 , which is not an effective measure of forecast accuracy ([Armstrong, 2001](#), p457). The research design of this study is not rigorous and the study is not reproducible. [Vile et al. \(2012\)](#) explore using a Singular Spectrum Analysis (SSA) method to generate forecasts of the EMS demand at the national level for 7-day, 14-day, 21-day and 28-day forecast horizons using data provided by an ambulance service in Great Britain. The performance of this approach is compared to Auto-Regressive Integrated Moving Average (ARIMA) and Holt-Winters time series methods using Root Mean Squared Error (RMSE). They concluded that point forecasts generated by SSA are more accurate for longer-term, but that ARIMA and Holt-Winters performance is superior for shorter-term horizons. [Vile et al. \(2016\)](#) further develop a decision support system to integrate forecasts generated by SSA. However, the study does not compare and contrast the performance of forecasting methods based on utility measures such as cost, resource utilization or response time. The tool contains options that allow generating forecasts at various levels of granularity, however, it ignores the hierarchical and grouped relationships structure, preventing aligned decision making and coordination.

[Haugsbø Hermansen and Mengshoel \(2021\)](#) investigate forecasting EMS demand in a high spatio-temporal resolution of 1km^2 spatial regions and 1-hr time intervals using total incidents in Oslo, Norway, from 1 January 2015 to 11 February 2019. They used multi-layer perceptron (MLP) and long short-term memory (LSTM) models to forecast the EMS demand, and compare the results to simple aggregation methods and baselines. The point forecast accuracy is evaluated using Mean Absolute Error (MAE) and Mean Squared Error (MSE), and the forecast distribution is measured by Categorical Cross-Entropy. They found that Neural Network models performed better in producing point forecasts, while a distribution baseline method based on spatial distribution of the incidents across all time steps provided more accurate forecast distributions. [Zhou \(2016\)](#) proposed three methods based on Gaussian mixture models, kernel density estimation, and kernel warping to predict hourly data 4 weeks ahead for a 1km^2 spatial region. Two years of incidents from Toronto, Canada (years 2007 and 2008 with 391,296 events) and Melbourne, Australia (years 2011 and 2012 with 696,975 events) were used to build the model and examine the performance on test data using mean negative log likelihood. They show that forecasts generated by the proposed methods were significantly more accurate than the current industry practice (a simple averaging formula). [Grekousis and Liu \(2019\)](#) investigated the combination of spatial analysis methods with data mining techniques based on an improved Hungarian algorithm and a MLP neural network to identify the most likely locations of future emergency events. The proposed approach was tested using data of 2851 events attended by the EMS in Athens, Greece, over 24 weeks. They showed that 23% of real emergency events lie within 50 meters of the predicted ones and nearly 70% of the real emergency events lie no further than 150 meters away, which is

rather accurate given the granularity of the problem at the city level.

We note a number of limitations in the literature of EMS forecasting, that encourage us to undertake this research. These limitations are summarized as following:

1. Current studies ignore the inherent hierarchical and/or grouped structure of the time series data, and the relationship between series at different levels of hierarchy. This may result in incoherent forecasts leading to misaligned planning and decision making. While the hierarchical forecasting methodology has been developed and applied in various domains over the past 10 years ([Panagiotelis et al., 2022](#)), it has never been explored in this area.
2. Current research is mainly concerned with generating point forecasts at a single level of hierarchy. There is a lack of studies considering the entire forecast distribution of daily ambulance demand for the whole hierarchy to inform the decision-making process and to better represent the uncertainty of future demand, providing a risk management tool for planners.
3. Reproducibility is still a major challenge in EMS forecasting, as it is unlikely that any reader can reproduce prior studies without the help of the authors of those papers.
4. Another limitation is related to the generated forecasts not being on the sample space of non-negative counts. Since actual ambulance counts cannot be negative or non-integer, ambulance demand forecast distributions should reflect the data. Of course, point forecasts represent means, so they should be non-negative, but may be non-integer. While this might not be an issue when producing forecasts at a single level, producing non-negative count forecasts in a hierarchical/grouped structure is challenging and requires further investigation in the future.

This paper concerns the problem of hierarchical forecasting in EMS and generates and evaluates both point and probabilistic forecast across different levels of the hierarchy, hence addressing some important gaps identified in the literature.

3. Experiment setup

We are interested in generating forecasts to inform the planning horizon of 42 days, required by planners in the ambulance service. The forecast horizon in this study is $2 \times 42 = 84$ days ahead, because the planning is generally frozen for 42 days and so forecasts of the next 42 days is not particularly helpful for planning. While forecasts are generated for 84 days ahead, performance evaluation is only assessed based on the last 42 days and not the whole forecast period. The forecasts are produced for various training and test sets using time series cross-validation ([Hyndman and Athanasopoulos, 2021](#)).

In the following section, we discuss the dataset, describe the forecasting methods used to generate base forecasts, and present the point and probabilistic accuracy measures.

3.1. Data

The dataset used in this study is from a major ambulance service in Great Britain. It contains information relating to the daily number of attended incidents from 1 October 2015 to 31 July 2019, disaggregated by nature of incidents, priority, the health board managing the service and the control area (or region). Figure 1 depicts both the hierarchical and grouped structure of the data. Figure 1a illustrates the nested hierarchical structure based on control area and health board and Figure 1b shows the grouped structure by priority and the nature of incident.

Table 1 also displays the structure of data with the total number of series at each level. At the top level, we have the total attended incidents for the country. We can split these total attended incidents by control area, by health board, by priority or by nature of incident. There are 3 control areas breakdown by 7 local health boards. Attended incident data are categorized into 3 priority classes of red, amber and green. There are

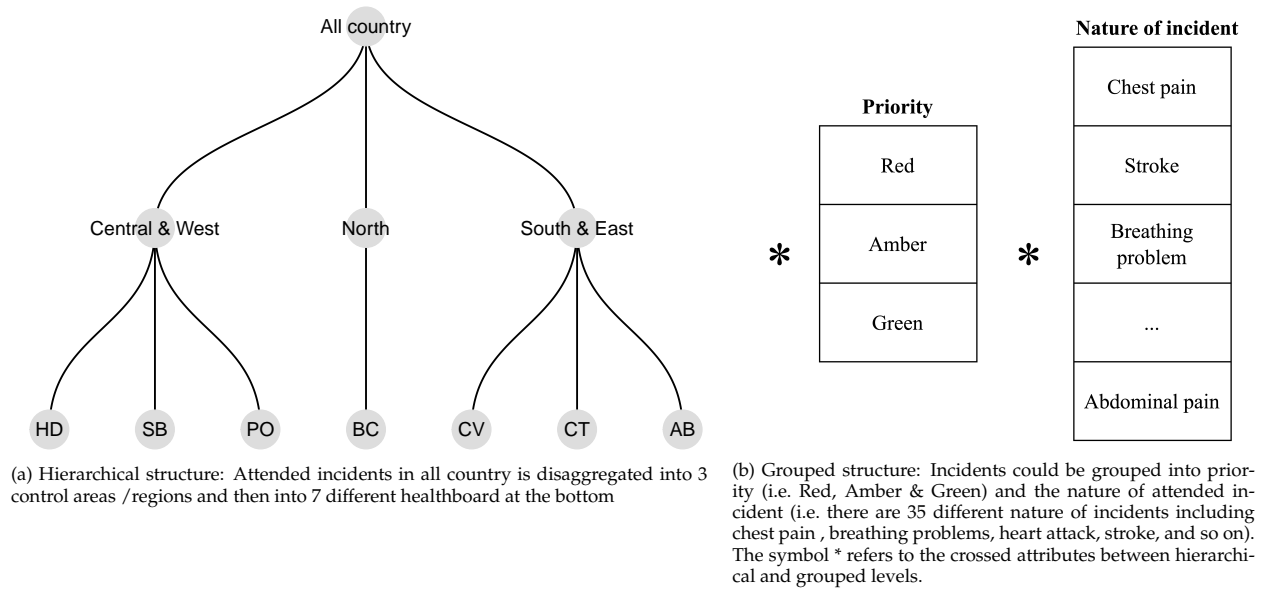


Figure 1: The hierarchical and grouped structure of attended incidents (ambulance demand).

Table 1: Number of time series in each level for the hierarchical & grouped structure of attended incidents

Level	Number of series
All country	1
Control	3
Health board	7
Priority	3
Priority * Control	9
Priority * Health board	21
Nature of incident	35
Nature of incident * Control	105
Nature of incident * Health board	245
Priority * Nature of incident	104
Control * Priority * Nature of incident	306
Control * Health board * Priority * Nature of incident (Bottom level)	691
Total	1530

also 35 different nature of incidents such as chest pain, stroke, breathing problem, etc. In total, across all levels of disaggregation, there are 1530 time series.

Given the total number of time series, direct visual analysis is infeasible. Therefore, we first compute features of all 1530 time series (Kang et al., 2017) and display the strength of trend and weekly seasonality strength in Figure 2. Each point represents one time series with the strength of trend in x-axis and the strength of seasonality in y-axis. Both measures are on a scale of $[0,1]$. It is clear that there are some series showing strong trends and/or seasonality, corresponding to series at the higher levels of the hierarchy. The majority of series show low trend and seasonality. These are time series belonging to the bottom series, series related to the nature of incidents for a given control, health board and priority level. Bottom series are dominated by noise with little or no systematic patterns.

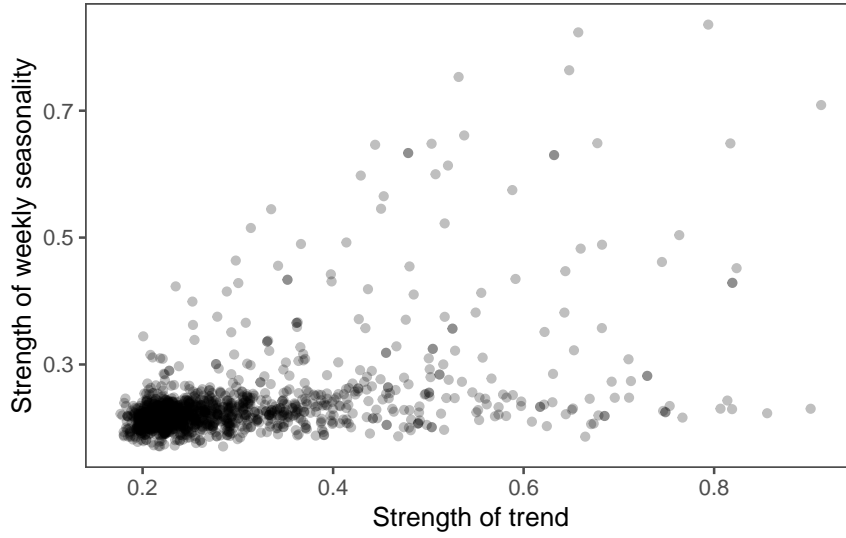


Figure 2: Time series features of attended incidents across all levels (1530 series)

In addition to displaying the trend and seasonality features, we also visualize few time series at various levels of the aggregation. Figure 3 reveals different information such as trend, seasonality and noise. For example, some series depict seasonality and trend, whereas some other series report low volume of attended incidents and entropy, making them more volatile and difficult to forecast. At the level on nature of incidents combined with categories of other levels, there are many series that contain zeros with low counts. As such, the data set represents a diverse set of daily time series patterns.

We consider several forecasting models that account for the diverse patterns of the time series across the entire hierarchy. In developing the forecasting models, the time series of holidays are also used in addition to the attended incidents. We use public holidays, school holidays and Christmas Day and New Year's Day as predictors of incident attended. These type of holidays will affect peoples' activities and may increase or decrease the number of attended incidents.

3.2. Forecasting methods

Given the presence of various significant patterns in the past attended incidents, we consider three different forecasting models to generate the base forecasts. Once the base forecasts are produced, hierarchical and grouped time series methods are used to reconcile them across all levels. We briefly discuss forecasting models in the following sections, and the hierarchical forecasting methods are discussed in Section 4.

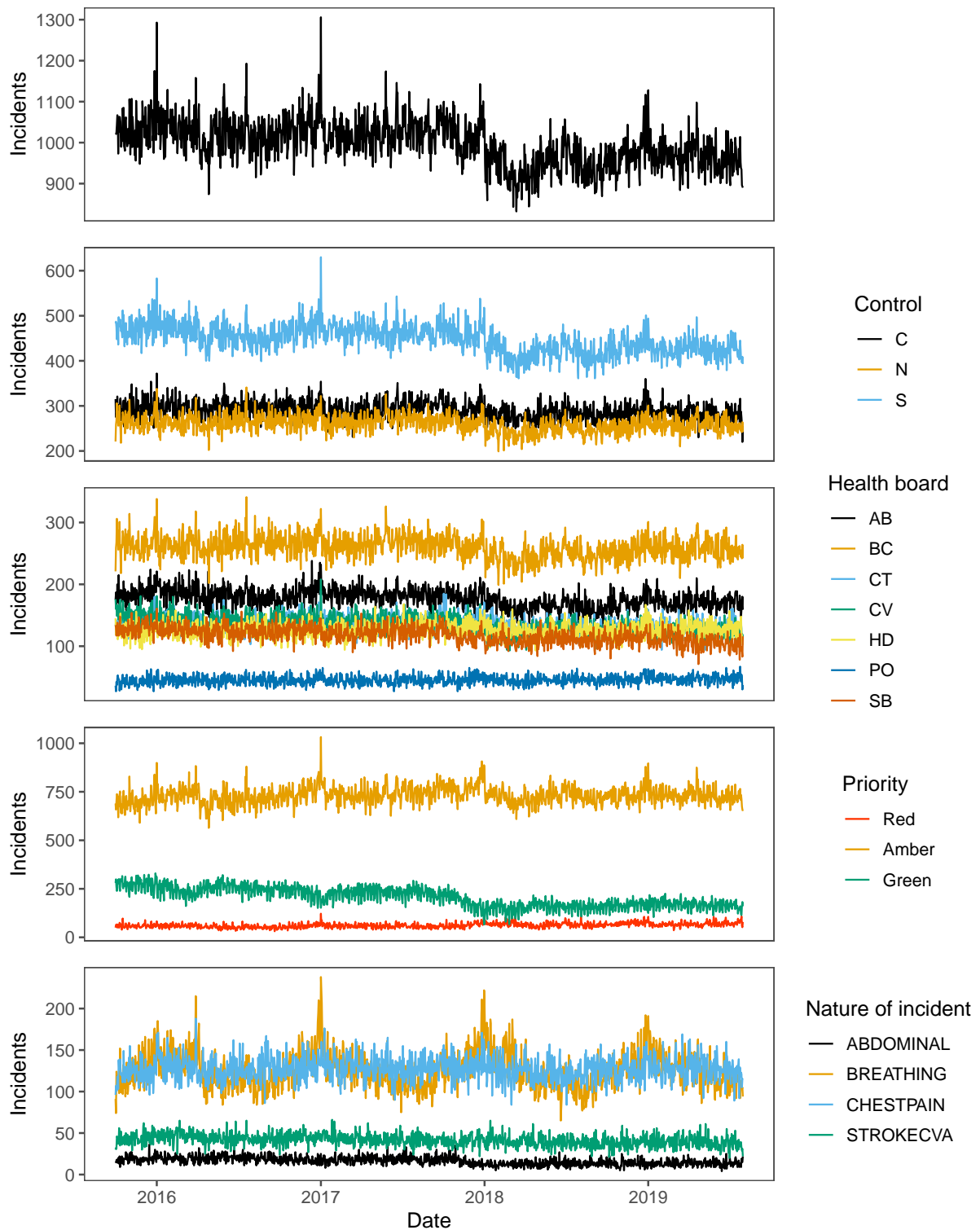


Figure 3: Time series of attended incidents at various levels. The panels show data from the whole country, by control area, by health board, by priority level, and by nature of incident. Only four of the 35 nature of incident categories are shown to avoid too much overplotting.

Naive: We start with a simple forecasting approach, assuming that the future days will be similar to past days. We use the empirical distribution of the past daily attended incidents to create the forecast distribution of future attended incidents.

Exponential Smoothing State Space model (ETS): ETS models (Hyndman and Athanasopoulos, 2021) can combine trend, seasonality and error components in a time series through various forms that can be additive, multiplicative or mixed. The trend component can be none (“N”), Additive (“A”) or damped (“Ad”); the seasonality can be none (“N”), Additive (“A”), or multiplicative (“M”); and the error term can be additive (“A”) or multiplicative (“M”). To forecast the attended incidents at each level, we use the `ETS()` function in the `fable` package (O’Hara-Wild et al., 2022) in R. To identify the best model for a given time series, the ETS function uses the corrected Akaike’s Information Criterion (AICc).

Despite the popularity and the relevance of automatic ETS in this study, it may produce forecast distributions that are non-integer and include negative values, although the number of attended incidents is always integer and non-negative. When using ETS, a time series transformation approach could be used to generate strictly positive forecasts, although forecast distributions will still be non-integer. An alternative is to use forecasting models that produce integer, non-negative forecasts. In the following section we present Generalized Linear Models (GLMs) and Poisson time series regression to produce count base forecasts.

Generalized Linear Model (GLM): GLMs are a family of models developed to extend the concept of linear regression models to non-Gaussian distributions (Faraway, 2016). They model the response variable as a particular member of the exponential family, with the mean being a transformation of a linear function of the predictors. One of the models that is frequently used in practice to generate count forecasts is Poisson regression. We will consider forecasting attended incidents using the covariates spline trend, day of the week dummy variables (from Monday to Sunday), Fourier terms to capture yearly seasonality, public holidays (1 when is public holiday, 0 otherwise), school holidays (1 when is school holiday, 0 otherwise) and Christmas Day (1 when is Christmas Day, 0 otherwise) and New Year’s Day (1 when is New Year’s Day, 0 otherwise). We fit a Poisson regression model using the function `glm()` from the `stats` package in R, with the argument `family = poisson` to specify that we wish to fit a Poisson regression model with a log link function.

Poisson Regression using `tscount` (TSGLM): We also consider another Poisson regression model that takes into account serial dependence, using the `tsglm()` function in the `tscount` package in R (Liboschik et al., 2017) to model the attended incidents. Again, the logarithmic link function is used to ensure that the parameter of Poisson distribution is always positive. This model captures the short range serial dependence by including three autoregressive terms, in addition to the same covariates that were used in the GLM model. To distinguish this from the previous GLM model, we will refer to this model as “TSGLM”.

Ensemble method: Finally, we use an ensemble method that combines the forecasts generated from the Naive, ETS, GLM and TSGLM models to form a mixture distribution (Wang et al., 2022).

3.3. Performance evaluation

To evaluate the performance of the various forecasting approaches, we split the data into a series of ten training and test sets. We use a time series cross-validation approach (Hyndman and Athanasopoulos, 2021), with a forecast horizon of 84 days, and each training set expanding in 42-day steps. The first training set uses all data up to 2018-04-25, and the first test set uses the 84 days beginning 2018-04-26. The second training set uses all data up to 2018-06-06, with the second test set using the following 84 days. The largest training set ends on 2019-05-09, with the test set ending on 2019-07-31. Model development and hyperparameter tuning is performed using the training data only. While we compute forecast errors for the entire 12 weeks, we are most interested in the last 42 days of each test set, because that corresponds to how forecasts are generated for planning in practice. Forecasting performance is evaluated using both point and probabilistic error measures.

Point forecast accuracy is measured via the Mean Squared Scaled Error (MSSE) and the Mean Absolute Scaled Error (MASE). The Mean Absolute Scaled Error (MASE) (Hyndman and Koehler, 2006) is calculated as:

$$\text{MASE} = \text{mean}(|q_j|),$$

where

$$q_j = \frac{e_j}{\frac{1}{T-m} \sum_{t=m+1}^T |y_t - y_{t-m}|},$$

and e_j is the point forecast error for forecast horizon j , $m = 7$ (as we have daily seasonal series), y_t is the observation for period t , and T is the sample size (the number of observations used for training the forecasting model). The denominator is the mean absolute error of the seasonal naive method in the fitting sample of T observations and is used to scale the error. Smaller MASE values suggest more accurate forecasts. Note that the measure is scale-independent, thus allowing us to average the results across series.

A related measure is MSSE (Hyndman and Athanasopoulos, 2021), which uses squared errors rather than absolute errors:

$$\text{MSSE} = \text{mean}(q_j^2),$$

where,

$$q_j^2 = \frac{e_j^2}{\frac{1}{T-m} \sum_{t=m+1}^T (y_t - y_{t-m})^2},$$

Again, this is scale-independent, and smaller MSSE values suggest more accurate forecasts.

To measure the forecast distribution accuracy, we calculate the Continuous Rank Probability Score (Gneiting and Katzfuss, 2014). It rewards sharpness and penalizes miscalibration, so it measures overall performance of the forecast distribution.

$$\text{CRPS} = \text{mean}(p_j),$$

where

$$p_j = \int_{-\infty}^{\infty} (G_j(x) - F_j(x))^2 dx,$$

where $G_j(x)$ is the forecasted probability distribution function for forecast horizon j , and $F_j(x)$ is the true probability distribution function for the same period.

4. Hierarchical and grouped time series forecasting techniques

There are many applications in healthcare, and in particular in EMS, where a collection of time series is available. These series are generally hierarchically organized based on multiple levels such as area/region, health board and/or are aggregated at different levels in groups based on nature of demand, priority of demand, or some other attributes. While series could be strictly hierarchical or only grouped based on some attributes, in many situation a more complex structures arise when attributes of interest are both nested and crossed, having hierarchical and grouped structure. This is also the case for our application as discussed in Section 3.1.

4.1. Independent (base forecast)

A common practice in healthcare (and EMS) to predict hierarchical and grouped series relies on producing independent forecasts, also refereed to as base forecasts, typically by different teams as the need for such forecasts arise. We observe n time series at time t , across the entire hierarchical and grouped structure, written as y_t . The base forecasts of y_{T+h} given data y_1, \dots, y_T are denoted by \hat{y}_h for h steps-ahead for all n series ($n = 1530$ in this study). Forecasts generated in this way are not coherent.

4.2. Reconciliation methos

Traditionally, approaches to produce coherent forecasts for hierarchical and grouped time series involve using bottom-up and top-down methods by generating forecasts at a single level and then aggregating or disaggregating. Top-down methods require having a unique hierarchical structure to disaggregate forecasts generated at the top level by proportions. However, given that we have multiple grouped attributes combined with the hierarchical structure, there is no unique way to disaggregate top forecasts. Hence the top-down cannot be used in our application. The recommended approach is to use forecast reconciliation (Hyndman et al., 2011). In the following sections, we first discuss some notation, and then present bottom-up and forecast reconciliation approaches used in this study to generate coherent forecasts.

4.2.1. Notations

Let \mathbf{b}_t be a vector of n_b “bottom-level” time series at time t , and let \mathbf{a}_t be a corresponding vector of $n_a = n - n_b$ aggregated time series, where

$$\mathbf{a}_t = \mathbf{A}\mathbf{b}_t,$$

and \mathbf{A} is the $n_a \times n_b$ “aggregation” matrix specifying how the bottom-level series \mathbf{b}_t are to be aggregated to form \mathbf{a}_t . The full vector of time series is given by

$$\mathbf{y}_t = \begin{bmatrix} \mathbf{a}_t \\ \mathbf{b}_t \end{bmatrix}.$$

This leads to the $n \times n_b$ “summing” or “structural” matrix given by

$$\mathbf{S} = \begin{bmatrix} \mathbf{A} \\ \mathbf{I}_{n_b} \end{bmatrix}$$

such that $\mathbf{y}_t = \mathbf{S}\mathbf{b}_t$.

4.2.2. Bottom-up (BU) and linear reconciliation methods

Bottom-Up is a simple approach to generate coherent forecasts. It involves first creating the base forecasts for the bottom level series (i.e., the most disaggregated series). These forecasts are then aggregated to the upper levels which naturally results in coherent forecasts. The BU approach can capture the dynamics of the series at the bottom level, but these series may be noisy and difficult to forecast. The approach using only the data at the most disaggregated level, and so does not utilize all the information available across the hierarchical and grouped structure.

Forecast reconciliation approaches fill this gap by combining and reconciling all the base forecasts in order to produce coherent forecasts. Linear reconciliation methods can be written (Wickramasuriya et al., 2019) as

$$\hat{\mathbf{y}}_h = \mathbf{S}(\mathbf{S}'\mathbf{W}^{-1}\mathbf{S})^{-1}\mathbf{W}^{-1}\hat{\mathbf{y}}_h,$$

where \mathbf{W} is an $n \times n$ positive definite matrix, and $\hat{\mathbf{y}}_h$ contains the h -step forecasts of \mathbf{y}_{T+h} given data to time T . Different choices for \mathbf{W} lead to different solutions such as Ordinary Least Squares (OLS), Weighted Least Squares (WLS) and Minimum Trace (MinT). We use the implementation of these methods in the fable package in R in the experiment.

5. Results and discussion

In this section, we compare the forecasting performance of the Naive, ETS, GLM and TSGLM models along with the ensemble, using base forecast and Minimum Trace (MinT) reconciliation methods. We have

also computed the forecast accuracy for Ordinary Least Square (OLS) and Weighted Least Square (WLS) approaches, along with bottom up forecasting. However, they are not reported here because their accuracy is outperformed by MinT. We should also note that forecasts, and consequently their corresponding errors, are generated for the entire hierarchy and they could be reported at any level, if required. But to save space, we have reported only the top level (Total), the bottom level, and the levels corresponding to Control areas and Health boards. The latter are chosen because this is where decision-making takes place, so these forecasts are the most important.

The overall forecasting performance is reported in Table 2, in which the average forecast accuracy over horizons 43–84 days (corresponding to the planning horizon) is presented per model, method and the hierarchical level. Reported forecast accuracy are averaged across all forecast horizons, rolling origins and series at each level. Table 2 presents both point and probabilistic forecast accuracy at total, control area, health board and bottom level series. Point forecast performance are reported using MASE and MSSE in Table 2a and Table 2b, respectively. Probabilistic forecast accuracy is reported using CRPS in Table 2c. The bold entries in each table identify a combination of method and model that performs best for the corresponding level (i.e. each column), based on the smallest values of accuracy measures.

Table 2a and 2b show that forecast reconciliation (i.e. MinT) improves forecast accuracy at the higher levels of the hierarchy including total, control area and health board. However, it does not result in accuracy improvement at the bottom level series, for which base forecasts are more accurate. This might be due to the noisy structure of time series at the bottom level, and possibly due to very different patterns in the aggregated series. It is also clear from Table 2a that the ensemble method improves forecast accuracy at total, control area and health board. However, this does not remain valid for bottom series where different individual methods perform best, depending on the accuracy measure.

Table 2c presents the accuracy of the forecast distributions measures by CRPS, which considers both forecasting reliability and interval sharpness. The smaller the value of CRPS, the better the comprehensive performance. We observe that forecast reconciliation results in forecast improvement, regardless of the hierarchical level. The ensemble method is also more accurate for higher levels, but ETS performs slightly better at the bottom level. While these results show that forecast reconciliation does not improve *point* forecasts at the bottom level, Table 2c indicates that it generates more accurate *distributional* forecasts than the base method. This is probably due to the reconciliation method giving improved forecast accuracy in the tails of the forecast distribution, which are critical for managing risks.

Overall, our results indicate that forecast reconciliation using the MinT method provides reliable forecasts, improves upon the base (unreconciled) forecasts all levels except the bottom level series. But even there, forecast reconciliation using MinT improves accuracy in the tails of the distribution.

In addition to the overall forecast accuracy presented in Table 2, we also report the point and probabilistic forecast accuracy measures for each forecast horizon in Figure 4. The figure focuses on the hierarchical levels important for decision-making including total, control area and health board; however the accuracy could be calculated for any level. We only illustrate the results of the MinT method, given its strong performance described in Table 2. For illustration purposes, we report the average weekly forecast accuracy instead of the daily forecast horizon, as this reduces the visual noise in the figure. Thus the x-axis shows horizons from week 1 ($h = 1, \dots, 7$) to week 12 ($h = 78, \dots, 84$). The forecast horizon from week 7 to week 12 corresponds to the upcoming planning horizon, which is used by planners and decision making. For both the point forecast and distributional accuracy we can see that the ensemble approach performs best across almost all horizons, with the biggest differences at the highest levels of aggregation. It is important to highlight that, all forecasting models outperform the Naive empirical distribution that is used as a benchmark for both point and probabilistic forecasts.

Table 2: Average forecast performance calculated on the test sets at forecast horizons $h = 43, \dots, 84$ days, with time series cross validation applied to attended incident data. The best approach is highlighted in bold.

(a) Point forecast accuracy using MASE					
Method	Model	MASE			
		Total	Control areas	Health boards	Bottom
Base	Naive	1.139	1.059	1.047	1.019
Base	ETS	0.963	0.930	0.899	1.038
Base	GLM	0.910	0.940	0.923	1.002
Base	TSGLM	0.911	0.939	0.924	1.005
Base	Ensemble	0.782	0.856	0.876	1.008
MinT	Naive	1.138	1.059	1.047	2.651
MinT	ETS	0.877	0.916	0.915	1.289
MinT	GLM	0.848	0.901	0.902	2.493
MinT	TSGLM	0.852	0.903	0.903	2.513
MinT	Ensemble	0.753	0.844	0.872	2.260
(b) Point forecast accuracy using MSSE					
Method	Model	MSSE			
		Total	Control areas	Health boards	Bottom
Base	Naive	1.169	1.056	1.062	1.031
Base	ETS	0.979	0.875	0.816	0.975
Base	GLM	0.813	0.897	0.875	1.009
Base	TSGLM	0.822	0.901	0.875	1.050
Base	Ensemble	0.599	0.729	0.774	0.993
MinT	Naive	1.168	1.057	1.062	2.095
MinT	ETS	0.785	0.852	0.845	0.994
MinT	GLM	0.720	0.827	0.837	1.803
MinT	TSGLM	0.722	0.833	0.839	1.851
MinT	Ensemble	0.560	0.706	0.765	1.557
(c) Probabilistic forecast accuracy using CRPS					
Method	Model	CRPS			
		Total	Control areas	Health boards	Bottom
Base	Naive	30.387	10.882	5.500	0.302
Base	ETS	14.309	6.074	3.476	0.244
Base	GLM	15.396	6.253	3.576	0.244
Base	TSGLM	15.316	6.227	3.575	0.245
Base	Ensemble	12.978	5.727	3.430	0.243
MinT	Naive	30.368	10.902	5.498	0.313
MinT	ETS	13.515	5.967	3.547	0.243
MinT	GLM	13.839	5.917	3.453	0.246
MinT	TSGLM	14.000	5.947	3.455	0.248
MinT	Ensemble	12.585	5.728	3.426	0.247

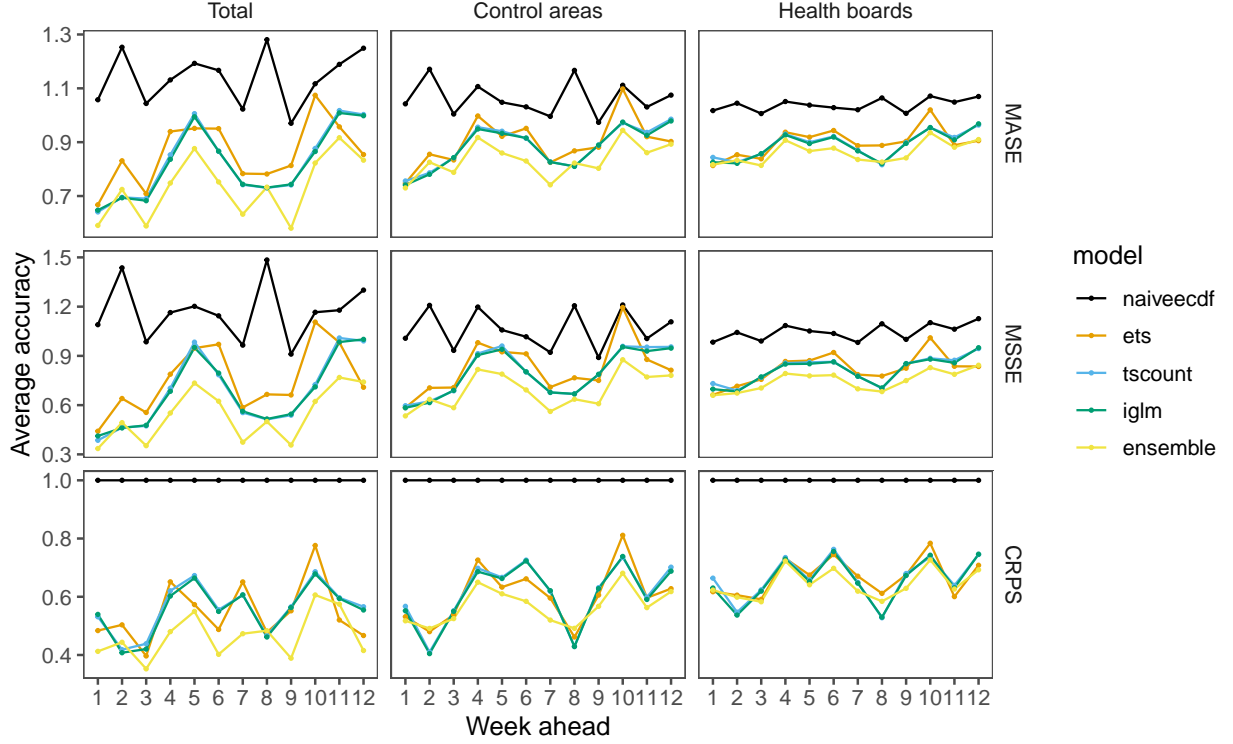


Figure 4: Average accuracy by week for 12 weeks using MinT reconciliation. CRPS is relative to a naive ECDF. MASE and MSSE are relative to the corresponding values for the training set.

6. Conclusion

Forecasting problems at Emergency Medical Services often have inherent hierarchical and grouped structures. For example, looking at time series of arrival calls in a clinical desk service, verified incidents, or attended incidents in a country, they could be disaggregated by various attributes of interest. Total demand in the country could be disaggregated by region, then within each region by health board, within each health board, by station/hospital, and so on down to the post code area. Alternative structures may arise when attributes of interest are crossed rather than nested. For example, the total demand could be disaggregated by priority (e.g., Red, Amber, Green) or by nature of incidents. It is also natural to have a mixed structures, for example the total demand could be disaggregated by priority and also by health board.

Despite the inherent hierarchical structure of the forecasting problem in EMS, the common practice is to produce point forecasts for each time series independently. This practice may lead to a lack of coordination and possibly undesirable and conflicting outcomes. Furthermore, due to the asymmetric impact of resource allocation in this area, quantifying forecast uncertainty through probabilistic forecasts is also of value as it enables planners to manage associated risks. In this paper, we investigate the application of hierarchical forecasting methods for producing probabilistic forecasts of daily incidents attended up to 84 days ahead, using different forecasting methods.

Our results indicate that forecast reconciliation in EMS can not only contribute to a more coordinated approach to decision making through producing coherent forecasts, but also it can result in forecast accuracy improvement. Our proposed forecasting models, combined with reconciliation approaches, outperform the empirical distribution benchmark. We show that a substantial forecast improvement can be achieved at

higher levels of aggregation by applying forecast reconciliation methods. When a point forecast is of interest at the bottom level of series, we observe that reconciliation may not improve the forecast accuracy, if the bottom series are noisy and lack systematic patterns. However, forecast reconciliation may result in more forecast results for bottom series, if we are interested in the tails of forecast distribution rather than just center measures like mean (i.e. point forecast). Producing consistent forecast are also crucial for informing planning activities, we also demonstrate that proposed models produce consistent forecasts across all forecast horizons. Therefore, we recommend that forecast reconciliation approaches to be adopted for routine use in EMS, whenever hierarchical and/or grouped time series data need to be forecasted. Moreover, we found that using an ensemble forecasting model, combining all models developed in this paper, instead of using each individually, works remarkably well for our mixed hierarchical & grouped structure.

Further research could investigate the practical benefits of probabilistic hierarchical forecasting in EMS. Linking forecasts with its utilities (response time, resource utilization, etc) can offer an opportunity to maximize benefits through more holistic planning approach. While, we generated count attend incident for the base forecast using Poisson regression models, however the reconciles forecast is not yet count. This could be also an avenue for further research.

References

- Aringhieri, R., Bruni, M.E., Khodaparasti, S., van Essen, J.T., 2017. Emergency medical services and beyond: Addressing new challenges through a wide literature review. *Computers & Operations Research* 78, 349–368.
- Armstrong, J.S., 2001. Evaluating forecasting methods, in: Armstrong, J.S. (Ed.), *Principles of forecasting: a handbook for researchers and practitioners*. Kluwer Academic Publishers. chapter 14, pp. 443–472.
- Boylan, J.E., Goodwin, P., Mohammadipour, M., Syntetos, A.A., 2015. Reproducibility in forecasting research. *International Journal of Forecasting* 31, 79–90.
- Ekström, A., Kurland, L., Farrokhnia, N., Castrén, M., Nordberg, M., 2015. Forecasting emergency department visits using internet data. *Annals of emergency medicine* 65, 436–442.
- Faraway, J.J., 2016. *Extending the linear model with R: generalized linear, mixed effects and nonparametric regression models*. 2nd edition ed., Chapman and Hall/CRC.
- Gneiting, T., Katzfuss, M., 2014. Probabilistic forecasting. *Annual Review of Statistics and Its Application* 1, 125–151.
- Grekousis, G., Liu, Y., 2019. Where will the next emergency event occur? predicting ambulance demand in emergency medical services using artificial intelligence. *Computers, Environment and Urban Systems* 76, 110–122.
- Gul, M., Celik, E., 2020. An exhaustive review and analysis on applications of statistical forecasting in hospital emergency departments. *Health Systems* 9, 263–284.
- Haugsbø Hermansen, A., Mengshoel, O.J., 2021. Forecasting ambulance demand using machine learning: A case study from oslo, norway, in: 2021 IEEE Symposium Series on Computational Intelligence (SSCI), pp. 01–10. doi:[10.1109/SSCI50451.2021.9659837](https://doi.org/10.1109/SSCI50451.2021.9659837).
- Hulshof, P.J., Kortbeek, N., Boucherie, R.J., Hans, E.W., Bakker, P.J., 2012. Taxonomic classification of planning decisions in health care: a structured review of the state of the art in or/ms. *Health systems* 1, 129–175.
- Hyndman, R.J., Ahmed, R.A., Athanasopoulos, G., Shang, H.L., 2011. Optimal combination forecasts for hierarchical time series. *Computational statistics & data analysis* 55, 2579–2589.
- Hyndman, R.J., Athanasopoulos, G., 2021. *Forecasting: principles and practice*. OTexts. URL: <https://otexts.com/fpp3>.
- Hyndman, R.J., Koehler, A.B., 2006. Another look at measures of forecast accuracy. *International Journal of Forecasting* 22, 679–688.
- Ibrahim, R., Ye, H., L'Ecuyer, P., Shen, H., 2016. Modeling and forecasting call center arrivals: A literature survey and a case study. *International Journal of Forecasting* 32, 865–874.
- Kang, Y., Hyndman, R.J., Smith-Miles, K., 2017. Visualising forecasting algorithm performance using time series instance spaces. *International Journal of Forecasting* 33, 345–358.
- Liboschik, T., Fokianos, K., Fried, R., 2017. tscount: An r package for analysis of count time series following generalized linear models. *Journal of Statistical Software* 82, 1–51. URL: <https://www.jstatsoft.org/index.php/jss/article/view/v082i05>, doi:[10.18637/jss.v082.i05](https://doi.org/10.18637/jss.v082.i05).
- O'Hara-Wild, M., Hyndman, R., Wang, E., Caceres, G., 2022. fable: Forecasting models for tidy time serie. URL: <https://fable.tidyverts.org/>. r package version 0.3.2.
- Panagiotelis, A., Gamakumara, P., Athanasopoulos, G., Hyndman, R.J., 2022. Probabilistic forecast reconciliation: Properties, evaluation and score optimisation. *European Journal of Operational Research*.
- Rostami-Tabar, B., Ziel, F., 2022. Anticipating special events in emergency department forecasting. *International Journal of Forecasting* 38, 1197–1213. URL: <https://www.sciencedirect.com/science/article/pii/S0169207020300017>, doi:<https://doi.org/10.1016/j.ijforecast.2020.01.001>.
- Sasaki, S., Comber, A.J., Suzuki, H., Brunson, C., 2010. Using genetic algorithms to optimise current and future health planning-the example of ambulance locations. *International journal of health geographics* 9, 1–10.
- Shi, M., Rostami-Tabar, B., Gartner, D., 2022. Forecasting for unplanned care services: A literature review. Working Paper.
- Stodden, V., Miguez, S., 2013. Best practices for computational science: Software infrastructure and environments for reproducible and extensible research. Available at SSRN 2322276.

- Vile, J.L., Gillard, J.W., Harper, P.R., Knight, V.A., 2012. Predicting ambulance demand using singular spectrum analysis. *Journal of the Operational Research Society* 63, 1556–1565.
- Vile, J.L., Gillard, J.W., Harper, P.R., Knight, V.A., 2016. Time-dependent stochastic methods for managing and scheduling Emergency Medical Services. *Operations Research for Health Care* 8, 42–52.
- Wang, X., Hyndman, R.J., Li, F., Kang, Y., 2022. Forecast combinations: an over 50-year review. Technical Report. URL: robjhyndman.com/publications/combinations/.
- Wargon, M., Guidet, B., Hoang, T., Hejblum, G., 2009. A systematic review of models for forecasting the number of emergency department visits. *Emergency Medicine Journal* 26, 395–399.
- Wickramasuriya, S.L., Athanasopoulos, G., Hyndman, R.J., 2019. Optimal forecast reconciliation for hierarchical and grouped time series through trace minimization. *Journal of the American Statistical Association* 114, 804–819. URL: <https://doi.org/10.1080/01621459.2018.1448825>, doi:10.1080/01621459.2018.1448825.
- Zhou, Z., 2016. Predicting ambulance demand: Challenges and methods. arXiv preprint arXiv:1606.05363 .
- Zhou, Z., Matteson, D.S., 2016. Predicting melbourne ambulance demand using kernel warping. *The Annals of Applied Statistics* 10, 1977–1996.