

Forecast reconciliation: A review

George Athanasopoulos

Monash University, VIC 3145, Australia

Email: george.athanasopoulos@monash.edu

Corresponding author

Rob J Hyndman

Monash University, VIC 3800, Australia

Email: rob.hyndman@monash.edu

Nikolaos Kourentzes

University of Skövde, 541 28 Skövde, Sweden

Email: nikolaos@kourentzes.com

Anastasios Panagiotelis

The University of Sydney, NSW 2006, Australia

Email: anastasios.panagiotelis@sydney.edu.au

24 August 2023

JEL classification: C10,C14,C53

Forecast reconciliation: A review

Abstract

Collections of time series that are formed via aggregation are prevalent in many fields. These are commonly referred to as hierarchical time series and may be constructed cross-sectionally across different variables, temporally by aggregating a single series at different frequencies, or may even be generalised beyond aggregation as time series that respect linear constraints. When forecasting such time series, a desirable condition is for forecasts to be coherent, that is to respect the constraints. The past decades have seen substantial growth in this field with the development of reconciliation methods that not only ensure coherent forecasts but can also improve forecast accuracy. This paper serves as both an encyclopaedic review of forecast reconciliation and an entry point for researchers and practitioners dealing with hierarchical time series. The scope of the article includes perspectives on forecast reconciliation from machine learning, Bayesian statistics and probabilistic forecasting as well as applications in economics, energy, tourism, retail demand and demography.

Keywords: Aggregation, Coherence, Cross-temporal, Hierarchical time series, Grouped time series, Temporal aggregation

1 Introduction

In time series forecasting, aggregation occurs in a variety of settings. For example, regional level tourism demand aggregates to national tourism demand; total revenue from the sale of individual stock keeping units aggregates to total revenue from all stock keeping units; the Gross Domestic Product of an economy is an aggregate of individual components; time series data measured at a quarterly frequency can be aggregated to data at annual frequency. While hierarchical time series will be defined more formally in [Section 2](#), the notion of hierarchical forecasting can be understood via the simple example where there is a time series X , a time series Y and a time series $Z = X + Y$, and we are interested in forecasts of X , Y and Z .

In practice, it is important to acknowledge that our variables X , Y and Z may be forecast in isolation from one another. This may occur when each forecast is obtained using a different time series model, or when forecasts are produced by separate organisational silos ([Cha2013](#)). In such cases, it will typically be the case that adding the forecast of X to the forecast of Y will

not be equal to the forecast of Z . Indeed, even where forecasts for X , Y and Z are produced jointly, it is not typically the case that forecasts aggregate in the correct fashion¹. This leads to two fundamental questions facing the forecaster of hierarchical time series:

Question 1: How best to adjust forecasts so that they agree with the known aggregation structure?

and

Question 2: Does adjusting forecasts in this manner lead to improvements in forecast accuracy?

These questions have motivated a growing and fruitful area of research, particularly over the past decade. The top panel of [Figure 1](#) shows the growth in Google Scholar items by the search terms “Hierarchical forecasting” and “Forecast Reconciliation” (the latter to be defined in [Section 3](#)). The bottom panel tracks the occurrence of the terms “hierarchy”² and “reconcil” in the book of abstracts of the International Symposium of Forecasters, the leading conference on forecasting. Both measures, while crude, pick up on the growing interest in the topic, especially in academic circles. Further, the reference list of this paper will attest to the multidisciplinary nature of the field with breakthroughs in hierarchical forecasting being published in top-tier journals in statistics, econometrics, operations research and machine learning.

The impact of methods for forecasting hierarchical time series has not been limited to academia, with industry also showing a strong interest. We are aware of many organizations using modern hierarchical forecast methods in practice, including Amazon, the International Monetary Fund, the Bank of New York Mellon, IBM, Huawei, H&M, and Volkswagen. Methods have been implemented in leading analytics software platforms such as SAP, SAS, ForecastPro and Fiddlehead technologies, not to mention numerous open source packages in R and Python (see [Section 7](#)). Among the broader forecasting community including academics and practitioners, hierarchical data have featured as part of the M5 competition ([MakElAl2020](#); [SeaBow2021](#)) and the Global Energy Forecasting competition ([Gefcom2017](#)).

The growth and impact of hierarchical forecasting make a review paper timely. Throughout our focus will be on forecasting, although where there are similarities between hierarchical forecasting methods and other literature, we will discuss seminal papers (for example in [Section 3.1](#)). Methods for dealing with hierarchical time series often involve first generating forecasts of all series in the hierarchy. Throughout this review we will not focus on the models and methods used to obtain these original forecasts. Also, while we will focus in [Section 6](#) on some specific

¹Some rare, and necessarily restrictive exceptions are discussed in [Section 3.9](#).

²We acknowledge that this search term may pick up instances of the word hierarchical related to hierarchical models rather than hierarchical time series, although the use of hierarchical models in a conference primarily on time series is rare.

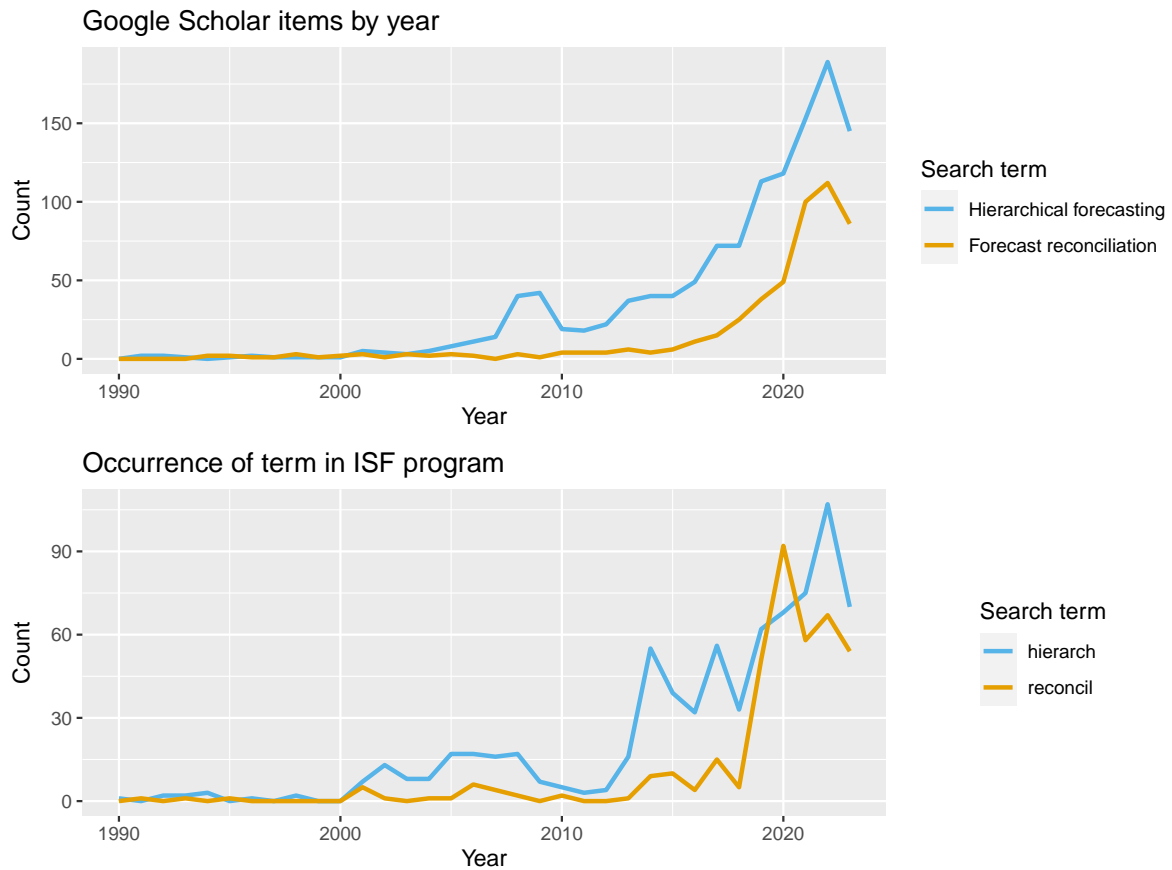


Figure 1: Search term results in Google Scholar and the book of abstracts for the International Symposium on Forecasting during ISFs (1990-2023). The Google Scholar search is as of August 20, 2023.

application areas where hierarchical forecasting methods have been used extensively, the methods have been applied so widely that not every applied paper can lie within the scope of this review and instead our focus is on papers that make important methodological contributions.

The remainder of the paper is organised as follows. [Section 2](#) provides the basic setting for hierarchical forecasting introducing notation and terminology and covering important historical background. [Section 3](#) covers forecast reconciliation which has been the hierarchical forecasting method to garner by far the most attention over the past decade. [Section 4](#) covers the special case of temporal aggregation of a single time series, which (despite a separate historical development) has since adopted methods from (and now merged with) cross-sectional aggregation. [Section 5](#) covers approaches to probabilistic forecasting, an area that, while previously ignored, has in recent years seen some important breakthroughs. Significant applications in tourism, macroeconomics, energy, demography, retail and healthcare are covered in [Section 6](#). Finally, in [Section 8](#), we look to the future of the field and point to some open questions in hierarchical forecasting.

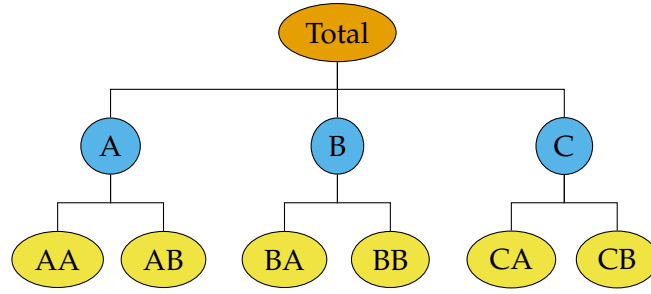


Figure 2: The diagram shows a 2-level hierarchical tree structure.

2 The setting

2.1 Hierarchical and grouped time series

We define a hierarchical time series as a multivariate time series, $\mathbf{y}_1, \dots, \mathbf{y}_T$, that adheres to some known linear constraints.³ For example, Figure 2 shows a 2-level hierarchical structure. Let $y_{\text{Tot},t}$ be the total (level 0) of all series at time t ; and let $y_{i,t}$ be the value of the time series at node i and time t . Let $\mathbf{y}_t \in \mathbb{R}^n$ be a vector comprising observations at time t of all time series in the hierarchy, and $\mathbf{b}_t \in \mathbb{R}^{n_b}$ ($n_b < n$) be a vector comprising the observations at time t of only the most disaggregated *bottom-level* series. The remaining $n_a = n - n_b$ aggregated series can be written as

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

for an appropriate $n_a \times n_b$ aggregation matrix \mathbf{A} , and the full set of time series can be written for all t as

$$\mathbf{y}_t = \mathbf{S}\mathbf{b}_t,$$

where $\mathbf{y}_t = \begin{bmatrix} \mathbf{a}_t \\ \mathbf{b}_t \end{bmatrix}$ and $\mathbf{S} = \begin{bmatrix} \mathbf{A} \\ \mathbf{I}_{n_b} \end{bmatrix}$ is the $n \times n_b$ *summing* or *structural* matrix.

For example, for the hierarchical structure of Figure 2, $n = 10$, $n_b = 6$, $n_a = 4$, $\mathbf{b}_t = [y_{AA,t}, y_{AB,t}, y_{BA,t}, y_{BB,t}, y_{CA,t}, y_{CB,t}]'$, $\mathbf{a}_t = [y_{\text{Tot},t}, y_{A,t}, y_{B,t}, y_{C,t}]'$ and

$$\mathbf{A} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix}.$$

The aggregation matrix, \mathbf{A} , describes how the bottom-level series aggregate to the series above. Hence, the columns of \mathbf{S} span the linear subspace of \mathbb{R}^n for which the linear constraints hold.

³This paper follows the recommendations of <https://robjhyndman.com/hyndsight/reconciliation-notation.html>.

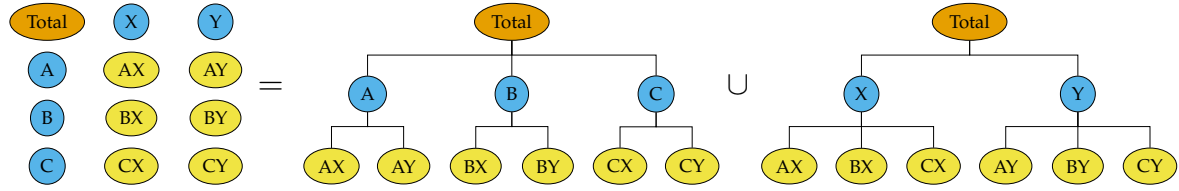


Figure 3: *The diagram shows that a 2-level grouped structure can be thought of as the union of two hierarchical tree structures with common top and bottom level series.*

We refer to this as the *coherent subspace* and denote it by \mathfrak{s} . We refer to the property that data adhere to these linear constraints as *coherence*.

Figure 3 shows a *grouped* structure in which the bottom-level series are aggregated by attributes of interest that are crossed, in contrast to the *hierarchical* (nested) structure shown in Figure 2. For this grouped example, the bottom-level series $\mathbf{b}_t = (y_{AX,t}, y_{AY,t}, y_{BX,t}, y_{BY,t}, y_{CX,t}, y_{CY,t})'$, aggregate into $y_{A,t}$, $y_{B,t}$ and $y_{C,t}$, and also into $y_{X,t}$ and $y_{Y,t}$. Hence, in contrast to hierarchical time series, grouped time series do not naturally aggregate (or disaggregate) in a unique manner. However, for simplicity, when we refer to hierarchical time series we mean both hierarchical and grouped structures. We will highlight the difference when it is important to do so.

2.2 Other representations

The *structural representation*, based on the summing matrix \mathbf{S} in the form shown above, is not the only way to write the constraints for the time series \mathbf{y}_t .

First, the ordering of the series within \mathbf{y}_t is arbitrary, and there is no requirement for the bottom-level series to appear below the aggregated series. An alternative order is sometimes more convenient, and then the rows of \mathbf{S} can be permuted to match the order of \mathbf{y}_t .

The coherent structure can also be expressed via a *constraint matrix* such that

$$\mathbf{C}\mathbf{y}_t = \mathbf{0}.$$

If we start with the structural representation shown above, then we can write $\mathbf{C} = [\mathbf{I}_{n_a} \quad -\mathbf{A}]$. It is often more convenient to work with this zero-constrained representation, rather than the structural representation. In fact, we can simply start with a general constraint matrix \mathbf{C} , that may not be of full rank, without defining an aggregation or summing matrix (Di_FonGir2022a).

Each of these representations has been used in the forecast reconciliation literature, and we will return to them in subsequent sections.

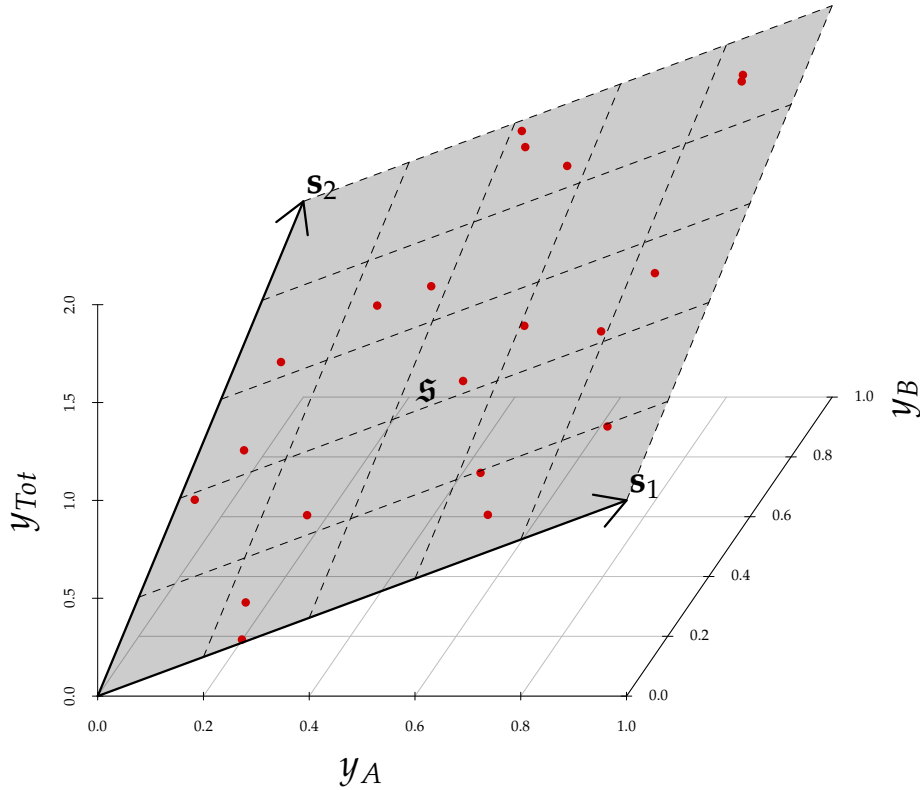


Figure 4: Depiction of a three dimensional hierarchy with $y_{Tot} = y_A + y_B$. The gray coloured two dimensional plane depicts the coherent subspace \mathfrak{s} where $\vec{s}_1 = (1, 1, 0)'$ and $\vec{s}_2 = (1, 0, 1)'$ are basis vectors that span \mathfrak{s} . The red points in \mathfrak{s} represent realisations or coherent forecasts.

There is no requirement for the S , A or C matrices to contain only 0s and ± 1 s. They can include any real values, specifying linear constraints that apply to the available time series (AthEtAl2020). For an example of this, see LiHyn2021 where the elements of the S matrix are proportions. Furthermore, this brings to light the full generality of so-called hierarchical time series. The methods discussed below can be applied to problems that need not be hierarchical and need not be time series (GirdiF2023).

2.3 Coherent forecasts

When forecasting hierarchical time series, we require the forecasts to adhere to the same linear constraints as the data; i.e., to aggregate in the same manner, or to follow the same linear constraints. We define a set of h -step-ahead point forecasts $\tilde{y}_{t+h|t} \in \mathbb{R}^n$ as *coherent* if $\tilde{y}_{t+h|t} \in \mathfrak{s}$. Figure 4 presents an example of the simplest possible hierarchy for which $y_t \in \mathbb{R}^3$, $b_t \in \mathbb{R}^2$ and $y_{Tot,t} = y_{A,t} + y_{B,t}$. The coherent subspace is shown as a grey 2-dimensional plane within a 3-dimensional space. Note that the columns of S , $\vec{s}_1 = (1, 1, 0)'$ and $\vec{s}_2 = (1, 0, 1)'$, span the coherent subspace; i.e., $\mathfrak{s} = \text{span}(S)$. The red points in \mathfrak{s} represent realisations or coherent forecasts.

pritularga2021stochastic note that this definition of coherence implicitly assumes that the measurement of the data occurs at a given level, typically the lowest. In practice, this may not be the case and due to measurement errors, different data collection methodologies, or otherwise, there may be discrepancies in the coherence. Therefore, they propose to add to the aggregation a statistical discrepancy term δ_t :

$$y_t = Sb_t + \delta_t.$$

Equivalently, this can be expressed as a slack term in the coherence constraints. When the data collection is done perfectly, naturally this term is zero. **AthEtAl2020** and **kourentzes2021visitor** provide examples where a time series is included in the hierarchy to deal with this discrepancy.

2.4 Single-level approaches

Traditionally, forecasts of hierarchical time series have involved selecting one level of aggregation, generating forecasts for that level, and then linearly combining these to obtain a set of coherent forecasts for the rest of the structure. These methods are usually classified as bottom-up, top-down or middle-out. Bottom-up approaches require generating forecasts at the bottom-level and then aggregating these up. Generating forecasts at the most disaggregate level implies that no information is lost due to aggregation. On the other hand, bottom-level data can be very noisy, or even intermittent, and hence more challenging to forecast. Top-down approaches require forecasts for only one time series at the most aggregate level and then disaggregating these down. Forecasting at the most aggregate level implies a large loss of information, and it can also be challenging to disaggregate these forecasts down. The disaggregation becomes even more challenging when the structure is grouped, as then the disaggregation paths are not unique. Further, **HynEtAl2011** and **PanEtAl2021** show that top-down approaches can introduce bias, even if the forecasts for the top level are unbiased. Middle-out approaches require forecasts at some intermediate-level and then aggregating these up and also disaggregating them down. In general, single-level approaches are limited to using information from a single-level and potentially ignoring valuable information from all other levels.

Another consideration comes from the number of forecasting models used in the hierarchy. In the top-down case, all predictions in the hierarchy are anchored to a single forecasting model at the top level (although forecasts at other levels may be used to compute the proportions for disaggregation). Similarly, in bottom-up, all forecasts are anchored to the bottom level. This introduces modelling and estimation risks, where the few forecasts that are used to populate the rest of the hierarchy may be of poor quality. For example, **kourentzes2017demand** show that

even with full knowledge of the data generating process, estimation errors can substantially reduce the quality of the resulting forecasts on other levels of the hierarchy.

In this section we concentrate on the implementation of single-level approaches for generating point forecasts. Related approaches for generating coherent probabilistic forecasts are discussed in [Section 5](#). Methods in the purely temporal setting are reviewed in [Section 4.1](#).

The vast majority of the literature, prior to the introduction of the concept of forecast reconciliation, almost exclusively focused on comparing bottom-up and top-down methods. **OrcEtAl1968** and **EdwOrc1969** are from the early works arguing that information loss is substantial and therefore it is important to work with the most disaggregate data available. **Kin1971** found that disaggregated earnings' data by market segments resulted in more accurate forecasts than when firm-level data were used. Building on this result, **Col1976** compared segmented econometric models with aggregate models for a group of 96 firms, and found that the segmented models produced more accurate forecasts for both sales and profit. **DunEtAl1976** show that forecasts aggregated from a lower level for modeling telephone demand are more accurate than the top-down method, although the comparison was based on only nine series. **ShlWol1979** concluded that the bottom-up method is preferable under some conditions on the structure of the hierarchy and the forecast horizon.

SchEtAl1988 looked at the bias and robustness of the two methods and concluded that the bottom-up method is better except when there are missing or unreliable data at the lowest levels. **DanMor1992** construct 15,000 artificial 2-level hierarchies using the M-competition data with two series at the bottom level. They found that bottom-up forecasts were more accurate, especially when the two bottom-level series were highly correlated. **ZelTob2000** used annual GDP growth rates from 18 countries and found that disaggregation provided better forecasts, results in line with earlier perspectives expressed by **Esp1994**. Another comparison is that of **wanke2007top** who compare the two approaches for safety inventory levels. **WanEtAl2013** analyse aggregate versus disaggregate forecasts for international arrivals into Hong Kong considering alternative bottom-up approaches, arguing that these take advantage of the heterogeneity across the disaggregate series, and show that the traditional bottom-up approach is more accurate compared to directly forecasting at the aggregate level.

There are fewer studies that find clear evidence and argue for a top-down approach, in contrast to a bottom-up approach. **GruGri1960** argue that disaggregated data are error prone and that top-down forecasts may therefore be more accurate. **StrEtAl2008** in an estimation setting show superiority of a top-down estimator. **HudEtAl2015** [find evidence of using a top down method](#)

when forecasting noisy geographic time series in many applications. **AthEtAl2009** propose two new top-down approaches based on forecast proportions rather than historical proportions which show promising performance. These approaches can lead to negative weights, which perhaps stretches the definition of disaggregation. **GroSoh1990** and **AthEtAl2009** provide a summary of top-down approaches.

FliMab1992 experiment with how different groupings (based on clustering) of time series, affect the forecast accuracy of traditional approaches. **Fli1999** argues that strong positive or negative correlation of sub-aggregate series, enhances forecast accuracy of the aggregate series whether a bottom-up or a direct approach is used for forecasting the aggregate and vice versa (low correlation in the bottom-level series diminishes forecast accuracy of the aggregate series). He further concludes that direct forecasts of an aggregate variable are more accurate than using a bottom-up approach. **Lut1984** and **IIm1990** show that it might be preferable to forecast the aggregate variable directly rather than using a bottom-up approach. **Hub2005** also concludes that using a bottom-up approach to forecast inflation for the Euro area. He attributes this to shocks affecting components of inflation in a similar way over the evaluation period and therefore forecast bias is increased when aggregating subcomponent forecasts. **KreEtAl2016** examine bottom-up versus direct forecasts for the aggregate through a behavioural lens and argue for advantages and disadvantages for a bottom-up judgemental approach which depend to a large degree on the underlying correlation structure at the bottom level.

ZotEtAl2005 and **ZotKal2007** argue that forecast accuracy is highly correlated to the choice of aggregation level which depends on the underlying data generation process. **WidEtAl2008** compare top-down to bottom-up, in a restricted simulation setting, and find that the difference in forecast accuracy between these is insignificant when the correlation between the sub-aggregate components is small or moderate. **SbrSil2013** extend these results and conclude that neither top-down nor the bottom-up approach should be preferred a priori in any empirical analysis. **WilWal2011** compare top-down to bottom-up demand forecasts. They conclude that the superiority of the methods depends on whether shared weekly point-of-sale data are used.

Fli2001 also reviewed these approaches, and discuss their advantages and disadvantages. He notes that different forecasting methods may be better suited to different aggregation levels, and this may affect the choice of which level to use for forecasting.

Kah1998 highlighted the need for a method that would enjoy and combine the good features of single-level approaches. **A simple solution of averaging top-down and bottom-up forecasts was**

suggested by Lap1998 and later built upon by RehEtAl2022. The call was taken up in full by HynEtAl2011 and AthEtAl2009, who introduced the concept of forecast reconciliation.

3 Forecast reconciliation

3.1 Least squares reconciliation outside of forecasting

The concept of least squares reconciliation has appeared in several contexts outside of the forecasting domain. As far back as the 1940s, reconciliation procedures were being used for national economic accounts. The national economic account is disaggregated into production, income and outlay, and capital transactions, which are further disaggregated by various factors. The aim is to coherently estimate the national account for all disaggregated and aggregated levels. StoEtAl1942 formulated the problem using simultaneous linear equations (similar to the zero-constrained form). Stone61 proposed a constrained estimation approach to balancing national accounts, where the constrained estimates are a weighted linear combination of initial estimates. This underpinned the work for which Richard Stone later won the 1984 Nobel Prize in Economics. Byron (Byron1978; Byron1979) formalized and extended Stone's work using more computationally efficient procedures. Suppose the national accounts are expressed as a vector y which need to satisfy the constraint $Cy = 0$, and let the original (incoherent) account estimates be denoted by \hat{y} . Then the reconciled estimates \tilde{y} are found by solving the constrained generalized least squares (GLS) problem

$$\tilde{y} = \arg \min_y (y - \hat{y})' W^{-1} (y - \hat{y}), \quad s.t. Cy = 0.$$

Assuming C is full rank, Byron (Byron1978; Byron1979) provide the solution $\tilde{y} = M\hat{y}$, where

$$M = I - WC'(CWC')^{-1}C$$

is a projection matrix, and W is a positive definite matrix. See also Wea1992, SmiEtAl1998 bikker2013benchmarking for a more modern treatment of this approach.

Later, the same idea was applied to reconciling other time series produced by national statistics offices. A review of some of this work is provided by Dagum_Cholette2006. To take just one example, seasonally adjusted time series require reconciliation. While the original time series data are coherent (e.g., national and state employment numbers), after each series is seasonally adjusted, they become incoherent. The same least squares solution is used for this problem (di2011simultaneous; corona2021optimal; QueFor2012).

Temporal reconciliation is also of interest to national statistics offices, ensuring monthly or quarterly estimates sum to the annual estimates (**chow1971best**). Simultaneous least squares reconciliation of time series estimates in both cross-sectional and temporal dimensions was introduced by **Di_Fonzo1990**, building on **rossi1982note**.

Least squares reconciliation has also found its way into chemical process measurement. Chemical process data are inherently noisy, and data reconciliation methods allow adjustment of measured values to satisfy specific material and energy constraints (**Romagnoli2000**).

In the engineering literature, a related problem involves optimal vehicle tracking where roads provide locally linear constraints on the position of a vehicle (**simon2002kalman**). This line of research is summarised in **simon2006optimal** and **simon2010kalman**.

3.2 First attempts at reconciliation in forecasting

To our knowledge, the earliest published work that applied least squares reconciliation in a forecasting context was the PhD thesis of Roman **Ahm2009**, working under the supervision of Rob Hyndman and George Athanasopoulos. The main methodological contributions from this thesis eventually appeared as **HynEtAl2011**. First, they showed that all of the existing bottom-up, middle-out and top-down methods could be expressed as

$$\tilde{y}_h = SG_h \hat{y}_h \quad (1)$$

for a suitably chosen $n_b \times n$ matrix G_h . (We will drop the subscript h when G does not depend on the forecast horizon, h .) Here, G_h maps the base forecasts \hat{y}_h into the bottom level, and so can be thought of as a forecast combination that combines all base forecasts to form bottom-level reconciled forecasts. In the special case of bottom-up forecasting, $G = [O_{n_b \times n_a} \quad I_{n_b}]$, while for top-down forecasts, the first column of G contains the proportions for each of the bottom-level series, while the remaining columns are all zero.

HynEtAl2011 showed that if the base forecasts \hat{y}_h are unbiased with covariance W_h , and $SG_h S = S$, then the reconciled forecasts \tilde{y}_h are also unbiased and have covariance $SG_h W_h G_h' S'$. Notably, the condition $SG_h S = S$ is generally not satisfied for top-down methods, with the exception of those discussed in [Section 3.6](#).

To find the optimal matrix G_h , **HynEtAl2011** formulated the problem as a regression of the form

$$\hat{y}_h = S\beta_h + \varepsilon_h$$

where ε_h is the reconciliation error with covariance V_h . This led to the GLS solution

$$G_h = (S'V_h^{-1}S)^{-1}S'V_h^{-1}. \quad (2)$$

The covariance V_h is unknown (**WicEtAl2019**), but under some conditions, **HynEtAl2011** showed that (2) collapses to an OLS solution where V_h is replaced by an identity matrix, giving $G = (S'S)^{-1}S'$.

Another key contribution of **HynEtAl2011** was to propose using sparse matrix algebra to greatly speed up the computations for large systems of time series. Simultaneously and independently, **di2011simultaneous** also proposed sparse matrix algebra for reconciling historical time series.

The first application of these new ideas was **AthEtAl2009**, which appeared two years earlier due to delays in publishing **HynEtAl2011**. There, the OLS reconciliation was compared to various top-down and bottom-up methods, using some quarterly Australian tourism data disaggregated by a geographic hierarchy and purpose of travel. Variations of these Australian tourism data have since become ubiquitous for benchmarking forecast reconciliation methods.

The `hts` R package (**Rhts1**) implementing the OLS reconciliation method appeared on CRAN in 2010, and led to the method quickly becoming popular in business and industry, long before the methodological paper actually appeared.

An early explanation of the method intended for practitioners appeared as **HynAth2014**, while the ideas made their way into an undergraduate textbook in **fpp2018** and **fpp2021**.

3.3 Scaled reconciliation methods

One obvious drawback of the OLS approach is that it weights all series equally, whether they are aggregates or disaggregates, and whether their base forecasts are good or bad. An early recognition of this issue is in **kourentzes2014improving** who treat the aggregate time series to bring all series on the same scale. The same issue prompted **HynEtAl2016** to propose a weighted least squares (WLS) solution, where the series are weighted by the inverse variances of the base forecasts, later referred to as “variance scaling”. If W_h is the covariance matrix $\text{Var}(\mathbf{y}_{T+h|h} - \hat{\mathbf{y}}_h)$, then the WLS solution is $G_h = (S'\Lambda_h^{-1}S)^{-1}S'\Lambda_h^{-1}$ and $\Lambda = \text{diag}(W_h)$. **HynEtAl2016** was also the first forecasting paper to note that the methods applied to grouped time series as well as to strictly hierarchical structures.

3.4 Minimum trace reconciliation

WicEtAl2019 provided theoretical insights into the problem by taking an optimization approach rather than a regression approach. They formulated the problem as minimizing the trace (MinT) of the covariance matrix $\text{Var}(\mathbf{y}_{T+h|h} - \tilde{\mathbf{y}}_h)$, equal to the sum of the variances of all the reconciled forecasts, and showed that the solution is given by

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

where \mathbf{W}_h is the covariance matrix $\text{Var}(\mathbf{y}_{T+h|h} - \hat{\mathbf{y}}_h)$. Equivalently, $\tilde{\mathbf{y}}_h = \mathbf{M}_h\hat{\mathbf{y}}_h$, where

$$\mathbf{M}_h = \mathbf{S}(\mathbf{S}'\mathbf{W}_h^{-1}\mathbf{S})^{-1}\mathbf{S}'\mathbf{W}_h^{-1}.$$

This MinT solution is equivalent to GLS, and has the WLS and OLS solutions as special cases. **WicEtAl2019** also showed that there is an equivalent and more computationally efficient solution given by

$$\mathbf{M}_h = \mathbf{I}_n - \mathbf{W}_h\mathbf{C}'(\mathbf{C}\mathbf{W}_h\mathbf{C}')^{-1}\mathbf{C}$$

where $\mathbf{C} = [\mathbf{I}_{n_a} \quad -\mathbf{A}]$, which matches the earlier work of Byron (**Byron1978; Byron1979**) for reconciling national accounts (although derived from a different perspective). **This inverts an $n_a \times n_a$ matrix rather than an $n \times n$ matrix.**

A difficulty with the MinT solution is in estimating the covariance matrix, \mathbf{W}_h , especially for $h > 1$. The sample covariance matrix of the base models' residuals provides an estimate of \mathbf{W}_1 , but it is often a very poor estimate, and may be singular, especially when $n > T$ (which is true for many real hierarchical time series). **WicEtAl2019** proposed a shrinkage estimator of \mathbf{W}_1 , where the off-diagonal elements are shrunk towards zero, and suggested approximating \mathbf{W}_h as a scalar multiple of \mathbf{W}_1 . The scalar multiple cancels when computing \mathbf{G}_h , and so it does not need to be estimated when computing point forecasts.

WicEtAl2019 also discuss a simple alternative approach to finding \mathbf{W}_h , first proposed by **AthEtAl2017**, based only on the structure of the hierarchy, and not on a statistical estimate. In this “structural scaling” approximation, $\mathbf{W}_h = \mathbf{\Lambda} \propto \text{diag}(\mathbf{S}\mathbf{1}_{n_b})$, where $\mathbf{1}_{n_b}$ is an n_b -vector of 1s. That is, \mathbf{W}_h is a diagonal matrix with entries proportional to the row sums of \mathbf{S} . This is the covariance matrix that would arise if all the most disaggregated series were uncorrelated with each other and they had the same forecast variance.

Observe that (1) implies a combination of forecasts in $G_h \hat{y}_h$. Motivated by this, [Pritularga2021stochastic](#) investigated the implications of estimation uncertainty in forecast reconciliation. They showed that uncertainties in both the forecasting models generating \hat{y}_h and in G_h will influence the quality of the reconciled forecasts. As S is fixed, the approximation of W_h carries the reconciliation uncertainty in G_h , where each element that needs to be estimated can potentially increase the forecast error of the reconciled forecasts. This can help explain the surprising performance of “structural scaling” ([AthEtAl2017](#)), where all elements of G are fixed, and the volatile performance of MinT, especially for short time series, where estimation errors in W_1 can dominate. **This bears some similarities with the forecast combination puzzle ([ClaEtAl2016](#)) where the performance of forecast combinations is hindered by imprecise estimation of covariances.** By introducing various approximations of increasing complexity of W_1 , [Pritularga2021stochastic](#) demonstrate that this choice can have significant effect on the quality of the reconciled forecasts, especially when it comes to having consistent performance over different forecast origins, a key element of trustworthiness in forecasting ([spavound2022making](#)). From the investigated approximations, retaining only the block-diagonal structure of the shrinkage estimator of W_1 was found to perform well in a variety of situations.

3.5 Other optimization approaches

[Van_ErvCug2015](#) took a game-theoretic approach to forecast reconciliation, and chose to find the solution to the minimax problem

$$V = \min_{\tilde{y} \in \mathcal{S}} \max_{y \in \mathcal{S}} \{ \ell(y, \tilde{y}) - \ell(y, \hat{y}) \},$$

where ℓ is a loss function, and \mathcal{S} is the coherent subspace. They demonstrate that $V \leq 0$, so that the reconciled forecasts are guaranteed to have smaller loss than the base forecasts. They further show that, when ℓ is L_2 loss, the minimax solution is equivalent to solving the constrained least squares problem, where the reconciled and base forecasts are as close as possible subject to the reconciled forecasts being coherent, leading to the closed form solution of (2) for G_h .

[PanEtAl2021](#) unify, and in certain cases generalise, the results of [Van_ErvCug2015](#) and [WicEtAl2019](#), providing a geometric intuition. **In particular, they consider a loss function of the form $\ell(y, \tilde{y}) = (y - \tilde{y})' \Psi (y - \tilde{y})$, where Ψ can be any symmetric positive definite matrix, \tilde{y} is either \hat{y} or \tilde{y} and derive two main results. The first is that the reconciled forecast $\tilde{y} = S(S' \Psi S)^{-1} S' \Psi \hat{y}$ will always improve upon the base forecast in the sense that $\ell(y, \tilde{y}) \leq \ell(y, \hat{y})$, where the strict inequality holds if \hat{y} is incoherent. This generalises the result of ([Van_ErvCug2015](#)) to non-diagonal Ψ**

The second is that the MinT solution $\tilde{\mathbf{y}} = \mathbf{S}(\mathbf{S}'\mathbf{W}_h^{-1}\mathbf{S})^{-1}\mathbf{S}'\mathbf{W}_h^{-1}\hat{\mathbf{y}}$ will optimise loss in expectation for any choice of Ψ (WicEtAl2019). Note that the second result does not consider the estimation uncertainty in G_h (pritularga2021stochastic), which can be especially prominent for short time series, explaining cases in the literature where this result seems to be violated.

If we are willing to drop the unbiased condition, and allow both base and reconciled forecasts to be biased, a different least squares solution emerges, as shown by Ben_TaiEtAl2019. They use an expanding window approach, applying the base forecasting method iteratively to the training data, computing $\hat{\mathbf{y}}_{t+h|t}$, the h -step-ahead base forecast of \mathbf{y}_{t+h} based on training data $\mathbf{y}_1, \dots, \mathbf{y}_t$, for $t = T_1, \dots, T - h$. They consider the regularized empirical risk minimization problem

$$\min_{\mathbf{G}} L_T(\mathbf{G}),$$

where

$$L_T(\mathbf{G}) = \frac{1}{Nn} \|\mathbf{Y} - \hat{\mathbf{Y}}\mathbf{G}'\mathbf{S}'\|_F + \lambda \|\text{vec}\mathbf{G}\|_1,$$

$N = T - T_1 - h + 1$, $\|\cdot\|_F$ is the Frobenius norm, $\mathbf{Y} = [\mathbf{y}_{T_1+h}, \dots, \mathbf{y}_T]'$, $\hat{\mathbf{Y}} = [\hat{\mathbf{y}}_{T_1+h|T_1}, \dots, \hat{\mathbf{y}}_{T|T-h}]'$, and λ is a regularization parameter. The first term contains the errors of the reconciled forecasts, while the second shrinks the elements of \mathbf{G} to zero, providing some regularization of the amount of reconciliation involved. When $\lambda = 0$, they show that its solution is

$$\hat{\mathbf{G}} = \mathbf{B}'\hat{\mathbf{Y}}(\hat{\mathbf{Y}}'\hat{\mathbf{Y}})^{-1}.$$

where $\mathbf{B} = [\mathbf{b}_{T_1+h}, \dots, \mathbf{b}_T]'$. When $\hat{\mathbf{Y}}'\hat{\mathbf{Y}}$ is non-invertible, the solution is not unique, and a generalized inverse can be used.

Inspired by this development, Wic2021 proposed minimizing the trace of the forecast error covariance matrix without an unbiased constraint, to create an unconstrained version of MinT which she called “MinT-U”. She also derived an estimate of the resulting \mathbf{G} matrix in the case where the series are jointly weakly stationary, dubbing the resulting method “EMinT-U” (empirical MinT unconstrained).

3.6 Adding optimization constraints

Any approach to reconciliation based on optimisation uses a form of constrained optimisation since reconciled forecasts must lie on the coherent subspace. However, at times additional constraints may be implemented. The first is the case where reconciled forecasts must be non-negative. In general, even if base forecasts are constrained to be positive (which can be

achieved by modelling on the log scale and back-transforming), there is no guarantee that the usual reconciliation approaches such as OLS and MinT will maintain the non-negativity of forecasts. To address this issue, the usual optimisation problem can be augmented with non-negativity constraints on the reconciled forecasts. Such optimisation problems can be solved using quadratic programming, with **WicEtAl2020** providing an early example for forecast reconciliation, and **di2023spatio** a more recent example.

KouAth2021 also consider the case of non-negative reconciled forecasts. However, instead of using a constrained optimisation approach, they propose a heuristic to iteratively adjust the reconciled predictions to be non-negative. Although this does not guarantee optimal solutions, their proposed algorithm has the interesting feature that it distributes adjustments of forecasts across the hierarchy, which can be useful in a variety of situations, such as the application of judgemental adjustments on specific nodes of the hierarchy.

di2023spatio also discuss an effective nonnegative heuristic called “set-negative-to-zero”, whereby the negative reconciled forecasts at the bottom level are set to zero, and the remaining forecasts computed via aggregation.

Another constraint of interest is where some particular base forecasts remain unchanged. For instance, **HolEtAl2021** consider the case of reconciliation where the top-level base forecast is retained. This differs from truly top-down approaches in that it can be done while also preserving the unbiasedness of base forecasts. To briefly illustrate the main idea, for a three variable hierarchy where $y_{Tot,t} = y_{A,t} + y_{B,t}$, either setting

$$\begin{pmatrix} \tilde{y}_{Tot,t} \\ \tilde{y}_{A,t} \\ \tilde{y}_{B,t} \end{pmatrix} = \begin{pmatrix} \hat{y}_{Tot,t} \\ \hat{y}_{A,t} \\ \hat{y}_{Tot,t} - \hat{y}_{A,t} \end{pmatrix} \quad \text{or} \quad \begin{pmatrix} \tilde{y}_{Tot,t} \\ \tilde{y}_{A,t} \\ \tilde{y}_{B,t} \end{pmatrix} = \begin{pmatrix} \hat{y}_{Tot,t} \\ \hat{y}_{Tot,t} - \hat{y}_{B,t} \\ \hat{y}_{B,t} \end{pmatrix}$$

will lead to coherent forecasts that preserve unbiasedness. Any average between these two solutions will have the same properties. **HolEtAl2021** generalise this idea to more complex hierarchies, and the properties of their methods are investigated by **Di_FonGir2022b**. **ZhaEtAl2023** further generalise this idea to a setting where reconciliation can be carried out while keeping a subset of base forecasts unchanged and not just the top level. Conditions on how a set of such “immutable” series can be selected are also provided by **ZhaEtAl2023**.

3.7 Machine learning and regularization

Machine Learning (ML), including Artificial Intelligence (AI), methods have been used to provide various modifications of the optimal combination approach of **HynEtAl2011**. Most of

the contributions, attempt to replace the linear regression formulation with a less restrictive method to obtain combinations of forecasts from the various hierarchical levels. Coherence is achieved via a bottom-up approach, or by embedding coherence in the ML training. A minority of contributions focus on other aspects of hierarchical forecasting, such as selecting the best reconciliation method.

QiaHua2018 focus on earnings forecasting, and rely on the problem structure of hierarchical forecasting to address this, recognising the forecast combination at its core. However, instead of combining forecasts across hierarchical levels, they combine forecasts across alternative hierarchical mappings, and then proceed to achieve coherence using a bottom-up approach. The different mappings are the product of the different ways one can do the accounting of the different components that contribute to the net earnings. This raises the question of which hierarchy is best to use, but also how to efficiently search across the different hierarchies. They resolve the construction of the hierarchy using a genetic algorithm, to avoid the computationally infeasible greedy search across mappings of the hierarchy. Although they do not discuss this in their work, this approach could be used to relax the conventionally rigid hierarchies, and identifying re-mappings that can potentially improve the final result. Forecasts are generated by LSTM networks, for each time series and each different mappings of the hierarchy. The forecasts are then combined across these mappings to give final prediction, with encompassed forecasts being rejected from the combination.

SpiEtAl2021 rely on random forests and gradient boosting machines, specifically XGBoost, to facilitate the combination of forecasts implied by hierarchical forecasting. They show superior performance to the linear approach, however it is unclear whether the gains are due to the nonlinear capabilities of the ML methods, or due to the differential combination over various forecasts horizons that are considered and are typically omitted by the linear counterparts. Furthermore, it should be noted that the objective of the training of the ML methods is obtaining the minimum forecast combination errors, rather than minimum reconciliation errors. Coherent forecasts across the complete hierarchy are obtained via bottom-up aggregation. **BurChe2021** propose an alternative use of ML to achieve coherent forecasts. They recast the reconciliation step as an encoder-decoder setup, where base forecasts are processed by a trainable encoder to produce the reconciled bottom-level forecasts. These are then decoded using the summing matrix, as with a standard bottom-up setup. The encoder is implemented using a shallow feed-forward neural network. They find that this approach demonstrates increasing gains for deeper hierarchies.

Gle2020 attempts to overcome the lack of focus on coherence by adjusting the objective function. Using neural network forecasts, he includes a regularisation term that penalises incoherences in the generated forecasts. This follows from **mishchenko2019self** who proposed a similar regularisation term to obtain reconciled forecasts directly from the base forecasts. The disadvantage of these regularisation approaches is that they result in soft constraints that do not guarantee coherence. **Gle2020** provides two alternatives for the regularisation term and shows that the resulting forecasts can outperform standard MinT reconciliation. However, when the regularisation is used in conjunction with MinT this results in both coherent and the most accurate forecasts. **HanEtAl2021** propose a similar approach, where a regularisation term is added in the loss function, again based on coherence constraints. They also consider a regularised loss for producing coherent quantile forecasts. The authors demonstrate the use of the proposed regularised loss on a variety of linear and ML models, and also empirically show the negative effect of regularisation on coherence.

ShiEtAl2020 introduces a regularisation term, based on the coherence constraints, in the objective function to push bottom-level forecasts to fit both on their target series and their aggregate counterparts. They forecast the bottom-level series but as the regularisation cannot ensure coherence, these are used in a bottom-up setting to produce coherent forecasts for the rest of the hierarchy. The authors demonstrate the efficacy of this in the context of neural networks. They find that these outperform conventionally trained networks whose forecasts are then reconciled, either with bottom-up or MinT.

ParEtAl2021 propose a regularised neural network with sequence-to-sequence architecture. Focusing on the hierarchical part of the contribution, a regularisation term is added to incorporate the coherence constraints. As with the previous work, this does not guarantee coherence, yet forces the final forecasts to be approximately coherent. The regularisation is embedded in the loss function of the network, achieving an integrated approach. In contrast to the previous work, the network outputs forecasts for all the levels.

The contribution by **AndLi2021** can be seen to belong loosely in the regularisation approaches. Focusing on the M5 competition dataset, the authors produce separate forecasts for the top- and the bottom-levels of the hierarchy, using different methods (NBEATS and LightGBM respectively). To achieve coherent forecasts they modify the loss function of the bottom-level method, where positive errors can be multiplied by a factor. This factor is identified by minimising the incoherence error between the summed bottom-level forecasts and the top-level

forecast. Note that this factor is calibrated by keeping the top-level forecasts fixed, and therefore any coherence is obtained by modifying the bottom-level forecasts.

RanEtAl2021 propose another way to achieve an integrated forecast-reconciliation mechanism. First, a global neural network is used to forecast the time series in the hierarchy, which with the assumption of a forecast distribution can produce sample forecasts. The sample forecasts are subsequently reconciled, obtaining distributions of coherent probabilistic forecasts. The calculation of the loss for the training of the model makes use of the reconciled forecasts, allowing an end-to-end approach that parametrises the model to achieve both accurate and coherent forecasts. Note that the methodology allows for various assumptions on the forecast distribution, and the relaxing of these. The resulting forecasts guarantee coherence, in contrast to the previous integrated approaches. The substantive difference with conventional hierarchical approaches, that post-processes base forecasts, is that the global network can model richer interconnections between the time series in the hierarchy for the generation of the forecasts. Furthermore, the authors compare the results of the proposed hierarchical model against a global learner without coherent forecasts and demonstrate on average better performance, suggesting that the integrated methodology offers gains beyond any achieved by the global learning.

WanEtAl2022 contribute with a similar formulation. The important differences are in that they use an autoregressive transformer to produce the base forecasts and that their approach does not require assuming particular predictive distributions. Instead, they rely on empirical estimation (conditional normalising flow) for the distributions. Furthermore, their approach focuses on obtaining bottom-level forecasts which are internally aggregated in a bottom-up fashion to the complete hierarchy. Similar to the **RanEtAl2021**, the errors of the final forecasts are used during training, realising an end-to-end hierarchical forecasting method. The authors demonstrate gains in performance over a non-hierarchical version of their method, as well as over various benchmarks.

Focusing on temporal hierarchies, **TheKou2021** provide an end-to-end neural network based method. Similarly to **BurChe2021** they explore a series of encoder-decoders to achieve reconciliation, considering fixed and trained decoder weights, using the complete or only the bottom-level of the hierarchy. These can be used instead of conventional hierarchical methods, demonstrating good performance in a global learning setting. To achieve an end-to-end integration, they pass the temporal hierarchy data through a convolutional layer and subsequently to an LSTM. The convolutional layer compresses the abstracts the hierarchical time series, and the LSTM models the dynamics over time. By appending an encoder-decoder to the outputs of the LSTM an

end-to-end methodology is obtained. The authors demonstrate the gains due to the various components, but also investigate the amount of time series that are needed to achieve good performance in a global training setting. By modifying the training loss function, **TheKou2020** extend the method to provide quantile forecasts.

AboEtAl2022 use ML to address the challenge of selecting the best method to perform the hierarchical reconciliation. To achieve this they rely on a meta-learning classifier, which is trained to identify given the time series features of a dataset what is the best reconciliation method to employ. The approach is quite flexible in terms of features and classifiers, as well as the forecasting models and reconciliation methods.

AboEtAl2019 use ML to improve the performance of non-combination hierarchical approaches. They consider various ML methods to obtain better proportions to decompose upper-level forecasts to lower-levels. **FenZha2020** provide a standard implementation of hierarchical methods using base forecasts from ML models, and find the hierarchically reconciled forecasts to be the most accurate. **PunEtAl2020** leverage on LSTM networks to produce forecasts for the cross-temporal case. Forecasts are generated for the various temporal aggregation levels, which are reconciled first by using temporal hierarchies and then cross-sectional. **SprEtAl2021** propose a method for probabilistic gradient boosting machines, and benchmark its performance on the hierarchical M5 dataset, demonstrating good performance against other non-hierarchical ML methods. Although the proposed method can provide probabilistic predictions for all time series of the hierarchy, coherence is not established. **ManEtAl2021** use a deep neural network to directly produce reconciled forecasts. The neural network captures the structure of the hierarchy, as well links the relationship between time series features extracted at any level of the hierarchy and explanatory variables into an end-to-end neural network. **SagEtAl2021** **propose a deep long short-term memory approach developed for the hierarchical time-series setting.**

3.8 Bayesian versions

Papers tackling hierarchical forecasting from a Bayesian angle focus on different aspects of the problem, but have much in common due to certain advantages of Bayesian inference. These are: the suitability of Bayesian inference for models with latent states such as state space models; the natural way uncertainty is propagated via Bayes' rule leading to probabilistic forecasts; and the use of priors to incorporate judgement into forecasts.

ParNas2014 propose a top-down approach, similar to **AthEtAl2009**, that forecasts the bottom-level series as proportions of the top-level forecast, rather than forecasting them directly. To this end, a state space model is proposed where latent states are mapped to proportions via the

softmax function. A variational approximation factorised into states and remaining parameters is employed with Evidence Lower Bound (ELBO) optimised via the EM algorithm. For the data considered, the method improves upon the top-down method of **AthEtAl2009**, and gives more accurate bottom-level forecasts than bottom-up and OLS. However, as a top-down method, there is no scope for top-level forecasts to be improved using bottom-level series, and for the data considered in the paper, both bottom-up and OLS generate more accurate forecasts for the top-level series.

RoqEtAl2021 also employ Bayesian modelling for forecasting hierarchical data in a way that differs from the two-step reconciliation approach. In particular, they decompose time series into a trend, modelled by a piecewise linear component, and a stationary component modelled as a sum of Gaussian Processes (GPs). Rather than modelling individual GPs for each series in the hierarchy, these are fit group-wise, with groups determined according to the hierarchical structure. The method outperforms MinT in an empirical study for the Australian prison population data, but not for the Australian tourism data.

Another strain of the literature brings a Bayesian approach to the regression model interpretation of forecast reconciliation. **novetal2017** recognise that the posterior of β_h can act as a probabilistic forecast for the bottom-level series. Using Markov chain Monte Carlo to obtain a sample from this posterior, and then aggregating, gives a probabilistic forecast for the entire hierarchy. **EckEtAl2021** also obtain a posterior on β_h , but their focus is on augmenting the reconciliation regression equation with a vector of intercepts that allow for base forecasts to be biased and evolve according to a state space representation. Both **novetal2017** and **EckEtAl2021** suggest that in a Bayesian setting, judgement can be incorporated via the prior, in the latter case via an explicit empirical example where prior information about a structural break in data classification can be exploited. Also, while both papers recognise the potential of Bayesian inference to obtain probabilistic forecasts, neither paper makes this the focus of empirical evaluation. **novetal2017** minimise loss functions over the posterior sample and then use this as a point forecast, while **EckEtAl2021** use maximum a posteriori (MAP) estimates as point forecasts.

Probabilistic forecast reconciliation is the motivation and focus of the Bayesian algorithm proposed by **CorEtAl2021**. In particular, a prior is placed on the bottom-level series with the mean set to point forecasts obtained in the first step of forecast reconciliation and a variance given by the variance-covariance matrix of one-step ahead errors. This prior is updated using the top-level forecasts obtained in the first stage of forecast reconciliation via Bayes' rule. The method generalises MinT in the sense that the posterior mean is equivalent to the usual MinT

approach. The necessary updates via Bayes' rule have parallels with the Kalman filter since the reconciliation problem is recast as a linear Gaussian model. The empirical results are evaluated using scoring rules for probabilistic forecasts including the CRPS and energy score. This approach has also recently been extending to the challenging case of forecast reconciliation for discrete data by **CorEtAl2022** and **Zambon2022**.

Bayesian methods are likely to continue to play an important role in the development of the forecast reconciliation literature. Some promising avenues will be incorporating information from the hierarchical structure via the prior and to use Bayesian methods to obtain non-Gaussian probabilistic forecasts. The challenges are likely to be computational in nature, as scalability of MCMC methods to large hierarchies may be difficult. The development of fast alternatives such as variational inference represent a promising way forward.

3.9 In-built coherence

So far, all the approaches discussed have involved two steps — first compute the base forecasts \hat{y}_h , and then reconcile them to produce \tilde{y}_h . The computationally slow part is producing the base forecasts, because they usually involve fitting models to each series individually, with the estimation requiring non-linear optimization. However, as shown by **AshEtAl2021**, if the base forecasts \hat{y} are produced using a linear regression model, the base forecasts and reconciliation can be combined, giving coherent forecasts directly in a single closed form equation. Further, the computation is extremely fast provided sparse matrix algebra is used.

Another approach which aims to produce coherent forecasts directly is due to **PenvanDal2017**, who propose the state space model

$$\begin{aligned} y_t &= S\mu_t + \varepsilon_t, & \varepsilon_t &\sim N(0, \Sigma_\varepsilon), \\ \mu_t &= \mu_{t-1} + \eta_t, & \eta_t &\sim N(0, \Sigma_\eta). \end{aligned} \tag{3}$$

Variations are also considered, including covariates in the measurement equation (3). Coherent forecasts arise naturally using the Kalman filter, as discussed by **simon2010kalman**. However, the covariance matrices are difficult to estimate with anything other than small hierarchies.

A related state space approach was proposed by **villegas2018supply**, who show that their formulation subsumes bottom-up, top-down, and some forms of forecast reconciliation and combination forecasting.

Inbuilt constraints to form coherent forecasts have also been considered in an ARIMA modelling context (Cho1982; guerrero1989optimal; de1993constrained) and for exponential smoothing forecasts (rosas1994restricted).

4 Temporal and cross-temporal reconciliation

4.1 Early temporal aggregation papers

Studying the effects of temporal aggregation on forecasts goes back to the seminal works of AmeWu1972, Tia1972, and Bre1973. Wei1979, Lut1984_JoE, Lut1986, StrWei1986, HotCar1993 and RosSea1995 study the effect of temporal aggregation on seasonal and non-seasonal ARIMA processes respectively, with aligned theoretical results. In general, these show that aggregation to the annual frequency simplifies dynamics of ARIMA processes generated at monthly or quarterly frequencies. They state that “quarterly data may be the best compromise among frequency of observation, measurement error, and temporal aggregation distortion”. Such observations are not unusual. A similar conclusion is reached by NijPal1990. SilVer2008 provide a detailed review of the literature to that point, with a focus on the implications on the model structure and identifiability for univariate ARIMA and multivariate GARCH processes. PinEtAl1987 study the temporal and contemporaneous aggregation of scale and vector ARIMA processes, and provide some general results and forecast comparisons.

In parallel to the investigations in the econometric literature up to the early 2000s, temporal aggregation was becoming popular in high-frequency time series forecasting, albeit implicitly. The relatively limited computational resources of that period had forced researchers to develop clever ways to handle the long time series appearing in applications such as daily electricity load forecasting. The dominant approach had become to split the daily time series into seven weekly series, each corresponding to a stock measurement at a specific week of the day (hippert2001neural).

From an accuracy standpoint, in the last decade there has been a revival of temporal aggregation in the forecasting literature. LunBal2011 study 111 weekly time series related to cash withdrawals. They examine forecasts generated by two top-down approaches and find considerable improvements compared to forecasts generated by the daily models directly. Motivated by intermittent demand forecasting, nikolopoulos2011aggregate proposed the ADIDA method, where a time series is temporally aggregated to a less intermittent level, forecasted, and subsequently disaggregated. The authors find promising accuracy gains, although there is only heuristic guidance for the level of temporal aggregation. spithourakis2011improving demonstrate the

benefits of using ADIDA in fast moving consumer goods, while **spithourakis2014systemic** attempt to develop the theoretical background for the method, investigating the aggregation and disaggregation mechanism, and trying to identify how to select well-performing temporal aggregation levels. **rostami2013demand** and **rostami2014note** derive the optimal temporal aggregation level for AR(1), MA(1), and ARMA(1,1), when simple exponential smoothing is used to produce the forecasts.

Notably the discussion so far has been on non-overlapping temporal aggregation. This form of aggregation acts as a moving average filter (**kourentzes2014improving**), while substantially reducing the available sample size. **boylan2016performance** investigate the effect of overlapping temporal aggregation, where a moving window is used to aggregate the time series, which is moved iteratively over the original series, including and dropping one observation at the time. They provide the conditions for which overlapping temporal aggregation outperforms its non-overlapping counterpart for independently and identically distributed demand processes. Earlier work on overlapping aggregation by **hotta1992effect** explore the effects on ARIMA, and by **mohammadipour2012forecast** on INARMA models. Similarly, **petropoulos2016another** motivated by intermittent demand problems, investigate empirically the usefulness of aggregating over unequal time periods, finding cases that can be beneficial. **BabEtAl2021** provide an extensive review of the aggregation literature.

The use of temporal aggregation in demand forecasting is natural, as we are typically interested in the demand of the lead time period, which lends itself to temporal aggregation, and remains an active research area (**rostami2019impact**; **saoud2022approximations**).

4.2 Temporal reconciliation

Motivated by the resurgence in the interest in temporal aggregation, following the work by **nikolopoulos2011aggregate**, **kourentzes2014improving** propose the Multiple Aggregation Prediction Algorithm (MAPA), where a time series is modelled independently at multiple temporal aggregation levels with a state-space model, such as exponential smoothing, and then these models are combined by state. Temporal aggregation filters high frequency components in the time series, making low frequency ones more prominent. MAPA takes advantage of this to provide forecasts that perform well both in the short- and long-term forecast horizons. A by-product of the algorithm is that the resulting forecasts are coherent at the various temporal aggregation levels, which is a benefit the authors stress from a decision-making perspective. In follow up work **kourentzes2016forecasting** investigate the impact of temporal aggregation on promotional indicator variables. They demonstrate the usefulness of the method in the presence of

promotional information. Similar benefits are seen in the case of intermittent demand when relying on multiple temporal aggregation levels instead of a single level (**petropoulos2015forecast**), as well as in inventory management (**barrow2016distributions**; **petropoulos2019inventory**), again alluding to the connection between temporal aggregation and supply chain management. **kourentzes2016forecasting** compare directly the use of optimally identified single temporal aggregation levels to using multiple levels. They find that the latter is preferable, even when the underlying data generating process is known. Using multiple aggregation levels results in combinations of forecasts, therefore reducing the modelling risk.

AthEtAl2017 combine the idea of using the multiple temporal aggregation levels of MAPA with hierarchical reconciliation, introducing the notion of temporal hierarchies. With temporal hierarchies, multiple levels of temporal aggregation of a time series (typically up to the annual level) are constructed from the original series and modelled independently. Subsequently, the forecasts are reconciled using the approaches outlined previously.

Let the original time series y_t , with $t = 1, \dots, T$, be observed at a sampling frequency $1/m$ (e.g., $m = 12$ for monthly data). The aggregation levels $\{k_1, \dots, k_p\}$ are the p factors of m in ascending order, where $k_1 = 1$ and $k_p = m$. For each factor k of m , the non-overlapping temporally aggregated time series is constructed as:

$$x_j^{[k]} = \sum_{t=t^*+(j-1)k}^{t^*+jk-1} y_t,$$

where $j = 1, \dots, \lfloor T/k \rfloor$ and $t^* = T - \lfloor T/m \rfloor m + 1$, ensuring that all aggregation levels have complete aggregation windows. Note that $x_j^{[1]} = y_t$. The complete hierarchy progresses at the observation index of the most aggregate level, which we define as τ (this corresponds to j at that level). For each aggregation level, we stack the observations in column vectors

$$\mathbf{x}_\tau^{[k]} = \begin{bmatrix} x_{m_k(\tau-1)+1}^{[k]} \\ x_{m_k(\tau-1)+2}^{[k]} \\ \vdots \\ x_{m_k\tau}^{[k]} \end{bmatrix},$$

where $m_k = m/k$, $\tau = 1, \dots, N$, and $N = T/m$. Collecting these in one column vector, we obtain

$$\mathbf{x}_\tau = \begin{bmatrix} x_\tau^{[m]} \\ x_\tau^{[m-1]} \\ \vdots \\ x_\tau^{[1]} \end{bmatrix}.$$

The structural representation becomes $\mathbf{x}_\tau = \mathbf{S}\mathbf{x}_\tau^{[1]}$ with $\mathbf{S} = \begin{bmatrix} \mathbf{A} \\ \mathbf{I} \end{bmatrix}$ and

$$\mathbf{A} = \begin{bmatrix} \mathbf{1}'_m \\ \mathbf{I}_{m/k_{p-1}} \otimes \mathbf{1}'_{k_{p-1}} \\ \vdots \\ \mathbf{I}_{m/k_2} \otimes \mathbf{1}'_{k_2} \end{bmatrix}.$$

We have used \mathbf{x}_τ in the notation to clearly distinguish a temporal hierarchy from a cross-sectional one that uses \mathbf{y}_t . An example for a quarterly time series ($m = 4$) is provided in Figure 5. If there are multiple seasonalities that are not integer multiples of each other, the resulting additional temporal aggregations can simply be stacked in \mathbf{x}_τ , and \mathbf{A} can be extended accordingly.

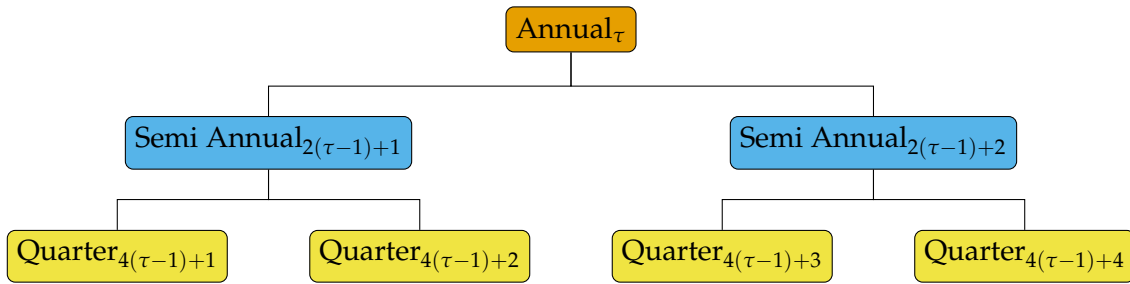


Figure 5: A temporal hierarchy for a quarterly time series at year τ , with $m = 4$.

An important difference between a temporal hierarchy and the data structure in its predecessor MAPA, is that the use of the p factors of m simplifies the hierarchical structure and allows the direct implementation of hierarchical reconciliation, with the major advantage being that now there are no model restrictions, and at each level different forecasting models/methods can be used. Further, there is added flexibility in the combination of the forecasts between levels, through \mathbf{G} . In contrast to cross-sectional hierarchies, temporal hierarchies can be constructed for any time series, requiring no additional data. However, this comes at the cost of estimation inefficiencies. Since there are only $N = T/m$ observations of the complete hierarchy, this significantly affects any estimation required in the approximation of \mathbf{W}_h , which motivated

AthEtAl2017 to propose the aforementioned structural scaling that requires no estimation. **NysEtAl2020** propose and demonstrate the benefits of more advanced approximations that take advantage of the autocorrelation in the forecast errors, when there is sufficient data, and demonstrate their merits in a short-term electricity load forecasting application. **NysEtAl2021** go one step further and propose an estimator based on an eigendecomposition of the temporal correlation matrix, which is able to perform well even with relatively limited data compared with the dimension of the temporal hierarchy.

AthEtAl2017 provide simulation and empirical evidence of the accuracy benefits of forecasting with temporal hierarchies. These benefits increase with added modelling uncertainty and at more temporally aggregate levels, echoing the arguments and evidence in **kourentzes2014improving** and **kourentzes2016forecasting**. The good performance of temporal hierarchies has been evidenced in various follow-up studies (**YanEtAl2017_temporal**; **JeoEtAl2019**; **NysEtAl2020**; **NysEtAl2021**; **kourentzes2021visitor**).

4.3 Cross-temporal reconciliation

Observe that cross-sectional hierarchies are described at time t for time series y_t , while temporal hierarchies at time τ for x_τ contain temporally aggregate views of y_t . Given that hierarchical methods are motivated to support forecasting at various levels of a hierarchy, cross-sectional and temporal hierarchies on their own may have limitations. For example, consider the case of forecasting for a grocery retailer. Let $y_t \in \mathbf{y}_t$ describe the sales of a particular ice cream product. Cross-sectionally this may be grouped with other similar products, or with product sales within a geographic demarcation, and so on. The further we aggregate, the less relevant a forecast becomes for the specific period t and at the granularity of y_t . Although we may be interested in the daily sales of a specific ice cream product, it is unlikely that we are interested in the daily sales at the top-level of the hierarchy describing the total company sales. Similarly, from the temporal point of view, it is unlikely that we are interested in the sales of y_t in time increments τ , for example the sales of a particular ice cream product in several years. A more aggregate view across products and time units is typically more relevant for decision makers, with many nodes in a hierarchy having the role of statistical devices that improve the quality of the overall coherent forecast, rather than being directly connected with some supported decision (**AthKou2021**).

Motivated by this, **KouAth2019** proposed the notion of cross-temporal hierarchies, where the hierarchy spans across both cross-sections and time, more accurately mapping the various

forecasts required by decision-makers and stakeholders. They show that sequential reconciliation across the cross-sectional and temporal dimensions, irrespective of order, does not always result in coherent forecasts, and address this by estimating all cross-sectional G_k , across the k temporal aggregation levels, which are then averaged in a common cross-temporal \bar{G} . They obtain holistically coherent forecasts, which are also the most accurate. In their experiments the temporal reconciliation provided the biggest accuracy gain.

Motivated by the sequential algorithm of **KouAth2019**, **Di_FonGir2022a** propose an iterative version whereby the forecasts are alternately reconciled in temporal and cross-sectional dimensions in a cyclic fashion, and find that it produces more accurate forecasts.

A number of contributions recognise that there are potential accuracy benefits in reconciling across both cross-sectional and temporal dimensions. **SpiEtAl2020**, **YagEtAl2019**, and **PunEtAl2020**, apply sequentially temporal and cross-sectional approaches, with **YagEtAl2019** experimenting with the order of reconciliation as well. They all identify accuracy benefits, but do not establish holistic coherence. This sequential approach is discussed and improved by **di2023spatio**.

Rather than separately reconciling the cross-sectional and temporal dimensions, **Di_FonGir2022a** proposed a single reconciliation step, using the full cross-temporal hierarchy. Following our notation, from the cross-sectional y_t at the most temporally disaggregate level, let $y_{i,t}$ denote its i th element, $i = 1, \dots, n$. For each i , we construct all the temporally aggregated variants, giving a vector of length p :

$$\mathbf{x}_{i,\tau} = \begin{bmatrix} x_{i,\tau}^{[m]} \\ \vdots \\ x_{i,\tau}^{[1]} \end{bmatrix}.$$

These can then be stacked into a long vector:

$$\mathbf{x}_\tau = \begin{bmatrix} \mathbf{x}_{1,\tau} \\ \vdots \\ \mathbf{x}_{n,\tau} \end{bmatrix}.$$

With S_{cs} and S_{te} denoting the structural matrices for the cross-sectional and temporal reconciliations respectively, the cross-temporal structural matrix is $S_{ct} = S_{cs} \otimes S_{te}$, so that

$$\mathbf{x}_\tau = S_{ct} \mathbf{b}_\tau^{[1]},$$

where the bottom-level series

$$\mathbf{b}_{\tau}^{[1]} = \begin{bmatrix} \mathbf{b}_{1,\tau}^{[1]} \\ \vdots \\ \mathbf{b}_{n_b,\tau}^{[1]} \end{bmatrix}.$$

Di_FonGir2022a develop optimal cross-temporal reconciliation and evaluate it against the heuristic approach of **KouAth2019** and variants. They report a relatively larger contribution to accuracy from the temporal side, and find that the optimal approaches tend to be outperformed by the heuristic approaches. We note that in their experiments all the approximations used for the cross-temporal \mathbf{W}_h required some estimation, which given the size of the cross-temporal matrices may explain the findings (**pritularga2021stochastic**).

Cross-temporally reconciled forecasts offer relatively limited accuracy gains compared to one-way reconciled forecasts (primarily temporally reconciled), with their major benefit being the qualitative difference of being coherent across both dimensions. This is impactful within a decision-making context and can be seen as a tool to sidestep organisational information silos, and achieve aligned plans across different functions in organisations (**KouAth2019**). **kourentzes2022toward** argues that cross-temporally coherent forecasts offer a pathway towards so-called “one-number” forecasts, enabling the integration of independent forecasts built for different functions and decisions that are typically based on different information, with different horizons, and purposes. If these forecasts remain disconnected, they can lead to misaligned decisions and organisational friction. In the cross-temporal case, some nodes of the hierarchy are by-products of the structure, rather than directly connected with some decision (**AthKou2021**). This raises questions how to best evaluate the quality of these forecasts, given that the relevant metrics for different decisions may vary in an organisational context.

5 Probabilistic hierarchical forecasting and reconciliation

The period that saw growth in the development of methods for hierarchical forecasting and reconciliation coincided with an increasing awareness of the importance of probabilistic forecasting. Therefore it is unsurprising that in recent years, a number of papers have attempted to tackle the problem of probabilistic hierarchical forecasting. As in the case for point forecasting, methods for obtaining probabilistic hierarchical forecasts can be split into bottom-up, top-down and reconciliation approaches, although some algorithms combine elements of more than one

approach. Note that Bayesian approaches to probabilistic forecast reconciliation are discussed in [Section 3.8](#).

Bottom-up methods of probabilistic hierarchical forecasting were introduced by **Ben_TaiEtAl2017** and subsequently expanded on in **Ben_TaiEtAl2021**. The algorithm is initialised by generating a Monte Carlo sample from the predictive distribution of each bottom-level variable, which by construction are independent. To induce dependence, these samples are first ranked, and then permuted series-wise. The permutations are designed to ensure that the samples from bottom-level predictive distributions have the same empirical copula as ‘in-sample’ forecast errors (residuals), taking care to exploit the hierarchical structure to avoid dealing with very high-dimensional copulas. **The rationale for matching the empirical copulas of the samples and residuals is to bring dependence information into the reconciliation procedure, since this is known to work well in point forecasting.** The samples are then aggregated to yield a sample from the predictive distribution of all top-level series. This is summarised with a simple example in [Figure 6](#).

$$\begin{array}{ccccc}
 A \downarrow & B \downarrow & & A \circlearrowleft & B \circlearrowleft \\
 \begin{bmatrix} 1.4 \\ 2.3 \\ 1.7 \\ 2.1 \end{bmatrix} & \begin{bmatrix} 3.6 \\ 5.3 \\ 2.2 \\ 6.4 \end{bmatrix} & \xRightarrow{\text{rank}} & \begin{bmatrix} 1.4 \\ 1.7 \\ 2.1 \\ 2.3 \end{bmatrix} & \begin{bmatrix} 2.2 \\ 3.6 \\ 5.3 \\ 6.4 \end{bmatrix} \\
 & & & \xRightarrow{\text{permute}} & \\
 A & B & & & \\
 \begin{bmatrix} 1.4 \\ 1.7 \\ 2.1 \\ 2.3 \end{bmatrix} & + & \begin{bmatrix} 3.6 \\ 2.2 \\ 6.4 \\ 5.3 \end{bmatrix} & \xRightarrow{\text{aggregate}} & T \\
 & & & & \begin{bmatrix} 5.0 \\ 3.9 \\ 8.5 \\ 7.6 \end{bmatrix}
 \end{array}$$

Figure 6: A toy example describing bottom-up approaches for probabilistic forecasting for a simple 3-variable hierarchy with $T = A + B$. A sample of size $K = 4$ has been drawn for two bottom-level series A and B . These are then ranked from smallest to largest, then permuted so that their empirical copula matches that of the residuals. The residuals are not shown here, but in this example, the smallest residual in A coincides with the second smallest residual in B , the second smallest residual in A coincides with the smallest residual in B and so on. Finally A and B are aggregated to give a sample from the predicted distribution of the total series, T .

While the bottom-up algorithm does not use a top-level forecast at all, both **Ben_TaiEtAl2017** and **Ben_TaiEtAl2021** propose an extension that incorporates top-level information, by adjusting the mean of each series to be equal to a reconciled point forecast; for example, in **Ben_TaiEtAl2021**, MinT is used. A shortcoming of the bottom-up approach and its extension is that the sample drawn from the predictive distribution and the sample of training data must be of equal size, making it ill-suited to problems with a small amount of training data. This is overcome by **PanZho2018** and **ZhaEtAl2019** who rather than using Monte Carlo, estimate predictive quantiles directly via quantile regression.

PanZho2018 also propose a top-down method for producing probabilistically coherent forecasts whereby quantile forecasts are first produced for all series. Proportions for top-down disaggregation are then found by taking the ratio of a forecast of a child node to the ratio of the forecast of the parent node. These are then applied to the original top-level forecasts. To the best of our knowledge the only other top-down method for coherent forecasting has been proposed by **DasEtAl2022** who model future proportions based on past proportions using a combination of an LSTM and multi-head self attention architecture. A sample from the predictive distribution of the top level is generated and each observation from this sample is disaggregated according to the forecast proportions.

Similar to the point forecasting case, progress has been made in extending the two-step reconciliation approach whereby probabilistic forecasts are produced from all series, and then reconciled to be coherent in a second step. **In the temporal reconciliation framework, JeoEtAl2019** propose drawing a sample from the predictive distribution of each series (both top- and bottom-level) and then stacking these into a matrix. The matrix can then be pre-multiplied by a projection matrix SG to obtain a sample from the coherent multivariate predictive distribution. One algorithm proposed by **JeoEtAl2019**, which they refer to as the ‘ranked sample’, orders the observations drawn from each predictive distribution before pre-multiplying by SG . This approach is described in [Figure 7](#) and corresponds to reconciling quantiles; an idea that has antecedents in **ShaHyn2017** who reconcile prediction intervals. Quantiles are only preserved under linear combinations when the data are perfectly dependent (**Kol2023**), however in cases where dependence between series is high, the method of **JeoEtAl2019** is shown to perform well.

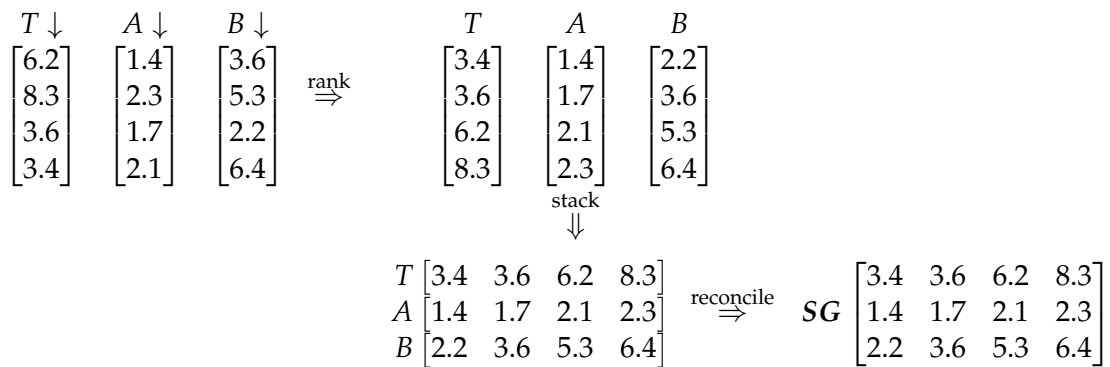


Figure 7: A toy example describing the ranked sample approach to probabilistic forecast reconciliation for a three variable hierarchy where $T = A + B$. A sample of size $K = 4$ has been drawn for the top-level series T and two bottom-level series A and B . These are then ranked from smallest to largest, then stacked, and reconciliation is applied (final reconciled forecasts depend on the choice of G , and are not shown here).

PanEtAl2020 make a number of contributions to probabilistic forecast reconciliation by providing formal definitions for coherence and reconciliation that justify Monte Carlo approaches as well as finding reconciled probabilistic forecasts for elliptical distributions (including the Gaussian). The **PanEtAl2020** framework allows any set of base forecasts, either univariate, or multivariate (the latter until then had not been considered in the hierarchical forecasting literature) to be reconciled using any reconciliation method. Reconciliation weights are trained by optimising with respect to a multivariate scoring rule. **RanEtAl2021** also optimise with respect to scoring rules for probabilistic forecasts, but rather than train reconciliation weights, they assume least squares reconciliation, with the novelty coming from training forecasting models and reconciling in an end-to-end fashion rather than via the usual two-step approach.

Despite very recent progress in forecast reconciliation, a number of open research questions remain. In the case of point forecasting, the optimality of certain reconciliation techniques such as MinT are now well understood. Similar results have not been derived for reconciliation methods in the probabilistic setting with the exception of **Wic2023** who derive the optimality of MinT in the Gaussian case. We expect a better understanding of the theoretical and empirical properties of different probabilistic forecast reconciliation techniques in coming years, including which methods produce forecasts that are well calibrated and whether their prediction intervals achieve correct coverage rates.

6 Significant applications

6.1 Tourism

The new concept of forecast reconciliation was first implemented on tourism data in **AthEtAl2009**. Tourism flows comprise aggregation structures across various dimensions. The most obvious of these are geographic hierarchies. At an international level, we have inbound travel from multiple countries or from regions within source countries, to multiple destination countries, airports or regions within the destination countries. The same applies for outbound travel while similar geographic divisions are natural for domestic travel. Further, policy makers and business planners are also interested in various characteristics of tourists. For example, the expenditure patterns of holiday makers are very different to those of business travelers or those visiting friends and relatives. Hence, grouped structures where geographic hierarchies are crossed with various attributes of interest also naturally arise. An attribute commonly observed for tourism flows is purpose of travel, which typically comprises holidays, visiting friends and relatives, business, and other.

AthEtAl2009 focus on two aggregation structures based on quarterly tourism flows. A grouped structure where geographic divisions (Australia, States, Capital city versus other) are crossed with purpose of travel (**HynEtAl2011**); and a pure geographic hierarchy where Australian domestic tourism flows are disaggregated by States, Zones and Regions. Coherent forecasts were generated from traditional bottom-up and top-down approaches based on historical proportions, as well a new top-down approach based on forecast proportions and forecast reconciliation. The paper found that the proposed top-down and reconciliation approaches improve forecast accuracy compared to the traditional approaches, and provides detailed forecasts for Australian domestic tourism flows, identifying some key features at both the aggregate and disaggregate levels that are crucial for informing policy makers. Variations of these Australian tourism data have since become ubiquitous for benchmarking forecast reconciliation methods. For example, **AboEtAl2022** propose what they refer to as conditional hierarchical forecasting, an approach based on machine learning classification methods that use time series features to select the reconciliation method for a hierarchy and evaluate the performance of the method based on the pure geographic hierarchy. **Gle2020** introduces an embedding reconciliation term that penalizes deviation from an aggregation structure, and uses the grouped structure to evaluate the forecasting performance of the proposed method, claiming improvements over MinT.

An updated and richer monthly Australian tourism data set was introduced in **WicEtAl2019**. The geographic hierarchy comprises 111 series. In particular the total tourism flow is disaggregated to 7 states, 27 zones and 76 regions. The hierarchical structure is crossed with the 4 purposes of travel, resulting in a total of 555 time series, of which 525 are unique (**Di_FonGir2022b**). This data set or close variations of it, is considered in several studies, including **KouAth2019**, **HolEtAl2021**, **SpiEtAl2021** and **Di_FonGir2022b**.

KarMal2019 study ten different data sets, including three related to tourism. The first one is a variation of the quarterly data set introduced in **AthEtAl2009**, using only the hierarchical structure based on the geographic divisions (**KarEtAl2023**). They also consider a weekly dataset of flight passengers between Melbourne and Sydney for an airline, disaggregated by flight ticket class type (first class, business, and economy), and monthly departures from Australia disaggregated into permanent, long-term, and short-term departures, and then between residents and visitors for the cases of long-term and short-term departures. The paper presents forecasting methods for hierarchical time series based on Support Vector Machines. These are compared to traditional single-level approaches with the conclusion that a major limitation of the proposed methods is the lack of data in a time series context.

AthEtAl2022 and **kourentzes2021visitor** focus on tourism flows amid the COVID-19 pandemic. **AthEtAl2022** model both international inbound and domestic flows for the case of Australia, while **kourentzes2021visitor** analyse international arrivals for the case of South Africa. The papers argue for the use of forecast reconciliation in order to generate robust forecasts for tourism flows during the pre-COVID period.

6.2 Macroeconomics

Since macroeconomics is the study of aggregate economic phenomena, it is not surprising that this field has provided fertile ground for hierarchical data. For example, Gross Domestic Product (GDP) is constructed as an aggregate of individual components. The expenditure method constructs GDP as an aggregate of expenditure characterized by consumption, investment, imports and exports, while the income approach aggregates variables such as the gross operating surplus of firms with employee compensation. Thus there are two hierarchies involved with the same aggregated series. The structural formulation of the reconciliation problem does not allow for this scenario, but the more general constraint formulation ([Section 2.2](#)) does.

As noted in [Section 3.1](#), reconciliation of estimates (as distinct from forecast reconciliation) has a long history in macroeconomics. For forecasting problems, forecasts of the disaggregate series may be of direct interest, otherwise the objective may be to improve accuracy by leveraging forecasts of the disaggregate series via reconciliation. **AthEtAl2020**, **BisEtAl2020** and **DFG22** find evidence in favour of forecast reconciliation in both of these settings for Australian GDP. In particular for both point and probabilistic forecasts, improvements in forecast accuracy over base and bottom-up methods can be achieved using forecast reconciliation. The MinT method is found to work best overall, while when attention is restricted to only the top-level series, weighted least squares is found to be more accurate. **Di_FonGir2022a** use the constraint representation with a cross-temporal framework, and show that it leads to better forecasts on this data set.

Another important macroeconomic variable that admits a hierarchical structure is the consumer price index (CPI). The CPI is constructed as a weighted price index of a basket of goods. As such the novel aspect to reconciliation in this setting is that rather than a summing matrix consisting only of ones and zeros, the S matrix includes weights that can change over time. In an early application of reconciliation methodology, **CapEtAl2010** find that OLS reconciliation can improve upon a bottom-up approach for some but not all components of Mexican CPI. For overall CPI, OLS improves upon a bottom-up approach, however the difference between the forecast accuracy of these methods is not found to be statistically significant. **Wei2018** considers

UK CPI and finds that reconciliation using OLS performs better than traditional approaches for 1-month-ahead forecasts, but that middle out approaches work better for longer horizons. **Wei2018** also considers inflation volatility; in this case, reconciliation methods do not outperform bottom-up.

While hierarchical structures are common in macroeconomics, certain details concerning the construction of these datasets suggest new directions in reconciliation methodology. One example, **KooEtAl2022** consider a geographic hierarchy of productivity, with a model that includes growth rates in national and regional output. However, since chain volume measures are used to construct output, with different price deflators for each region, regional growth rates only add up to the national growth rate as an approximation. Since the usual aggregation constraint only holds approximately, **KooEtAl2022** recommend shrinking towards the aggregation constraint via a Bayesian approach rather than imposing a hard constraint. Considering unemployment data from multiple labour force surveys in Brazil, **LiEtAl2022** introduce robust estimation in the reconciliation stage with primary aim to address measurement issues occurring in the original time series. These methods are likely to generalise to other hierarchical forecasting problems in macroeconomics and other disciplines.

6.3 Energy

Energy applications are widespread in the forecast reconciliation literature due to the natural geographic hierarchies that arise in energy distribution. For example, both the GEFCom2012 (**Gefcom2012**) and GEFcom2017 (**Gefcom2017**) energy forecasting competitions included hierarchical electricity load data, although none of the participants took advantage of the hierarchical structure to improve their forecasts.

One of the earliest uses of forecast reconciliation applied to energy data was **Van_ErvCug2015**, discussed earlier, who applied their methods to electricity demand data from Électricité de France, disaggregated into 17 tariff groups. **Other early uses of reconciliation methods for forecasting short-term (24h) electricity demand include AlmEtAl2016 who applied OLS reconciliation to a hierarchy disaggregated by grid supply point and voltage level, daSEtAl2019 who use Bayesian estimation, and MeiEtAl2023 who propose methods robust to the influence of outliers.** Further examples of application to load forecasting include **FenZha2020**, who compared point forecasts obtained from bottom-up, OLS and MinT on hourly load data from 13 buildings in Texas, USA and **NesEtAl2020** who applied reconciliation methods to electricity load data, comparing several of the probabilistic forecasting methods discussed in [Section 5](#) when applied to a small set of 24 power meters located in Rolle, Switzerland.

Three interesting examples of forecast reconciliation of electricity load are due to Ben Taieb and his coauthors. **Ben_TaiEtAl2017_AAAI** propose a regularized version of MinT reconciliation with penalties analogous to those used in LASSO and elastic-net regressions, giving sparse adjustments to the base forecasts. They apply the method to electricity consumption measured by 5701 smart meters with a rich geographic hierarchy. **Ben_TaiEtAl2017** and **Ben_TaiEtAl2021** each developed new probabilistic reconciliation methods (discussed in [Section 5](#)), and applied them to the same dataset.

ZhaEtAl2019 proposed a computationally simpler variation of the method of **Ben_TaiEtAl2017**, and applied it to two public data sets: the ISO New England data from **Gefcom2017**, and some Irish smart meter data aggregated to groups of 100 customers. **Roa2019** also studies the ISO New England data from the **Gefcom2017** and proposes a method based on generating reconciled quantile forecasts using a gradient boosted model which is shown to outperform the benchmark.

BreHua2021 also tackle load forecasting, but introduce an innovation by including an “aggregation algorithm” before reconciling the forecasts. This aggregation algorithm involves computing revised base forecasts that are linear combinations of the base forecasts of all series. The revised base forecasts are then reconciled using OLS reconciliation. This approach means that any cross-sectional relationships between series are modelled in the aggregation step rather than the reconciliation step. The authors applied this approach to the same dataset used by **Ben_TaiEtAl2017**.

Forecast reconciliation has also been used in solar generation forecasts. Here, power generated by distributed photovoltaics (PV) are naturally disaggregated in a geographic hierarchy such as transmission zones, distribution nodes, PV plants, subsystems and inverters. **YanEtAl2017_geography** explored the application of MinT reconciliation and some of its special cases to hourly generation data from 318 power plants in California, USA, while **YagEtAl2020** applies probabilistic (Gaussian) MinT reconciliation to the same data set. The data set is used again in **YanEtAl2017_temporal**, applying temporal reconciliation methods, and in **YagEtAl2019**, where both cross-sectional and temporal reconciliation are considered. They apply the cross-sectional and temporal reconciliations sequentially, rather than simultaneously. **Di_Fonzo2022** and **Di_FonGir2022a** critique this two-step approach and argue for simultaneous cross-temporal reconciliation, but show that the two approaches can be equivalent if the covariance matrices used in both steps are constant across levels and time granularities. They further show how the forecasts can be constrained to be non-negative using the simple approach of setting any negative bottom-level reconciled forecasts to zero, and then aggregating the results. Another

application to solar power is **PanZho2018**, who applied the probabilistic reconciliation methods they developed (discussed in [Section 5](#)) to 5-minute solar power data for about 6000 simulated PV plants in Florida, USA.

Forecast reconciliation has been extensively applied to wind power forecasting, see **JeoEtAl2019**, **BaiPin2019** and **HanEtAl2023**. Also of particular note is **GilEtAl2018** who produce probabilistic forecasts of wind power for individual turbines, and for the entire wind farm. An innovative feature is the use of a weighted aggregation based on an elastic net penalized regression.

Other applications of forecast reconciliation in the general field of energy include **BuzEtAl2021** who consider probabilistic forecasting of loads at electric vehicle charging stations, and **BerEtAl2021** who consider heat load forecasting. The latter propose an adaptive reconciliation method for temporal hierarchies by allowing for time-varying weights.

6.4 Mortality rates

Applications of forecast reconciliation in mortality started with **ShaHyn2017**, who forecast Japanese mortality rates disaggregated by age, sex, and by a geographic hierarchy of 47 prefectures within 8 regions. The base forecasts were obtained using a functional data method. Because mortality rates do not sum directly, they proposed an aggregation matrix A comprising population ratios. For example, for the mortality rate of 50-year-old females within a region, the non-zero values of the corresponding row of the A matrix contains the 50-year-old female population of each prefecture divided by the total 50-year-old female population of the region. Thus, the aggregation matrix A is time-varying, and the values for the future time periods were forecast using univariate time series models. WLS reconciliation was used. **ShaHab2017** provide an application to annuity pricing using an identical approach to the same data set. **LiHyn2021** used a similar formulation, but applied MinT reconciliation to US age-sex-specific mortality rates, and explored future mortality inequality.

LiEtAl2019 consider forecasting mortality due to different causes of death and show that forecast accuracy is improved by reconciliation. Although the data form a simple two-level hierarchy with bottom-level series and their aggregate, the authors combine individual causes of death into middle level series. This is done via hierarchical clustering on the data. Augmenting the hierarchy with middle level series in this fashion, leads to further improvement in forecast accuracy.

LiTan2019 use a forecast reconciliation approach to forecast a longevity divergence index (LDI), used to compute the Swiss Re Kortis bond. The most disaggregated series in the hierarchy are

age-specific mortality improvement rates in the US and UK, while the most aggregated series are LDI values, expressed as linear combinations of these disaggregated series. They use a MinT approach to generate probabilistic forecasts, following [JeoEtAl2019](#).

6.5 Retail demand & supply chain

Forecasting for demand planning has been an attractive application for hierarchical forecasting and therefore it has attracted attention in the literature ([DanEtAl2013](#)). [rostami2015non](#) look at the conditions where top-down or bottom-up is favourable assuming that the disaggregate series are ARIMA(0,1,1) processes, using the reconciliation approach as a benchmark. The latter is found to be more accurate overall, with the bottom-up method outperforming it in some cases, particularly for the bottom-level time series. [YanEtAl2016](#) using the publicly available Dominick's Finer Food dataset find that the shrinkage estimator for W_1 performed better than classic hierarchical methods (bottom-up and top-down). [OliRam2019](#) find the same conclusion on data from a Portuguese supermarket. [MirEtAl2021](#) look at sales of a major European brewery and find that hierarchical reconciliation performs better than base forecasts. They also propose combining the bottom-level forecasts of different hierarchical forecasting methods and then construct the forecasts for the rest of the hierarchy using a bottom-up approach. They argue that this method has the advantage that it eliminates the need to select a hierarchical approach and find small gains over the reconciliation method. [KarMal2019](#) explore the performance of hierarchical forecasting on sales in the travel retail industry. They do not consider the reconciliation approach and find bottom-up to be best.

[villegas2018supply](#) propose encapsulating cross-sectional hierarchical reconciliation in a state-space formulation, together with the forecasting model. Using simulations and an empirical investigation on a Spanish grocery retailer, they find that the standard reconciliation performs well for short horizons (1–3 days), while for longer horizons (4–7 days) the state-space based reconciliation is best and is also the overall most accurate method.

Other applications include the contributions of [abolghasemi2022model](#), [AboEtAl2022](#), and [SpiEtAl2021](#) who investigate the application of hierarchical forecasting on data from a food manufacturer in Australia. Finally, the M5 forecasting competition used data from Walmart, a major US retailer, with a grouped time series structure, providing a test bed for reconciliation methods ([MakEtAl2020](#)). We provide further details of these works in [Section 3.7](#).

6.6 Intermittent demand

Hierarchical forecasting has seen some application in intermittent demand modelling. However, in reviewing these papers it is useful to consider that the literature has progressed substantially

over the last years in terms of how to evaluate forecasts of intermittent demand time series (**kourentzes2014intermittent**; **kolassa2016evaluating**; **AthKou2021**), recognising that classic error metrics, especially those based on absolute errors, are often inappropriate. Direct evaluation on decision metrics, such as inventory metrics, or the predictive distribution are preferable.

One of the earliest works using hierarchical forecasting for intermittent demand is by **moon2012development** who looked at top-down hierarchical forecasts, against combination and base forecasts for predicting spare parts for the South Korean navy. In the reported inventory cost, the top-down approach offers benefits in some settings, but overall is outperformed by combination methods. It should be noted that the hierarchy that was used in this work was constructed by the researchers and its eventual structure may have been significant for the findings.

petropoulos2015forecast provide a translation of the temporal MAPA algorithm (**kourentzes2014improving**) for intermittent demand. It is used to forecast spare parts for the UK Royal Air Force. They find it outperforms various benchmarks, including ADIDA that relies on a single temporal aggregation level (**nikolopoulos2011aggregate**), and various combinations of forecasts. However, the empirical evaluation lacks decision or predictive distribution related metrics.

LiLim2018 provide a top-down-like hierarchical forecasting method for predicting intermittent demand in fashion retailing. Their approach produces separately daily forecasts of the aggregate demand across multiple items, and a forecast of the inter-demand interval and demand size for each individual item. The latter forecasts are used to prorate the total forecast into the individual items. Although the proposed algorithm performs well against benchmarks, the empirical evaluation lacks other cross-sectional reconciliation benchmarks and strong performance metrics.

KouAth2021 investigate the use of temporal hierarchies for intermittent demand forecasting for aerospace spare parts. Their motivation is that a method that predicts well the intermittent pattern should be able to demonstrate the various patterns (such as seasonality and trend) that may appear when the data are explored at lower sampling frequency. They demonstrate that the use of temporal hierarchies allows the capture of these patterns at higher temporal aggregation levels and combining this information with intermittent demand forecasts of the original time series, results in reconciled forecasts that dominate the base forecasts on a variety of metrics and horizons. They obtain prediction intervals using the empirical distribution of the reconciled forecasts and provide a heuristic to ensure non-negative forecasts. **The use of temporal hierarchies in an intermittent demand setting has more recently been considered by**

SanEtAl2023 who focus on the estimation of forecasting downward trends in a scenario of obsolescence.

6.7 Healthcare, accidents & emergencies

AthEtAl2017 applied temporal hierarchies to predict weekly admissions for Accidents & Emergency wards in UK hospitals. The volume of patients relates to different decisions in the operations of the wards, from staff scheduling, to procuring consumables, training and hiring of staff, etc. They showed that temporally coherent forecasts dominated base forecasts in all cases.

pritularga2021stochastic looked at weekly Accidents & Emergency cases for a specific hospital, and were interesting in producing cross-sectionally coherent forecasts across various patient demarcations. They compared a variety of approximations for W_h , controlling for the sample size, and found that hierarchical forecasts were universally better than the base forecasts. They also showed that the complexity of the approximation is important, with simpler ones performing best at smaller sample sizes, and more complex ones gaining a substantial advantage when there were sufficiently long time series.

Wei2018 investigated improving the staffing for a large UK teaching hospital. Hierarchical forecasts were found to provide more accurate forecasts than the benchmark used by the hospital. Further, when the forecasts were used in a staffing model, they resulted in cheaper operations and less under-staffing.

GibEtAl2021 forecast weighted influenza-like illness (wILI), with weights corresponding to the population size of different U.S. regions. Enforcing probabilistic coherence increases forecast skill for most models when tested over multiple flu seasons. This provides further evidence of the benefits of using forecast reconciliation with geographical hierarchies.

7 Open-source software implementations

The first available open-source implementation of forecast reconciliation methods was the `hts` package for R (**Rhts1**), which helped popularise cross-sectional point forecast reconciliation methods in business and industry. The `hts` package has continued to be developed, and its latest version (**Rhts**) includes implementations of **WicEtAl2019** and **WicEtAl2020**. The game-theoretic approach of **Van_ErvCug2015** is implemented in the `gtop` package for R (**Rgtop**). Temporal point forecast reconciliation is provided in the `thief` package for R (**Rthief**). Cross-sectional, temporal, and cross-temporal point forecast reconciliation with optional non-negativity constraints is provided by the `FoReco` package for R (**RFoReco**). The score optimization approach

of **PanEtAl2020** is implemented in the ProbReco package (**RProbReco**). Probabilistic cross-sectional forecast reconciliation is also included in the `fabletools` package for R (**Rfabletools**), with a very simple user interface for specifying complicated hierarchical and grouping structures.

There are also several Python implementations of the methods, including the `pyhts` package of **pyhts**, which is a python translation of **Rhts**, and Darts (**pydarts**; **pydarts_jmlr**) which provides similar functionality. The `hierarchicalforecast` package (**pyhierarchicalforecast**; **Olivares2023**) provides a more comprehensive suite of functions covering both point and probabilistic forecast reconciliation including the methods of **Ben_TaiEtAl2017**, **Ben_TaiEtAl2019**, **WicEtAl2019**, and **PanEtAl2020**. The method of **RanEtAl2021** is implemented in `GluonTS` (**GluonTS**).

8 Conclusion

Research into hierarchical time series and forecast reconciliation has seen great success and impact, particularly over the last decade. At this juncture we wish to speculate on what the next decade holds in store by identifying some key open questions. We anticipate growth in the following areas.

The sheer size of some hierarchical forecasting problems can lead to computational difficulties. It is common to have n_b , the number of the most disaggregated series, to be well over 1 million, and with many grouping factors, the number of aggregated series can be even larger. This leads to very large matrices that need to be inverted, even when using the more efficient constraint matrix approach of [Section 2.2](#). Sparse matrix algebra has helped considerably with this problem, and when coupled with the Lanczos algorithm (**paige1982algorithm**), very large problems can be handled (**HynEtAl2011**). However, further improvements may be possible using dimension reduction methods, along the lines of **wang2020improved**.

A related problem arises with the MinT approach where the covariance matrix of the base forecasts needs to be computed. Here, shrinkage methods have been used to good effect (**WicEtAl2019**), but these don't scale to very large matrices. Alternative sparse or low-rank approaches (**lam2020high**) would be a welcome addition to the literature.

The breadth of data to which reconciliation methods need to be applied necessitates the extension of methods to non-standard domains. This includes non-negative data (**WicEtAl2020**), discrete data (for which **CorEtAl2021**, **OliEtAl2021** and **Zambon2022** give early attempts to address the

issues), and finally mixtures of discrete and continuous data. The latter could be potentially useful for zero inflated data (which arise in intermittent sales data) or where there are hierarchies where some bottom-level series are best modelled as discrete, while the top -level series are best modelled as continuous. The development of algorithms to handle these cases, as well as understanding the theoretical properties of reconciliation in such settings, represents a bold research agenda.

While recent progress has been made in probabilistic forecasting, as discussed in [Section 5](#), there are a number of open questions on the properties of probabilistic forecast reconciliation. For example, what are the coverage properties of prediction intervals derived from reconciled forecasts, how do these depend on the coverage properties of the base forecasts themselves and does reconciliation even improve coverage relative to base forecasts? These questions are particularly vexed since reconciled probabilistic forecasts are defined on a domain that is a linear subspace.

An often stated advantage of coherent forecasts is that they have the potential to lead to aligned decisions. However the attempts to quantify this effect have been limited. Where gains in forecast accuracy due to forecast reconciliation have been established, forecast evaluation is often based on general purpose metrics such as RMSE and MAE as well as scaled versions thereof. These metrics do not explicitly penalise incoherence in forecasts particularly when, as is often the case, the forecasts of different variables are evaluated individually. The disconnect between metrics of forecast evaluation and the operational considerations of hierarchical forecasting may explain why hierarchical methods have not been popular in forecasting competitions such as the M5, even where the data follow a hierarchical structure. Therefore we anticipate the development of new forecast evaluation metrics that account for the multivariate and hierarchical nature of the data. Further, we concur with the view of [AthKou2021](#) that forecast evaluation must be integrated with the decisions made by agents at different levels of the hierarchies.

This issue opens up additional questions related to the game theoretic aspects of hierarchical forecasting. Most empirical work to date shows that even where reconciliation improves forecast accuracy, these improvements do not occur across all levels of the hierarchy. In some applications, improvements may be seen in forecasts of bottom-level series after reconciliation, while base forecasts at the top level still outperform the top-level reconciled forecast. In other applications, the reverse may be true. This has a number of interesting implications in an organisational setting where those making forecasts at different levels can be treated as separate agents. Can

reconciliation methods be found that lead to Pareto improvements across the hierarchy meaning all agents gain from reconciliation? If not, in a cooperative setting, can forecast reconciliation, and the decisions agents make based on these forecasts be aligned to improve overall welfare of the organisation? Also, in a competitive setting, can compensation mechanisms be developed to encourage agents at different levels of the hierarchy, each making forecasts based on their own information sets to share base forecasts for reconciliation? These open questions should stimulate research for years to come.

Acknowledgements

We thank Tommaso Di Fonzo, Xiaoqian Wang and Daniele Girolimetto for providing helpful comments on an earlier draft of this paper. Rob J Hyndman was funded by the Australian Government through the Australian Research Council Industrial Transformation Training Centre in Optimisation Technologies, Integrated Methodologies, and Applications (OPTIMA), Project ID IC200100009.