

Hierarchical Forecasting

Puwasala Gamakumara, Anastasios Panagiotelis, George Athanasopoulos and
Rob J Hyndman

1 Introduction

Accurate forecasting of key macroeconomic variables such as Gross Domestic Product (GDP), inflation, industrial production, has been at the forefront of economic research over many decades. Early approaches involved univariate models or at best low dimensional multivariate systems. The era of big data has led to the use of regularization and shrinkage methods such as dynamic factor models, Lasso, LARS, Bayesian VARs, in an effort to exploit the plethora of potentially useful predictors now available. These predictors commonly also include the components of the variables of interest. For instance, GDP is formed as an aggregate of consumption, government expenditure, investment and net exports, with each of these components also formed as aggregates of other economic variables. While the macroeconomic forecasting literature regularly uses such sub-indices as predictors, it does so in ways that fail to exploit accounting identities that describe known deterministic relationships between macroeconomic variables.

In this chapter we take a different approach. Over the past decade there has been a growing literature on forecasting collections of time series that follow aggregation constraints, known as hierarchical time series. Initially the aim of this literature was

Puwasala Gamakumara
Monash University, Clayton VIC 3800, Australia, e-mail: Puwasala.Gamakumara@monash.edu

Anastasios Panagiotelis
Monash University, Caulfield VIC 3162, Australia e-mail: Anastasios.Panagiotelis@monash.edu

George Athanasopoulos
Monash University, Caulfield VIC 3162, Australia e-mail: George.Athanasopoulos@monash.edu

Rob J Hyndman
Monash University, Clayton VIC 3800, Australia e-mail: Rob.Hyndman@monash.edu

to ensure that forecasts adhered to aggregation constraints thus ensuring aligned decision making. However in many empirical settings the forecast reconciliation methods designed to deal with this problem have also been shown to improve forecast accuracy. Examples include forecasting accidents and emergency admissions (Athanasopoulos et al. 2017), mortality rates (Shang & Hyndman 2017), prison populations (Athanasopoulos et al. 2019), retail sales (Villegas & Pedregal 2018), solar energy (Yang et al. 2017, Yagli et al. 2019), tourism demand (Athanasopoulos et al. 2009, Hyndman et al. 2011, Wickramasuriya et al. 2018), wind power generation (Zhang & Dong 2019). Since both aligned decision making and forecast accuracy are key concerns for economic agents and policy makers we propose the application of state of the art forecast reconciliation methods to macroeconomic forecasting. To the best of our knowledge the only application of forecast reconciliation methods to macroeconomics is the PhD thesis of (Weiss 2018) which focuses on point forecasting for inflation.

Also Capistrán et al. (2010)

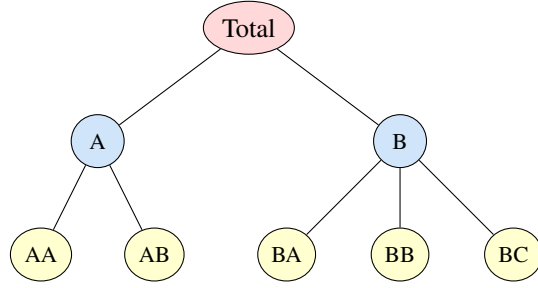
The remainder of the chapter is set out as follows. Section 2 introduces the concept of hierarchical time series, i.e. collections of time series with known linear constraints, with a particular emphasis on macroeconomic examples. Section 3 describes state-of-the-art forecast reconciliation techniques for point forecasts, while Section 4 describes the more recent extension of these techniques to probabilistic forecasting. Section 5 describes the data used in our empirical case study, namely Australian GDP data, that is represented using two alternative hierarchical structures. Section 6 provides details on the setup of our empirical study including metrics used for the evaluation of both point and probabilistic forecasts. Section 7 presents results and Section 8 concludes providing future avenues for research that are of particular relevance to the empirical macroeconomist.

2 Hierarchical time series

To simplify the introduction of some notation we use the simple two-level hierarchical structure shown in Figure 1. Denote as $y_{\text{Tot},t}$ the value observed at time t for the most aggregate (Total) series corresponding to level 0 of the hierarchy. Below level 0, denote as $y_{i,t}$ the value of the series corresponding to node i , observed at time t . For example, $y_{A,t}$ denotes the t th observation of the series corresponding to node A at level 1, $y_{AB,t}$ denotes the t th observation of the series corresponding to node AB at level 2, and so on.

Let $\mathbf{y}_t = (y_{\text{Tot},t}, y_{A,t}, y_{B,t}, y_{AA,t}, y_{AB,t}, y_{BA,t}, y_{BB,t}, y_{BC,t})'$ denote a vector containing observations across all series of the hierarchy at time t . Similarly denote as $\mathbf{b}_t = (y_{AA,t}, y_{AB,t}, y_{BA,t}, y_{BB,t}, y_{BC,t})'$ a vector containing observations only for the bottom-level series. In general, $\mathbf{y}_t \in \mathbb{R}^n$ and $\mathbf{b}_t \in \mathbb{R}^m$ where n denotes the number of total series in the structure, m the number of series at the bottom level, and $n > m$ always. In the simple example of Figure 1, $n = 8$ and $m = 5$.

Fig. 1 A simple two-level hierarchical structure.



Aggregation constraints dictate that $y_{\text{Tot}} = y_{A,t} + y_{B,t} = y_{AA,t} + y_{AB,t} + y_{BA,t} + y_{BB,t} + y_{BC,t}$, $y_{A,t} = y_{AA,t} + y_{AB,t}$ and $y_B = y_{BA,t} + y_{BB,t} + y_{BC,t}$. Hence we can write

$$\mathbf{y}_t = \mathbf{S}\mathbf{b}_t, \quad (1)$$

where

$$\mathbf{S} = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 \\ \mathbf{I}_5 \end{pmatrix}$$

is an $n \times m$ matrix referred to as the *summing matrix* and \mathbf{I}_m is an m -dimensional identity matrix. \mathbf{S} reflects the linear aggregation constraints and in particular how the bottom-level series aggregate to levels above. Thus, columns of \mathbf{S} span the linear subspace of \mathbb{R}^n for which the aggregation constraints hold. We refer to this as the *coherent subspace* and denote it by \mathfrak{s} . Notice that pre-multiplying a vector in \mathbb{R}^m by \mathbf{S} will result in an n -dimensional vector that lies in \mathfrak{s} .

Property 1. A hierarchical time series has observations that are *coherent*, i.e., $\mathbf{y}_t \in \mathfrak{s}$ for all t . We use the term coherent to describe not just \mathbf{y}_t but any vector in \mathfrak{s} .

Structures similar to the one shown in Figure 1 can be found in macroeconomics. For instance, in Section 5 we consider two alternative hierarchical structures for the case of GDP and its components. However, while this motivating example involves aggregation constraints, the mathematical framework we use can be applied for any general linear constraints, examples of which are ubiquitous in macroeconomics. For instance, the trade balance is computed as exports minus imports, while the consumer price index is computed as a weighted average of sub-indices, which are in turn weighted averages of sub-sub-indices, and so on. These structures can also be captured by an appropriately designed \mathbf{S} matrix.

An important alternative aggregation structure, also commonly found in macroeconomics, is one for which the most aggregate series is disaggregated by attributes of interest that are crossed, as distinct to nested which is the case for hierarchical time series. For example, industrial production may be disaggregated along the lines of geography or sector or both. We refer to this as a *grouped* structure. Figure 2 shows a simple example of such a structure. The Total series disaggregates into $y_{A,t}$

and $y_{B,t}$, but also into $y_{X,t}$ and $y_{Y,t}$, at level 1, and then into the bottom-level series, $\mathbf{b}_t = (y_{AX}, y_{AY}, y_{BX}, y_{BY})'$. Hence, in contrast to hierarchical structures, grouped time series do not naturally disaggregate in a unique manner.

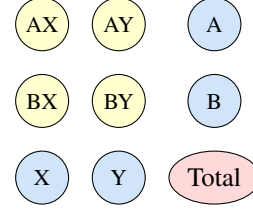


Fig. 2 A simple two-level grouped structure.

An important implementation of aggregation structures are *temporal hierarchies* introduced by Athanasopoulos et al. (2017). In this case the aggregation structure spans the time dimension and dictates how higher frequency data (e.g., monthly) are aggregated to lower frequencies (e.g. quarterly, annual). There is a vast literature that studies the effects of temporal aggregation, going back to the seminal work of Zellner & Montmarquette (1971), Amemiya & Wu (1972), Tiao (1972), Brewer (1973) and others, including Hotta (1993), Hotta & Cardoso Neto (1993), Marcellino (1999), Silvestrini et al. (2008). The main aim of this work is to find the single best level of aggregation for modelling and forecasting time series. In this literature, the analyses, results (whether theoretical or empirical) and inferences, are extremely heterogeneous making it very challenging to reach a consensus or some concrete conclusions. For example, Rossana & Seater (1995) who study the effect of aggregation on several key macroeconomic variables state,

“Quarterly data do not seem to suffer badly from temporal aggregation distortion, nor are they subject to the construction problems affecting monthly data. They therefore may be the optimal data for econometric analysis.”

A similar conclusion is reached by Nijman & Palm (1990). Silvestrini et al. (2008) consider forecasting French cash state deficit and provide empirical evidence of forecast accuracy gains from forecasting with the aggregate model rather than aggregating forecasts from the disaggregate model.

The vast majority of this literature concentrates on a single level of temporal aggregation (although there are some notable exceptions such as Andrawis et al. 2011, Kourentzes et al. 2014). Athanasopoulos et al. (2017) show that considering multiple levels of aggregation via temporal hierarchies and implementing forecast reconciliation approaches rather than single level approaches results in substantial gains in forecast accuracy across all levels of temporal aggregation.

3 Point forecasting

A requirement when forecasting hierarchical time series is that the forecasts adhere to the same aggregation constraints as the observed data; i.e., they are coherent.

Definition 1. A set of h -step-ahead forecasts $\tilde{y}_{T+h|T}$, stacked in the same order as y_t and generated using information up to and including time T , are said to be *coherent* if $\tilde{y}_{T+h|T} \in \mathfrak{s}$.

Hence, coherent forecasts of lower level series aggregate to their corresponding upper level series and vice versa.

Let us consider the smallest possible hierarchy with two bottom-level series, depicted in Figure 3, where $y_{\text{Tot}} = y_A + y_B$. While base forecasts could lie anywhere in \mathbb{R}^3 , the realisations and coherent forecasts lie in a two dimensional subspace $\mathfrak{s} \subset \mathbb{R}^3$.

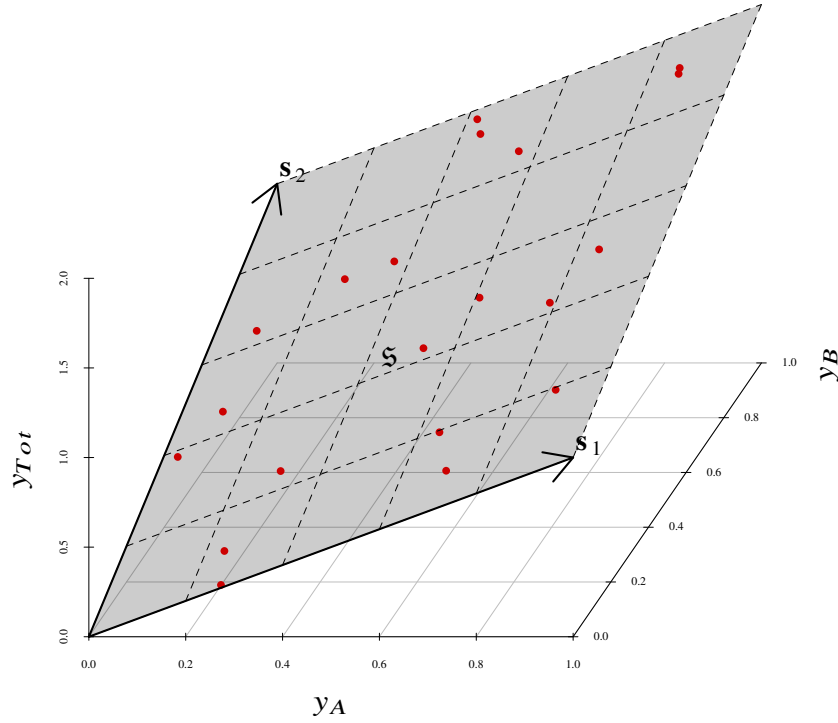


Fig. 3 Representation of a coherent subspace in a three dimensional hierarchy where $y_{\text{Tot}} = y_A + y_B$. The coherent subspace is depicted as a gray two dimensional plane labelled \mathfrak{s} . Note that the columns of $\vec{s}_1 = (1, 1, 0)'$ and $\vec{s}_2 = (1, 0, 1)'$ form a basis for \mathfrak{s} . The red points lying on \mathfrak{s} can be either realisations or coherent forecasts.

3.1 Single-level approaches

A common theme across all traditional approaches for forecasting hierarchical time series is that a single-level of aggregation is first selected and forecasts for that level are generated. These are then linearly combined to generate a set of coherent forecasts the rest of the structure.

3.1.1 Bottom-up

In the *bottom-up* approach, forecasts for the most disaggregate level are first generated. These are then aggregated to obtain forecasts for all other series of the hierarchy (Dunn et al. 1976). In general, this consists of first generating $\hat{\mathbf{b}}_{T+h|T} \in \mathbb{R}^m$, a set of h -step-ahead forecasts for the bottom-level series. For the simple hierarchical structure of Figure 1, $\hat{\mathbf{b}}_{T+h|T} = (\hat{y}_{AA,T+h|T}, \hat{y}_{AB,T+h|T}, \hat{y}_{BA,T+h|T}, \hat{y}_{BB,T+h|T}, \hat{y}_{BC,T+h|T})$, where $\hat{y}_{i,T+h|T}$ is the h -step-ahead forecast of the series corresponding to node i . A set of coherent forecasts for the whole hierarchy is then given by

$$\tilde{\mathbf{y}}_{T+h|T}^{\text{BU}} = \mathbf{S} \hat{\mathbf{b}}_{T+h|T}.$$

Generating bottom-up forecasts has the advantage of no information being lost due to aggregation. However, bottom-level data can potentially be highly volatile or very noisy and therefore challenging to forecast.

3.1.2 Top-down

In contrast, *top-down* approaches involve first generating forecasts for the most aggregate level and then disaggregating these down the hierarchy. In general, coherent forecasts generated from top-down approaches are given by

$$\tilde{\mathbf{y}}_{T+h|T}^{\text{TD}} = \mathbf{S} \mathbf{p} \hat{y}_{Tot,T+h|T},$$

where $\mathbf{p} = (p_1, \dots, p_m)'$ is an m -dimensional vector consisting of a set of proportions which disaggregate the top-level forecast $\hat{y}_{Tot,T+h|T}$ to forecasts for the bottom-level series; hence $\mathbf{p} \hat{y}_{Tot,T+h|T} = \hat{\mathbf{b}}_{T+h|T}$. These are then aggregated by the summing matrix \mathbf{S} .

Traditionally, proportions have been calculated based on the observed historical data. Gross & Sohl (1990) present and evaluate twenty-one alternative approaches. The most convenient attribute of these approaches is their simplicity. Generating a set of coherent forecasts involves only modelling and generating forecasts for the most aggregate top-level series. In general, such top-down approaches seem to produce quite reliable forecasts for the aggregate levels and they are useful with low count data. However, a significant disadvantage is the loss of information due to aggregation. A limitation of such top-down approaches is that characteristics

of lower level series cannot be captured. To overcome this, Athanasopoulos et al. (2009) introduced a new top-down approach which disaggregates the top-level based on proportions of forecasts rather than the historical data and showed that this method outperforms the conventional top-down approaches. However, a limitation of all top-down approaches is that they introduce bias to the forecasts even when the top-level forecast itself is unbiased. We discuss this in detail in Section 3.2.

3.1.3 Middle-out

A compromise between bottom-up and top-down approaches is the middle-out approach. It entails first forecasting the series of a selected middle-level. For series above the middle-level, coherent forecasts are generated using the bottom-up approach by aggregating the middle-level forecasts. For series below the middle level, coherent forecasts are generated using a top-down approach by disaggregating the middle-level forecasts. Since the middle-out approach involves generating top-down forecasts, it also introduces bias to the forecasts.

3.2 Point forecast reconciliation

All approaches discussed so far are limited to only using information from a single-level of aggregation. Furthermore, these ignore any correlations across levels of a hierarchy. An alternative framework that overcomes these limitations is one that involves forecast *reconciliation*. In a first step, forecasts for all the series across all levels of the hierarchy are computed, ignoring any aggregation constraints. We refer to these as *base* forecasts and denote them by $\hat{\mathbf{y}}_{T+h|T}$. In general, base forecasts will not be coherent, unless a very simple method has been used to compute them such as for naïve forecasts. In this case, forecasts are simply equal to a previous realisation of the data and they inherit the property of coherence.

The second step is an adjustment that reconciles base forecasts so that they become coherent. In general, this is achieved by mapping the base forecasts $\hat{\mathbf{y}}_{T+h|T}$ onto the coherent subspace \mathfrak{s} via a matrix \mathbf{SG} , resulting in a set of coherent forecasts $\tilde{\mathbf{y}}_{T+h|T}$. Specifically,

$$\tilde{\mathbf{y}}_{T+h|T} = \mathbf{SG}\hat{\mathbf{y}}_{T+h|T}, \quad (2)$$

where \mathbf{G} is an $m \times n$ matrix that maps $\hat{\mathbf{y}}_{T+h|T}$ to \mathbb{R}^m , producing new forecasts for the bottom-level, which are in turn mapped to the coherent subspace by the summing matrix \mathbf{S} . We restrict our attention to projections on \mathfrak{s} in which case $\mathbf{SGS} = \mathbf{S}$. This ensures that unbiasedness is preserved, i.e., for a set of unbiased base forecasts reconciled forecasts will also be unbiased.

Note that all single-level approaches discussed so far can also be represented by (2) using appropriately designed \mathbf{G} matrices, however not all of these will be

projections. For example for the bottom-up approach, $\mathbf{G} = (\mathbf{0}_{(m \times n-m)} \mathbf{I}_m)$ in which case $\mathbf{SGS} = \mathbf{S}$. For any top-down approach $\mathbf{G} = (\mathbf{p} \mathbf{0}_{(m \times n-1)})$, for which $\mathbf{SGS} \neq \mathbf{S}$.

3.2.1 Optimal MinT reconciliation

Wickramasuriya et al. (2018) build a unifying framework for much of the previous literature on forecast reconciliation. We present here a detailed outline of this approach and in turn relate it to previous significant contributions in forecast reconciliation.

Assume that $\hat{\mathbf{y}}_{T+h|T}$ is a set of unbiased base forecasts, i.e., $\mathbf{e}_{1:T}(\hat{\mathbf{y}}_{T+h|T}) = \mathbf{e}_{1:T}[\mathbf{y}_{T+h} \mid \mathbf{y}_1, \dots, \mathbf{y}_T]$, the true mean with the expectation taken over the observed sample up to time T . Let

$$\hat{\mathbf{e}}_{T+h|T} = \mathbf{y}_{T+h|T} - \hat{\mathbf{y}}_{T+h|T} \quad (3)$$

denote a set of base forecast errors with $\text{Var}(\hat{\mathbf{e}}_{T+h|T}) = \mathbf{W}_h$, and

$$\tilde{\mathbf{e}}_{T+h|T} = \mathbf{y}_{T+h|T} - \tilde{\mathbf{y}}_{T+h|T}$$

denote a set of coherent forecast errors. Lemma 1 in Wickramasuriya et al. (2018) shows that for any matrix \mathbf{G} such that $\mathbf{SGS} = \mathbf{S}$, $\text{Var}(\tilde{\mathbf{e}}_{T+h|T}) = \mathbf{SGW}_h\mathbf{S}'\mathbf{G}'$. Furthermore Theorem 1 shows that

$$\mathbf{G} = (\mathbf{S}'\mathbf{W}_h^{-1}\mathbf{S})^{-1}\mathbf{S}'\mathbf{W}_h^{-1} \quad (4)$$

is the unique solution that minimises the trace of $\mathbf{SGW}_h\mathbf{S}'\mathbf{G}'$ subject to $\mathbf{SGS} = \mathbf{S}$. MinT is optimal in the sense that given a set of unbiased base forecasts, it returns a set of best linear unbiased reconciled forecasts, using as \mathbf{G} the unique solution that minimises the trace (hence MinT) of the variance of the forecast error of the reconciled forecasts.

A significant advantage of the MinT reconciliation solution is that it is the first to incorporate the full correlation structure of the hierarchy via \mathbf{W}_h . However, estimating \mathbf{W}_h is challenging, especially for $h > 1$. Wickramasuriya et al. (2018) present possible alternative estimators for \mathbf{W}_h and show that these lead to different \mathbf{G} matrices. We summarise these below.

- Set $\mathbf{W}_h = k_h \mathbf{I}_n$ for all h , where $k_h > 0$ is a proportionality constant. This simple assumption returns $\mathbf{G} = (\mathbf{S}'\mathbf{S})^{-1}\mathbf{S}'$ so that the base forecasts are orthogonally projected onto the coherent subspace \mathfrak{s} minimising the Euclidean distance between $\hat{\mathbf{y}}_{T+h|T}$ and $\tilde{\mathbf{y}}_{T+h|T}$. Hyndman et al. (2011) come to same solution, however from the perspective of the following regression model

$$\hat{\mathbf{y}}_{T+h|T} = \mathbf{S}\boldsymbol{\beta}_{T+h|T} + \boldsymbol{\varepsilon}_{T+h|T},$$

where $\boldsymbol{\beta}_{T+h|T} = \mathbf{e}[\mathbf{b}_{T+h} \mid \mathbf{b}_1, \dots, \mathbf{b}_T]$ is the unknown conditional mean of the bottom-level series and $\boldsymbol{\varepsilon}_{T+h|T}$ is the coherence or reconciliation error with mean zero and variance \mathbf{V} . The OLS solution leads to the same projection matrix

$S(S'S)^{-1}S'$, and due to this interpretation we continue to refer to this reconciliation method as OLS. A disadvantage of the OLS solution is that the homoscedastic diagonal entries do not account for the scale differences between the levels of the hierarchy due to aggregation. Furthermore, OLS does not account for the correlations across series.

- Set $\mathbf{W}_h = k_h \text{diag}(\hat{\mathbf{W}}_1)$ for all h ($k_h > 0$), where

$$\hat{\mathbf{W}}_1 = \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{e}}_t \hat{\mathbf{e}}_t'$$

is the unbiased sample estimator of the in-sample one-step-ahead base forecast errors as defined in (3). Hence this estimator scales the base forecasts using the variance of the in-sample residuals and is therefore described and referred to as a weighted least squares (WLS) estimator applying variance scaling. A similar estimator was proposed by Hyndman et al. (2016).

An alternative WLS estimator is proposed by Athanasopoulos et al. (2017) in the context of temporal hierarchies. Here \mathbf{W}_h is proportional to $\text{diag}(\mathbf{S}\mathbf{1})$ where $\mathbf{1}$ is a unit column vector of dimension n . Hence the weights are proportional to the number of bottom-level variables required to form an aggregate. For example, in the hierarchy of Figure 1, the weights corresponding to the Total, series A and series B are proportional to 5, 2 and 3 respectively. This weighting scheme depends only on the aggregation structure and is referred to as structural scaling. Its advantage over OLS is that it assumes equivariant forecast errors only at the bottom-level of the structure and not across all levels. It is particularly useful in cases where forecast errors are not available; for example, in cases where the base forecasts are generated by judgemental forecasting.

- Set $\mathbf{W}_h = k_h \hat{\mathbf{W}}_1$ for all h ($k_h > 0$) to be proportional to the unrestricted sample covariance estimator for $h = 1$. Although this is relatively simple to obtain and provides a good solution for small hierarchies, it does not provide reliable results as m grows compared to T . This is referred to this as the MinT(Sample) estimator.
- Set $\mathbf{W}_h = k_h \hat{\mathbf{W}}_1^D$ for all h ($k_h > 0$), where $\hat{\mathbf{W}}_1^D = \lambda_D \text{diag}(\hat{\mathbf{W}}_1) + (1 - \lambda_D) \hat{\mathbf{W}}_1$ is a shrinkage estimator with diagonal target and shrinkage intensity parameter

$$\hat{\lambda}_D = \frac{\sum_{i \neq j} \hat{\text{Var}}(\hat{r}_{ij})}{\sum_{i \neq j} \hat{r}_{ij}^2},$$

where \hat{r}_{ij} is the (i, j) th element of $\hat{\mathbf{R}}_1$, the one-step-ahead sample correlation matrix as proposed by Schäfer & Strimmer (2005). Hence, off-diagonal elements of $\hat{\mathbf{W}}_1$ are shrunk towards zero while diagonal elements (variances) remain unchanged. This is referred to as the MinT(Shrink) estimator.

4 Hierarchical probabilistic forecasting

A limitation of point forecasts is that they provide no indication of uncertainty around the forecast. A richer description of forecast uncertainty can be obtained by providing a probabilistic forecast, also commonly referred to as a density forecast. For a review of probabilistic forecasts, and methods for evaluating such forecasts known as *scoring rules* see (Gneiting & Katzfuss 2014). In recent years, the use of probabilistic forecasts and their evaluation via scoring rules has become pervasive in macroeconomic forecasting, some notable (but non-exhaustive) examples are Geweke & Amisano (2010), Billio et al. (2013), Carriero et al. (2015) and Clark & Ravazzolo (2015).

The literature on hierarchical probabilistic forecasting is still an emerging area of interest. To the best of our knowledge the first attempt to even define coherence in the setting of probabilistic forecasting is provided by Taieb et al. (2017) who define a coherent forecast in terms of a convolution. An equivalent definition, provided by Gamakumara et al. (2018) defines a coherent probabilistic forecast as a probability measure on the coherent subspace \mathfrak{s} . Gamakumara et al. (2018) also generalise the concept of forecast reconciliation to the probabilistic setting.

Definition 2. Let \mathcal{A} be a subset¹ of \mathfrak{s} and let \mathcal{B} be all points in \mathbb{R}^n that are mapped onto \mathcal{A} after premultiplication by \mathbf{SG} . Letting $\hat{\nu}$ be a base probabilistic forecast for the full hierarchy, the coherent measure $\tilde{\nu}$ reconciles $\hat{\nu}$ if $\tilde{\nu}(\mathcal{A}) = \hat{\nu}(\mathcal{B})$ for all \mathcal{A} .

In practice this definition leads to two approaches. For some parametric distributions, for instance the multivariate normal, a reconciled probabilistic forecast can be derived analytically. However, in macroeconomic forecasting, non-standard distributions such as bimodal distributions are often required to take different policy regimes into account. In such cases a non-parametric approach based on bootstrapping in-sample errors proposed Gamakumara et al. (2018) can be used. These scenarios are now covered in detail.

4.1 Probabilistic forecast reconciliation in the Gaussian framework

In the case where the base forecasts are probabilistic forecasts characterised by elliptical distributions, Gamakumara et al. (2018) show that reconciled probabilistic forecasts will also be elliptical. This is particularly straightforward for the Gaussian distribution which is completely characterised by two moments. Letting the base probabilistic forecasts be $\mathcal{N}(\hat{\mathbf{y}}_{T+h|T}, \hat{\Sigma}_{T+h|T})$, then the reconciled probabilistic forecasts will be $\mathcal{N}(\tilde{\mathbf{y}}_{T+h|T}, \tilde{\Sigma}_{T+h|T})$, where

$$\tilde{\mathbf{y}}_{T+h|T} = \mathbf{SG}\hat{\mathbf{y}}_{T+h|T}, \quad (5)$$

¹ Strictly speaking \mathcal{A} is a Borel set

and

$$\tilde{\Sigma}_{T+h|T} = SG\hat{\Sigma}_{T+h|T}G'S'. \quad (6)$$

There are several options for obtaining the base probabilistic forecasts and in particular the variance covariance matrix $\hat{\Sigma}$. One option is to fit multivariate models either level by level or for the hierarchy as a whole leading respectively to a $\hat{\Sigma}$ that is block diagonal or dense. Another alternative is to fit univariate models for each individual series in which case $\hat{\Sigma}$ is a diagonal matrix. A third alternative, that we employ here, is to obtain $\hat{\Sigma}$ using in-sample forecast errors, in a similar vein to how \hat{W}_1 is estimated in the MinT method. Here the same shrinkage estimator described in Section 3.2 is used. The reconciled probabilistic forecast will ultimately depend on the choice of G ; the same choices of G matrices used in Section 3 can be used.

4.2 Probabilistic forecast reconciliation in the non-parametric framework

In many applications, including macroeconomic forecasting, it may not be reasonable to assume Gaussian predictive distributions. Therefore, non-parametric approaches have been widely used for probabilistic forecasts in different disciplines. For example, ensemble forecasting in weather applications (Gneiting 2005, Gneiting & Katzfuss 2014, Gneiting et al. 2008), and bootstrap based approaches (Manzan & Zerom 2008, Vilar & Vilar 2013). In macroeconomics, Cogley et al. (2005) discuss the importance of allowing for skewness in density forecasts and more recently Smith & Vahey (2016) discuss this issue in detail.

Due to these concerns, we employ the bootstrap method proposed by Gamakumara et al. (2018) that does not make parametric assumptions about the predictive distribution. An important result exploited by this method is that applying point forecast reconciliation to the draws from an incoherent base predictive distribution, results in a sample from the reconciled predictive distribution. We summarise this process below:

1. Fit univariate models to each series in the hierarchy over a training set from $t = 1, \dots, T$.
2. For each series compute h -step-ahead point forecasts, for $h = 1, \dots, H$. Collect these into an $n \times H$ matrix $\hat{Y} := (\hat{y}_{T+1|T}, \dots, \hat{y}_{T+H|T})$, where $\hat{y}_{T+h|T}$ is an n -vector of h -step-ahead point forecasts for all series in the hierarchy.
3. Compute one-step-ahead in-sample forecast errors. Collect these into an $n \times T$ matrix $\hat{E} = (\hat{e}_1, \hat{e}_2, \dots, \hat{e}_T)$, where the n -vector $\hat{e}_t = y_t - \hat{y}_{t|t-1}$. Here, $\hat{y}_{t|t-1}$ is a vector of forecasts made for time t using information up to and including time $t - 1$. These are called in-sample forecasts since they depend on past values but information from the entire training sample is used to estimate parameters that forecasts are based on.

4. Block bootstrap from $\hat{\mathbf{E}}$, that is choose H consecutive columns of $\hat{\mathbf{E}}$ at random, repeating this process B times. Denote the $n \times H$ matrix obtained at iteration b as $\hat{\mathbf{E}}^b$ for $b = 1, \dots, B$.
5. For all b , compute $\hat{\mathbf{Y}}^b := \hat{\mathbf{Y}} + \hat{\mathbf{E}}^b$. Each row of $\hat{\mathbf{Y}}^b$ is a sample path of h forecasts for a single series. Each column of $\hat{\mathbf{Y}}^b$ is a realisation from the joint predictive distribution at a particular horizon.
6. For each $b = 1, \dots, B$, select the h th column of $\hat{\mathbf{Y}}^b$ and stack these to form an $n \times B$ matrix $\hat{\mathbf{Y}}_{T+h|T}$.
7. For a given \mathbf{G} matrix and for each $h = 1, \dots, H$, compute $\tilde{\mathbf{Y}}_{T+h|T} = \mathbf{S}\mathbf{G}\hat{\mathbf{Y}}_{T+h|T}$. Each column of $\tilde{\mathbf{Y}}_{T+h|T}$ is a realisation from the joint h -step-ahead reconciled predictive distribution.

5 Australian GDP

In our empirical application we consider Gross Domestic Product (GDP) of Australia with quarterly data spanning the period 1984:Q4–2018:Q3. The Australian Bureau of Statistics (ABS) measures GDP using three main approaches namely Production, Income and Expenditure. The final GDP figure is obtained as an average of these three figures. Each of these measures are aggregates of economic variables which are also targets of interests for the macroeconomic forecaster. This suggests a hierarchical approach to forecasting could be used to improve forecasts of all series in the hierarchy including the headline GDP.

We concentrate on the Income and Expenditure approaches as nominal data are available only for these two. We restrict our attention to nominal data due to the fact that real data are constructed via a chain price index approach with different price deflators used for each series. As a result, real GDP data are not coherent — the aggregate series is not a linear combination of the disaggregate series. For similar reasons we do not use seasonally adjusted data; the process of seasonal adjustment results in data that are not coherent. Finally, although there is a small statistical discrepancy between each series and the headline GDP figure, we simply treat this statistical discrepancy, which is also published by the ABS, as a time series in its own right. For further of the details on the data please refer to Australian Bureau of Statistics (2018).

Income approach

Using the income approach, GDP is calculated by aggregating all income flows. In particular, GDP at purchaser's price is the sum of all factor incomes and taxes, minus subsidies on production and imports (?):

$$\begin{aligned}
 GDP = & \text{Gross operating surplus} + \text{Gross mixed income} \\
 & + \text{Compensation of employees} \\
 & + \text{Taxes less Subsidies on production and imports} \\
 & + \text{Statistical discrepancy (I)}.
 \end{aligned}$$

Figure 4 shows the full hierarchical structure capturing all components aggregated to form GDP using the income approach. The hierarchy has two levels of aggregation below the top-level, with a total of $n = 16$ series across the whole structure and $m = 10$ series at the bottom-level.

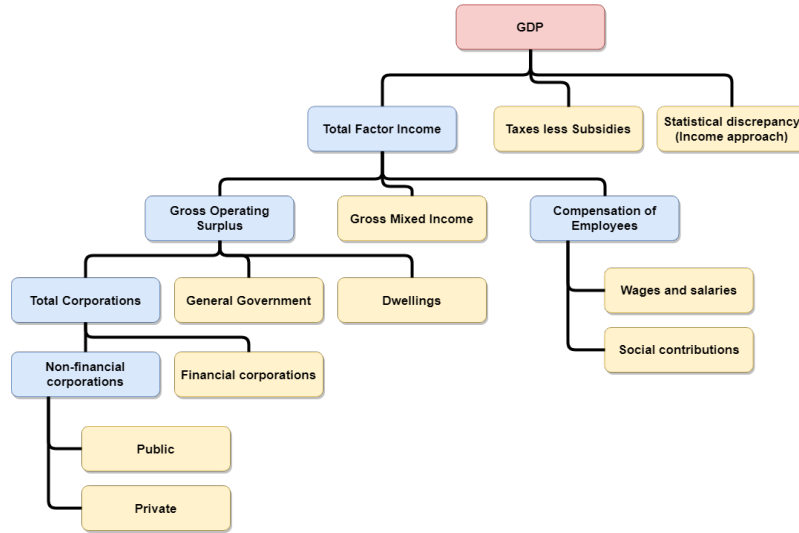


Fig. 4 Hierarchical structure of the income approach for GDP. The pink cell contains GDP the most aggregate series. The blue cells contain intermediate-level series and the yellow cells to the most disaggregate bottom-level series.

Expenditure approach

In the expenditure approach, GDP is calculated as the aggregation of final consumption expenditure, gross fixed capital formation (GFCF), changes in inventories of finished goods, work-in-progress and raw materials and the value of exports less imports of the goods and services (?). The underline equation is:

$$\begin{aligned}
 GDP = & \text{Final consumption expenditure} + \text{Gross fixed capital formation} \\
 & + \text{Changes in inventories} + \text{Trade balance} + \text{Statistical discrepancy (E)}.
 \end{aligned}$$

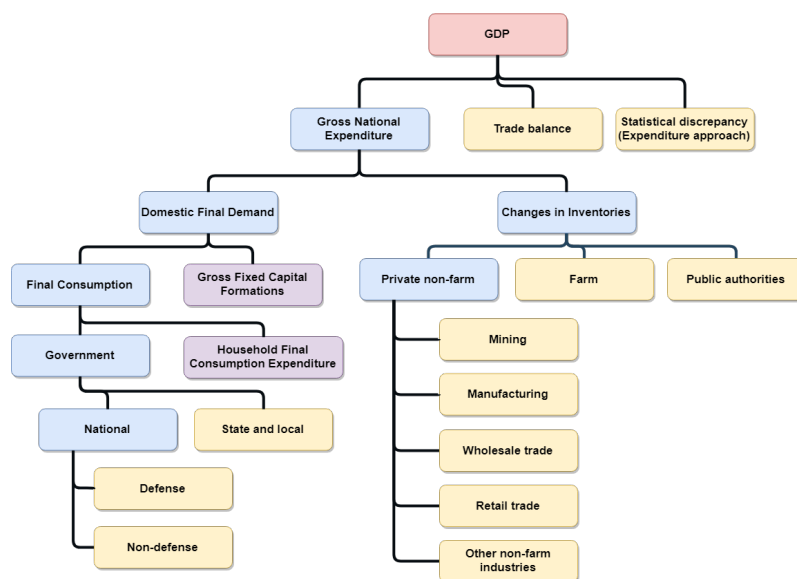


Fig. 5 Hierarchical structure of the expenditure approach for GDP. The pink cell contains GDP, the most aggregate series. The blue and purple cells contain intermediate-level series with the series in the purple cells further disaggregated in Figures 6 and 7. The yellow cells contain the most disaggregate bottom-level series.

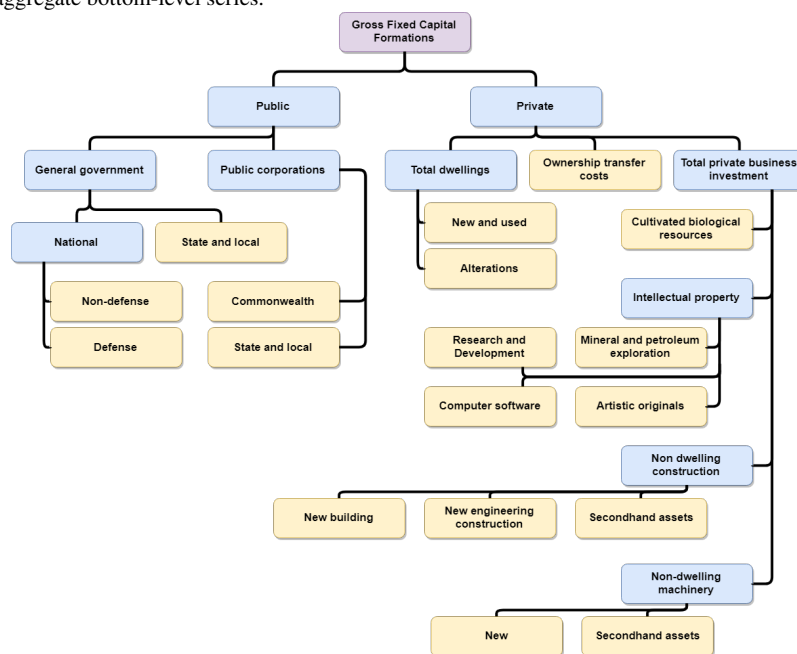


Fig. 6 Hierarchical structure for Gross Fixed Capital Formations under the expenditure approach for GDP, continued from Figure 5. Blue cells contain intermediate-level series and the yellow cells to the most disaggregate bottom-level series.

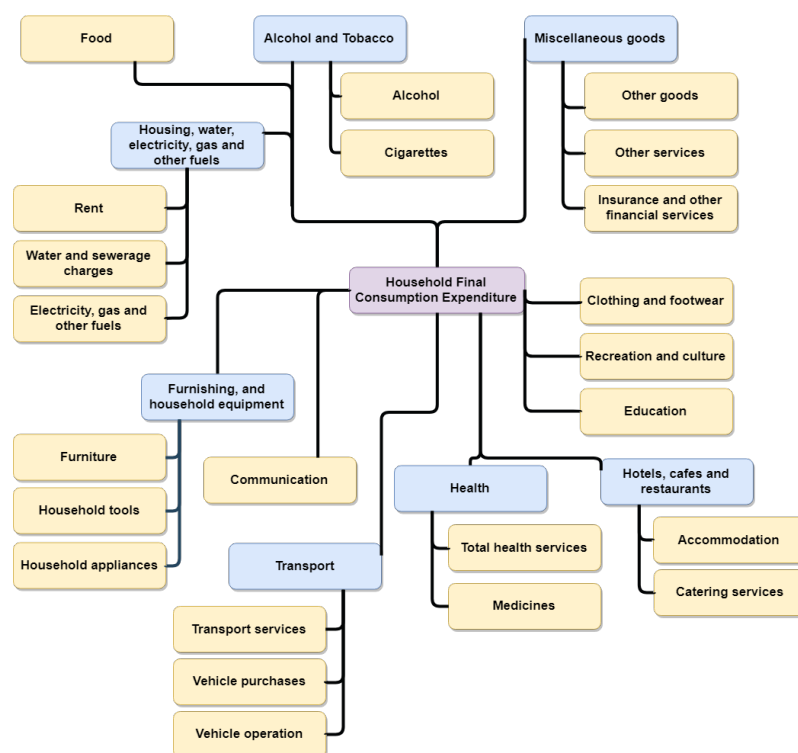


Fig. 7 Hierarchical structure for Household Final Consumption Expenditure under expenditure approach for GDP, continued from Figure 5. Blue cells contain intermediate-level series and the yellow cells to the most disaggregate bottom-level series..

Figures 5, 6 and 7 show the full hierarchical structure capturing all components aggregated to form GDP using the expenditure approach. The hierarchy has three levels of aggregation below the top-level, with a total of $n = 80$ series across the whole structure and $m = 53$ series at the bottom-level. Descriptions of each series in these hierarchies along with the series ID assigned by the ABS are given in the Tables 1, 2, 3 and 4 in the Appendix.

Figure 8 displays time series from the income and expenditure approaches. The top panel shows the most aggregate GDP series. The panels below show series from levels below for the income hierarchy (left panel) and the expenditure hierarchy (right panel). The plots show the diverse features of the time series with some displaying positive and others negative trending behaviour, some showing no trends but possibly a cycle, and some having a strong seasonal component. These highlight the need to account and model all information and diverse signals from each series in the hierarchy which can only be achieved through a forecast reconciliation approach.

GEORGE: I have added these to the Appendix. Will ask the Editors how they want these and that they can go to online supplementary materials.

Please save all figures as pdf not png. Also fix facet labels for Fig 8 and convert millions to billions.

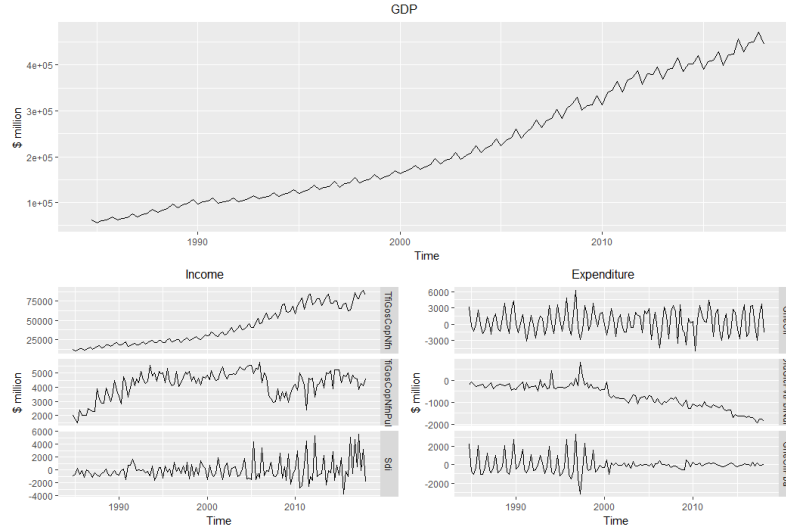


Fig. 8 Time plots for series from different levels of income and expenditure hierarchies.

6 Empirical application methodology

We now demonstrate the potential for reconciliation methods to improve forecast accuracy for Australian GDP. We consider forecasts from $h = 1$ quarter ahead up to $h = 4$ quarters ahead using an *expanding* window. First, the training sample is set from 1984:Q4 to 1994:Q3 and forecasts are produced for 1994:Q4 to 1995:Q3. Then the training window is expanded by one quarter at a time, i.e. from 1984:Q4 to 2017:Q4 with the final forecasts produced for the last available observation in 2018:Q1. This leads to 94 1-step-ahead, 93 2-steps-ahead, 92 3-steps-ahead and 91 4-steps-ahead forecasts available for evaluation.

6.1 Models

The first task in forecast reconciliation is to obtain base forecasts for all series in the hierarchy. In the case of the income approach this necessitates forecasting $n = 16$ separate time series while in the case of the expenditure approach forecasts for $n = 80$ separate time series must be obtained. Given the diversity in these time series discussed in Section 5, we focus on an approach that is fast but also flexible. We consider simple univariate ARIMA models, where model order is selected via a combination of unit root testing and the AIC using an algorithm developed by Hyndman & Khandakar (2008) and implemented in the `auto.arima()` function in Hyndman et al. (2019). A similar approach was also undertaken using the ETS framework to produce base forecasts (Hyndman et al. 2008). This had minimal impact

on our conclusions with respect to forecast reconciliation methods, and in most cases ARIMA forecasts were found to be more accurate than ETS forecasts. Consequently for brevity, we have excluded presenting the results for ETS models. However, these are available from [github²](#). We note that a number of more complicated approaches could have been used to obtain base forecasts including multivariate models such as vector autoregressions, and models and methods that handle a large number of predictors such as factor models or least angle regression. However, Panagiotelis et al. (2019) show that univariate ARIMA models are highly competitive for forecasting Australian GDP even compared to these methods, and in any case our primary motivation is to demonstrate the potential of forecast reconciliation.

Need to put them there in a clearly marked folder.

The hierarchical forecasting approaches we consider are bottom-up, OLS, WLS with variance scaling and the MinT(Shrink) approach. The MinT(Sample) approach was also used but due to the size of the hierarchy, forecasts reconciled via this approach were less stable. Finally, all forecasts (both base and coherent) are compared to a seasonal naïve benchmark (Hyndman & Athanasopoulos 2018); i.e. the forecast for GDP (or one of its components) is the realised GDP in the same quarter of the previous year. The naïve forecasts are by construction coherent and therefore do not need to be reconciled.

6.2 Evaluation

For evaluating point forecasts we consider two metrics, the Mean Squared Error (MSE) and the Mean Absolute Scaled Error (MASE). The absolute scaled error is defined as

$$q_{T+h} = \sum \frac{|\check{e}_{T+h|T}|}{(T-4)^{-1} \sum_{t=5}^T |y_t - y_{t-4}|},$$

where \check{e}_{t+h} is the difference between any forecast and the realisation³, and 4 is used due to the quarterly nature of the data. An advantage of using MASE is that it is a scale independent measure. This is particularly relevant for hierarchical time series, since aggregate series by their very nature are on a larger scale than disaggregate series. Consequently, scale dependent metrics may unfairly favour methods that perform well for the aggregate series but poorly for disaggregate series. For more details on different point forecast accuracy measures, refer to Chapter 3 of Hyndman & Athanasopoulos (2018).

Forecast accuracy of probabilistic forecasts can be evaluated using scoring rules (Gneiting & Katzfuss 2014). Let \check{F} be a probabilistic forecast and let $\check{y} \sim \check{F}$ where a breve is again used to denote that either base forecasts or coherent forecasts can be evaluated. The accuracy of multivariate probabilistic forecasts will be measured by the energy score given by

² URL???

³ Breve is used instead of a hat or tilde to denote that this can be the error for either a base or reconciled forecast.

$$eS(\check{F}_{T+h|T}, \mathbf{y}_{T+h}) = e_{\check{F}} \|\check{\mathbf{y}}_{T+h} - \mathbf{y}_{T+h}\|^\alpha - \frac{1}{2} e_{\check{F}} \|\check{\mathbf{y}}_{T+h} - \check{\mathbf{y}}_{T+h}^*\|^\alpha,$$

where \mathbf{y}_{T+h} is the realisation at time $T + h$, and $\alpha \in (0, 2]$. We set $\alpha = 1$, noting that other values of α give similar results. The expectations can be evaluated numerically as long as a sample from \check{F} is available, which is the case for all methods we employ. An advantage of using energy scores is that in the univariate case it simplifies to the commonly used cumulative rank probability score (CRPS) given by

$$\text{CRPS}(\check{F}_i, y_{i,T+h}) = e_{\check{F}_i} |\check{y}_{i,T+h} - y_{i,T+h}| - \frac{1}{2} e_{\check{F}_i} |\check{y}_{i,T+h} - \check{y}_{i,T+h}^*|,$$

where the subscript i is used to denote that CRPS measures forecast accuracy for a single variable in the hierarchy.

Alternatives to the energy score were also considered, namely log scores and variogram scores. The log score was disregarded since Gamakumara et al. (2018) prove that the log score is improper with respect to the class of incoherent probabilistic forecasts when the true DGP is coherent. The variogram score gave similar results to the energy score, but these results are omitted for brevity but these are available from [github](#)⁴

Need to set this up.

7 Results

7.1 Base forecasts

Due to the different features in each time series a variety of ARIMA and seasonal ARIMA models were selected for generating base forecasts. For example, in the income hierarchy, some series require seasonal differencing while other did not. Furthermore the AR orders vary from 0–3, the MA orders from 0–2, and their seasonal counterparts SAR from 0–2 and SMA from 0–1. Figure 9 compares the accuracy of the ARIMA base forecasts to the seasonal naïve forecasts over different forecast horizons. The panels on the left show results for the Income hierarchy while the panels on the right show the results for the Expenditure hierarchy. The top panels summarise results over all series in the hierarchy, i.e. we calculate the MSE for each series and then average over all series. The bottom panels show the results for the aggregate level GDP.

The clear result is that base forecasts are more accurate than the naïve forecasts, however as the forecasting horizon increases, the differences become smaller. This is to be expected since the naïve model here is a seasonal random walk, and for horizons $h < 4$ forecasts from an ARIMA model are based on more recent information. Similar results are obtained when MASE is used as the metric for evaluating forecast accuracy.

⁴ URL

One disadvantage of the base forecasts relative to the naïve forecasts is that base forecasts are not coherent. As such we now turn our attention to investigating whether reconciliation approaches can lead to further improvements in forecast accuracy relative to the base forecasts.

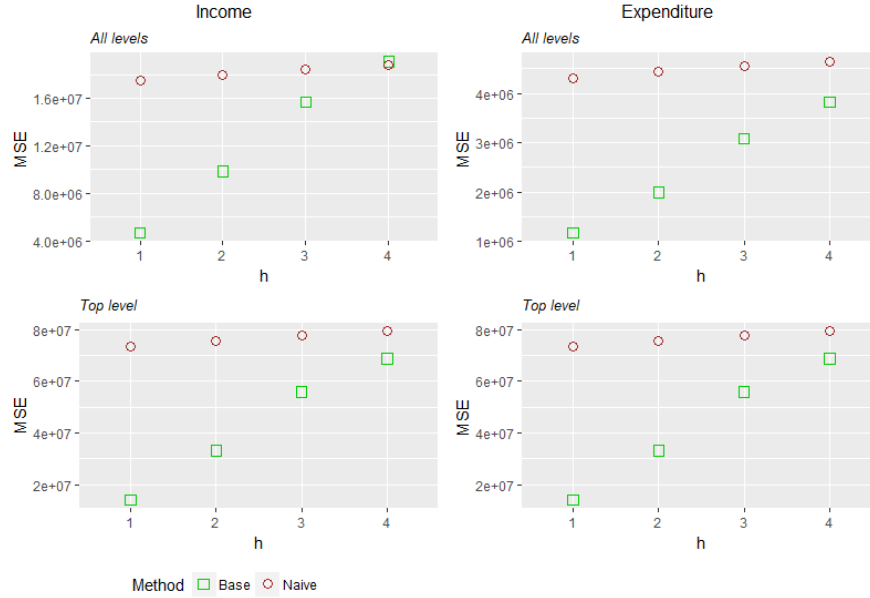


Fig. 9 Mean squared errors for naïve and ARIMA base forecasts. Top panels refer to results summarised over all series while bottom panels refer to results for the top-level GDP series. Left panels refer to the income hierarchy and right panels to the expenditure hierarchy.

7.2 Point Forecast Reconciliation

We now turn our attention to evaluating the accuracy of point forecasts obtained using the different reconciliation approaches as well as the single-level bottom-up approach. All results in subsequent figures are presented as the percentage changes in a forecasting metric relative to base forecasts, a measure known in the forecasting literature as *skill scores*. Skill scores are computed such that positive values represent an improvement in forecasting accuracy over the base forecasts while negative values represent a deterioration.

Figures 10 and 11 show skill scores using MSE and MASE respectively. The top row of each figure shows skill scores based on averages over all series. We conclude that reconciliation methods generally improve forecast accuracy relative to base forecasts regardless of the hierarchy used, the forecasting horizon, the forecast

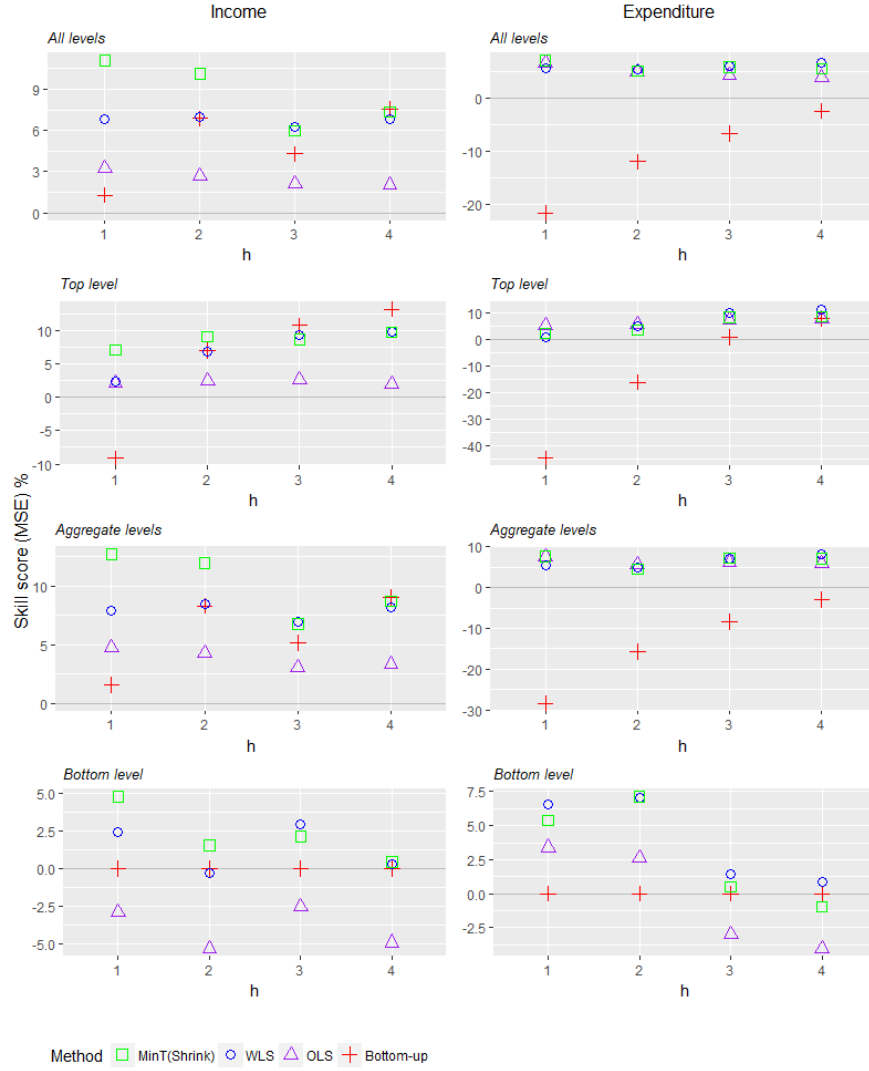


Fig. 10 Skill scores for point forecasts from alternative methods (with reference to base forecasts) using MSE. Left panels refer to the income hierarchy while the right panels refer to the expenditure hierarchy. The first row refers to results summarised over all series, second row to top-level GDP series, third row to aggregate levels and last row to the bottom-level.

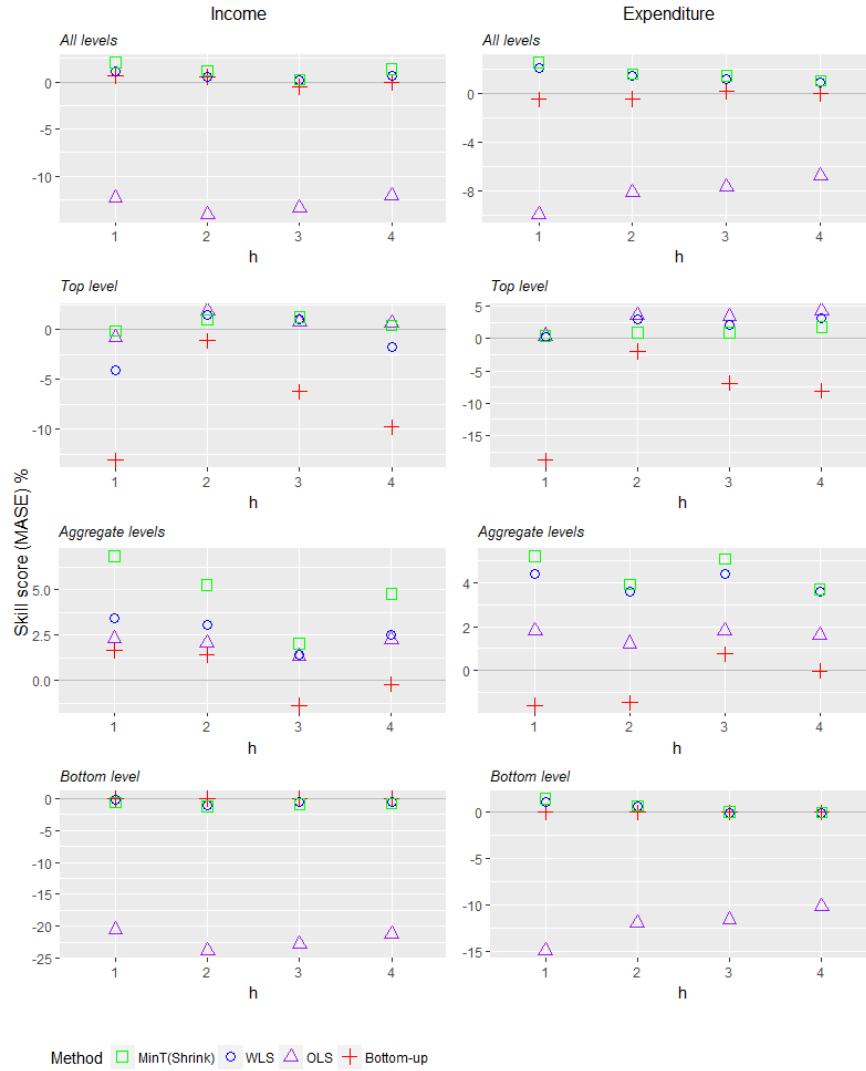


Fig. 11 Skill scores for point forecasts from different reconciliation methods (with reference to base forecasts) using MASE. Left two panels refer to income hierarchy and right two panels to expenditure hierarchy. First row refers to results summarised over all series, second row to top-level GDP series, third row to aggregate levels and last row to the bottom-level.

error measure or the reconciliation method employed. We do however note that while all reconciliation methods improve forecast performance, MinT(Shrink) is the best forecasting method in most cases.

To further investigate the results we break down the skill scores by different levels of each hierarchy. The second row of Figures 10 and 11 shows the skill scores for a single series, namely GDP which represents the top-level of both hierarchies. The third row shows results for all series excluding those of the bottom-level, while the final row shows results for the bottom-level series only. Here, we see two general features. The first is that OLS reconciliation performs poorly on the bottom-level series, and the second is that bottom-up performs relatively poorly on aggregate series. The two features are particularly exacerbated for the larger expenditure hierarchy. These results are consistent with other findings in the forecast reconciliation literature (see for instance Athanasopoulos et al. 2017, Wickramasuriya et al. 2018)

7.3 Probabilistic Forecast Reconciliation

We now turn our attention towards results for probabilistic forecasts. Figure 12 shows results for the energy score which as a multivariate score summarises forecast ac-

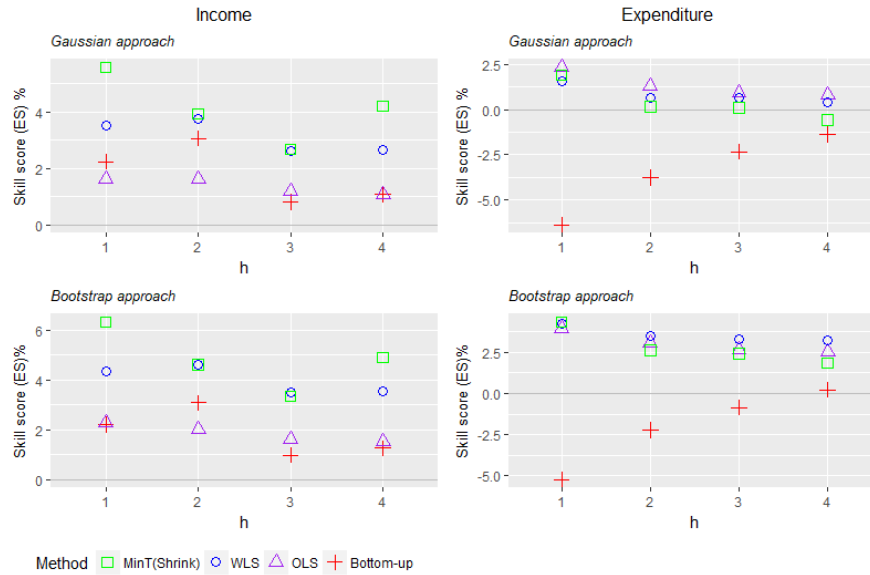


Fig. 12 Skill scores for multivariate probabilistic forecasts from different reconciliation methods (with reference to base forecasts) using energy score. Top panels refer to the results for Gaussian approach and bottom panels to the non-parametric bootstrap approach. Left panels refer to the income hierarchy and right panels to the expenditure hierarchy.

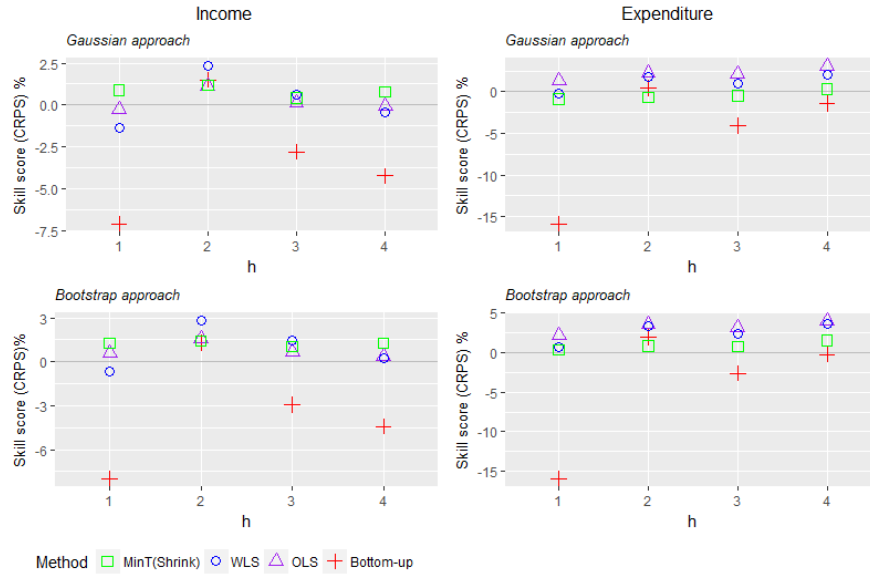


Fig. 13 Skill scores for probabilistic forecasts of top-level GDP from different reconciliation methods (with reference to base forecasts) using CRPS. Top panels refer to the results for Gaussian approach and bottom panels refer to the non-parametric bootstrap approach. Left panel refers to the income hierarchy and right panel to the expenditure hierarchy.

curacy over the entire hierarchy. Once again all results are presented as skill scores relative to base forecasts. The top panels refer to results assuming Gaussian probabilistic forecasts as described in Section 4.1 while the bottom panels refer to the non-parametric bootstrap method described in Section 4.2. The left panels correspond to the income hierarchy while the right panels correspond to the expenditure hierarchy. For the income hierarchy, all methods improve upon base forecasts at all horizons. In nearly all cases the best performing reconciliation method is MinT(Shrink), a notable result since the optimal properties for MinT have thus far only been established theoretically in the point forecasting case. For the larger expenditure hierarchy results are little more mixed. While bottom-up tends to perform poorly, all reconciliation methods improve upon base forecasts (with the single exception of MinT(Shrink) in the Gaussian framework four quarters ahead). Interestingly, OLS performs best under the assumption of Gaussianity - this may indicate that OLS is a more robust method under model misspecification but further investigation is required.

Finally, Figure 13 displays the skill scores based on the cumulative ranked probability score for a single series, namely top-level GDP. The cause of the poor performance of bottom-up reconciliation as a failure to accurately forecast aggregate series is apparent here.

8 Conclusions

In the macroeconomic setting, we have demonstrated the potential for forecast reconciliation methods to not only provide coherent forecasts, but to also improve overall forecast accuracy. This result holds for both point forecasts and probabilistic forecasts, for the two different hierarchies we consider and over different forecasting horizons. Even where the objective is to only forecast a single series, for instance top-level GDP, the application of forecast reconciliation methods improves forecast accuracy.

By comparing results from different forecast reconciliation techniques we draw a number of conclusions. Despite its simplicity, the single-level bottom-up approach can perform poorly at more aggregated levels of the hierarchy. Meanwhile, when forecast accuracy at the bottom-level is evaluated, OLS tends to break down in some instances. Overall, the WLS and MinT(Shrink) methods, and particularly the latter tend to yield the highest improvements in forecast accuracy. Similar results can be found in both simulations and the empirical studies of Athanasopoulos et al. (2017) and Wickramasuriya et al. (2018).

There are a number of open avenues for research in the literature on forecast reconciliation, some of which are particularly relevant to macroeconomic applications. First there is scope to consider more complex aggregation structures, for instance in addition to the hierarchies we have already considered, data on GDP and GDP components disaggregated along geographical lines are also available. This leads to a grouped aggregation structure. Also, given the substantial literature on the optimal frequency at which to analyse macroeconomic data, a study on forecasting GDP or other variables as a temporal hierarchy may be of interest. In this chapter we have only shown that reconciliation methods can be used to improve forecast accuracy when univariate ARIMA models are used to produce base forecasts. It will be interesting to evaluate whether such results hold when a multivariate approach, e.g. a Bayesian VAR or dynamic factor model, is used to generate base forecasts, or whether the gains from forecast reconciliation would be more modest. Finally, a current limitation of the forecast reconciliation literature is that it only applies to collections of time series that adhere to linear constraints. In macroeconomics there are many examples of data that adhere to non-linear constraints, for instance real GDP is a complicated but deterministic function of GDP components and price deflators. The extension of forecast reconciliation methods to non-linear constraints potentially holds great promise for continued improvement in macroeconomic forecasting.

References

- Amemiya, T. & Wu, R. Y. (1972), 'The effect of aggregation on prediction in the autoregressive model', *Journal of the American Statistical Association* **67**(339), 628–632.
- Andrawis, R. R., Atiya, A. F. & El-Shishiny, H. (2011), 'Combination of long term and short term forecasts, with application to tourism demand forecasting', *International Journal of Forecasting* **27**(3), 870–886.
- Athanasopoulos, G., Ahmed, R. A. & Hyndman, R. J. (2009), 'Hierarchical forecasts for Australian domestic tourism', *International Journal of Forecasting* **25**(1), 146–166.
- Athanasopoulos, G., Hyndman, R. J., Kourentzes, N. & Petropoulos, F. (2017), 'Forecasting with temporal hierarchies', *European Journal of Operational Research* **262**, 60–74.
- Athanasopoulos, G., Steel, T. & Weatherburn, D. (2019), Forecasting prison numbers: A grouped time series approach.
- Australian Bureau of Statistics (2018), Australian National Accounts: National income, expenditure and product, Technical report.
- Billio, M., Casarin, R., Ravazzolo, F. & Van Dijk, H. K. (2013), 'Time-varying combinations of predictive densities using nonlinear filtering', *Journal of Econometrics* **177**(2), 213–232.
- Brewer, K. (1973), 'Some consequences of temporal aggregation and systematic sampling for ARMA and ARMAX models', *Journal of Econometrics* **1**(2), 133–154.
- Capistrán, C., Constandse, C. & Ramos-Francia, M. (2010), 'Multi-horizon inflation forecasts using disaggregated data', *Economic Modelling* **27**(3), 666–677.
- Carriero, A., Clark, T. E. & Marcellino, M. (2015), 'Realtime nowcasting with a Bayesian mixed frequency model with stochastic volatility', *Journal of the Royal Statistical Society: Series A (Statistics in Society)* **178**(4), 837–862.
- Clark, T. E. & Ravazzolo, F. (2015), 'Macroeconomic forecasting performance under alternative specifications of time-varying volatility', *Journal of Applied Econometrics* **30**(4), 551–575.
- Cogley, T., Morozov, S. & Sargent, T. J. (2005), 'Bayesian fan charts for UK inflation: Forecasting and sources of uncertainty in an evolving monetary system', *Journal of Economic Dynamics and Control* **29**(11), 1893–1925.
- Dunn, D. M., Williams, W. H. & Dechaine, T. L. (1976), 'Aggregate versus subaggregate models in local area forecasting', *Journal of American Statistical Association* **71**(353), 68–71.
- Gamakumara, P., Panagiotelis, A., Athanasopoulos, G. & Hyndman, R. J. (2018), Probabilistic forecasts in hierarchical time series, Working paper 11/18, Monash University Econometrics & Business Statistics.
- Geweke, J. & Amisano, G. (2010), 'Comparing and evaluating Bayesian predictive distributions of asset returns', *International Journal of Forecasting* **26**(2), 216–230.

- Gneiting, T. (2005), 'Weather forecasting with ensemble methods', *Science* **310**(5746), 248–249.
- Gneiting, T. & Katzfuss, M. (2014), 'Probabilistic forecasting', *Annual Review of Statistics and Its Application* **1**, 125–151.
- Gneiting, T., Stanberry, L. I., Gneiting, E. P., Held, L. & Johnson, N. A. (2008), 'Assessing probabilistic forecasts of multivariate quantities, with an application to ensemble predictions of surface winds', *Test* **17**(2), 211–235.
- Gross, C. W. & Sohl, J. E. (1990), 'Disaggregation methods to expedite product line forecasting', *Journal of Forecasting* **9**(3), 233–254.
- Hotta, L. K. (1993), 'The effect of additive outliers on the estimates from aggregated and disaggregated ARIMA models', *International Journal of Forecasting* **9**(1), 85–93.
- Hotta, L. K. & Cardoso Neto, J. (1993), 'The effect of aggregation on prediction in autoregressive integrated moving-average models', *Journal of Time Series Analysis* **14**(3), 261–269.
- Hyndman, R. J., Ahmed, R. A., Athanasopoulos, G. & Shang, H. L. (2011), 'Optimal combination forecasts for hierarchical time series', *Computational Statistics and Data Analysis* **55**(9), 2579–2589.
- Hyndman, R. J. & Athanasopoulos, G. (2018), *Forecasting: Principles and Practice*, OTexts.
URL: <https://OTexts.com/fpp2>
- Hyndman, R. J., Athanasopoulos, G., Bergmeir, C., Caceres, G., Chhay, L., O'Hara-Wild, M., Petropoulos, F., Razbash, S., Wang, E., Yasmien, F., R Core Team, Ihaka, R., Reid, D., Shaub, D., Tang, Y. & Zhou, Z. (2019), *forecast: Forecasting Functions for Time Series and Linear Models*. Version 8.5.
URL: <https://CRAN.R-project.org/package=forecast>
- Hyndman, R. J. & Khandakar, Y. (2008), 'Automatic time series forecasting: the forecast package for R', *Journal of Statistical Software* **26**(3), 1–22.
- Hyndman, R. J., Koehler, A. B., Ord, J. K. & Snyder, R. D. (2008), *Forecasting with exponential smoothing: the state space approach*, Springer-Verlag, Berlin.
- Hyndman, R. J., Lee, A. J. & Wang, E. (2016), 'Fast computation of reconciled forecasts for hierarchical and grouped time series', *Computational Statistics and Data Analysis* **97**, 16–32.
- Kourentzes, N., Petropoulos, F. & Trapero, J. R. (2014), 'Improving forecasting by estimating time series structural components across multiple frequencies', *International Journal of Forecasting* **30**(2), 291–302.
- Manzan, S. & Zerom, D. (2008), 'A bootstrap-based non-parametric forecast density', *International Journal of Forecasting* **24**(3), 535–550.
- Marcellino, M. (1999), 'Some consequences of temporal aggregation in empirical analysis', *Journal of Business & Economic Statistics* **17**(1), 129–136.
- Nijman, T. E. & Palm, F. C. (1990), 'Disaggregate sampling in predictive models', *Journal of Business & Economic Statistics* **8**(4), 405–415.
- Panagiotelis, A., Athanasopoulos, G., Hyndman, R. J., Jiang, B. & Vahid, F. (2019), 'Macroeconomic forecasting for Australia using a large number of predictors', *International Journal of Forecasting* (forthcoming).

- Rossana, R. & Seater, J. (1995), 'Temporal aggregation and economic times series', *Journal of Business & Economic Statistics* **13**(4), 441–451.
- Schäfer, J. & Strimmer, K. (2005), 'A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics', *Statistical Applications in Genetics and Molecular Biology* **4**(1), 1–30.
- Shang, H. L. & Hyndman, R. J. (2017), 'Grouped functional time series forecasting: An application to age-specific mortality rates', *Journal of Computational and Graphical Statistics* **26**(2), 330–343.
- Silvestrini, A., Salto, M., Moulin, L. & Veredas, D. (2008), 'Monitoring and forecasting annual public deficit every month: the case of France', *Empirical Economics* **34**(3), 493–524.
- Smith, M. S. & Vahey, S. P. (2016), 'Asymmetric forecast densities for us macroeconomic variables from a gaussian copula model of cross-sectional and serial dependence', *Journal of Business & Economic Statistics* **34**(3), 416–434.
- Taieb, S. B., Taylor, J. W. & Hyndman, R. J. (2017), 'Hierarchical probabilistic forecasting of electricity demand with smart meter data', pp. 1–30.
- Tiao, G. C. (1972), 'Asymptotic behaviour of temporal aggregates of time series', *Biometrika* **59**(3), 525–531.
- Vilar, J. A. & Vilar, J. A. (2013), 'Time series clustering based on nonparametric multidimensional forecast densities', *Electronic Journal of Statistics* **7**(1), 1019–1046.
- Villegas, M. A. & Pedregal, D. J. (2018), 'Supply chain decision support systems based on a novel hierarchical forecasting approach', *Decision Support Systems* **114**, 29–36.
- Weiss, C. (2018), *Essays in Hierarchical Time Series Forecasting and Forecast Combination*, PhD thesis, University of Cambridge.
- Wickramasuriya, S. L., Athanasopoulos, G. & Hyndman, R. J. (2018), 'Optimal forecast reconciliation for hierarchical and grouped time series through trace minimization', *Journal of the American Statistical Association* **145**(9), 1–45.
- Yagli, G. M., Yang, D. & Srinivasan, D. (2019), 'Reconciling solar forecasts: Sequential reconciliation', *Solar Energy* **179**, 391–397.
- Yang, D., Quan, H., Disfani, V. R. & Liu, L. (2017), 'Reconciling solar forecasts: Geographical hierarchy', *Solar Energy* **146**, 276–286.
- Zellner, A. & Montmarquette, C. (1971), 'A study of some aspects of temporal aggregation problems in econometric analyses', *The Review of Economics and Statistics* **53**(4), 335–342.
- Zhang, Y. & Dong, J. (2019), 'Least squares-based optimal reconciliation method for hierarchical forecasts of wind power generation', *IEEE Transactions on Power Systems* (forthcoming).

Appendix

Table 1 Variables, Series IDs and their descriptions for the Income approach

Variable	Series ID	Description
Gdpi	A2302467A	GDP(I)
Sdi	A2302413V	Statistical discrepancy (I)
Tsi	A2302412T	Taxes less subsidies (I)
TfiCoeWns	A2302399K	Compensation of employees; Wages and salaries
TfiCoeEsc	A2302400J	Compensation of employees; Employers' social contributions
TfiCoe	A2302401K	Compensation of employees
TfiGosCopNfnPvt	A2323369L	Private non-financial corporations; Gross operating surplus
TfiGosCopNfnPub	A2302403R	Public non-financial corporations; Gross operating surplus
TfiGosCopNfn	A2302404T	Non-financial corporations; Gross operating surplus
TfiGosCopFin	A2302405V	Financial corporations; Gross operating surplus
TfiGosCop	A2302406W	Total corporations; Gross operating surplus
TfiGosGvt	A2298711F	General government; Gross operating surplus
TfiGosDwl	A2302408A	Dwellings owned by persons; Gross operating surplus
TfiGos	A2302409C	All sectors; Gross operating surplus
TfiGmi	A2302410L	Gross mixed income
Tfi	A2302411R	Total factor income

Table 2 Variables, Series IDs and their descriptions for Expenditure Approach

Variable	Series ID	Description
Gdpe	A2302467A	GDP(E)
Sde	A2302566J	Statistical Discrepancy(E)
Exp	A2302564C	Exports of goods and services
Imp	A2302565F	Imports of goods and services
Gne	A2302563A	Gross national exp.
GneDfdFceGvtNatDef	A2302523J	Gen. gov. - National; Final consumption exp. - Defence
GneDfdFceGvtNatNdf	A2302524K	Gen. gov. - National; Final consumption exp. - Non-defence
GneDfdFceGvtNat	A2302525L	Gen. gov. - National; Final consumption exp.
GneDfdFceGvtSnl	A2302526R	Gen. gov. - State and local; Final consumption exp.
GneDfdFceGvt	A2302527T	Gen. gov.; Final consumption exp.
GneDfdFce	A2302529W	All sectors; Final consumption exp.
GneDfdGfcPvtTdwNnu	A2302543T	Pvt.; Gross fixed capital formation (GFCF)
GneDfdGfcPvtTdwAna	A2302544V	Pvt.; GFCF - Dwellings - Alterations and additions
GneDfdGfcPvtTdw	A2302545W	Pvt.; GFCF - Dwellings - Total
GneDfdGfcPvtOtc	A2302546X	Pvt.; GFCF - Ownership transfer costs
GneDfdGfcPvtPbiNdcNbd	A2302533L	Pvt. GFCF - Non-dwelling construction - New building
GneDfdGfcPvtPbiNdcNec	A2302534R	Pvt.; GFCF - Non-dwelling construction - New engineering construction
GneDfdGfcPvtPbiNdcSha	A2302535T	Pvt.; GFCF - Non-dwelling construction - Net purchase of second hand assets
GneDfdGfcPvtPbiNdc	A2302536V	Pvt.; GFCF - Non-dwelling construction - Total
GneDfdGfcPvtPbiNdmNew	A2302530F	Pvt.; GFCF - Machinery and equipment - New
GneDfdGfcPvtPbiNdmSha	A2302531J	Pvt.; GFCF - Machinery and equipment - Net purchase of second hand assets
GneDfdGfcPvtPbiNdm	A2302532K	Pvt.; GFCF - Machinery and equipment - Total
GneDfdGfcPvtPbiCbr	A2716219R	Pvt.; GFCF - Cultivated biological resources
GneDfdGfcPvtPbiIprRnd	A2716221A	Pvt.; GFCF - Intellectual property products - Research and development
GneDfdGfcPvtPbiIprMnp	A2302539A	Pvt.; GFCF - Intellectual property products - Mineral and petroleum exploration
GneDfdGfcPvtPbiIprCom	A2302538X	Pvt.; GFCF - Intellectual property products - Computer software
GneDfdGfcPvtPbiIprArt	A2302540K	Pvt.; GFCF - Intellectual property products - Artistic originals
GneDfdGfcPvtPbiIpr	A2716220X	Pvt.; GFCF - Intellectual property products Total
GneDfdGfcPvtPbi	A2302542R	Pvt.; GFCF - Total private business investment
GneDfdGfcPvt	A2302547A	Pvt.; GFCF
GneDfdGfcPubPcpCmw	A2302548C	Plc. corporations - Commonwealth; GFCF
GneDfdGfcPubPcpSnl	A2302549F	Plc. corporations - State and local; GFCF
GneDfdGfcPubPcp	A2302550R	Plc. corporations; GFCF Total
GneDfdGfcPubGvtNatDef	A2302551T	Gen. gov. - National; GFCF - Defence
GneDfdGfcPubGvtNatNdf	A2302552V	Gen. gov. - National ; GFCF - Non-defence
GneDfdGfcPubGvtNat	A2302553W	Gen. gov. - National ; GFCF Total
GneDfdGfcPubGvtSnl	A2302554X	Gen. gov. - State and local; GFCF
GneDfdGfcPubGvt	A2302555A	Gen. gov.; GFCF
GneDfdGfcPub	A2302556C	Plc.; GFCF
GneDfdGfc	A2302557F	All sectors; GFCF

Table 3 Variables, Series IDs and their descriptions for Changes in Inventories - Expenditure Approach

Variable	Series ID	Description
GneCii	A2302562X	Changes in Inventories
GneCiiPfm	A2302560V	Farm
GneCiiPba	A2302561W	Public authorities
GneCiiPnf	A2302559K	Private; Non-farm Total
GneCiiPnfMin	A83722619L	Private; Mining (B)
GneCiiPnfMan	A3348511X	Private; Manufacturing (C)
GneCiiPnfWht	A3348512A	Private; Wholesale trade (F)
GneCiiPnfRet	A3348513C	Private; Retail trade (G)
GneCiiPnfOnf	A2302273C	Private; Non-farm; Other non-farm industries

Table 4 Variables, Series IDs and their descriptions for Household Final Consumption - Expenditure Approach

Variable	Series ID	Description
GneDfdHfc	A2302254W	Household Final Consumption Expenditure
GneDfdFceHfcFud	A2302237V	Food
GneDfdFceHfcAbt	A3605816F	Alcoholic beverages and tobacco
GneDfdFceHfcAbtCig	A2302238W	Cigarettes and tobacco
GneDfdFceHfcAbtAlc	A2302239X	Alcoholic beverages
GneDfdFceHfcCnf	A2302240J	Clothing and footwear
GneDfdFceHfcHwe	A3605680F	Housing, water, electricity, gas and other fuels
GneDfdFceHfcHweRnt	A3605681J	Actual and imputed rent for housing
GneDfdFceHfcHweWsc	A3605682K	Water and sewerage charges
GneDfdFceHfcHweEgf	A2302242L	Electricity, gas and other fuel
GneDfdFceHfcFhe	A2302243R	Furnishings and household equipment
GneDfdFceHfcFheFnt	A3605683L	Furniture, floor coverings and household goods
GneDfdFceHfcFheApp	A3605684R	Household appliances
GneDfdFceHfcFheTls	A3605685T	Household tools
GneDfdFceHfcHlt	A2302244T	Health
GneDfdFceHfcHltMed	A3605686V	Medicines, medical aids and therapeutic appliances
GneDfdFceHfcHltHsv	A3605687W	Total health services
GneDfdFceHfcTpt	A3605688X	Transport
GneDfdFceHfcTptPvh	A2302245V	Purchase of vehicles
GneDfdFceHfcTptOvh	A2302246W	Operation of vehicles
GneDfdFceHfcTptTsv	A2302247X	Transport services
GneDfdFceHfcCom	A2302248A	Communications
GneDfdFceHfcRnc	A2302249C	Recreation and culture
GneDfdFceHfcEdc	A2302250L	Education services
GneDfdFceHfcHcr	A2302251R	Hotels, cafes and restaurants
GneDfdFceHfcHcrCsv	A3605694V	Catering services
GneDfdFceHfcHcrAsv	A3605695W	Accommodation services
GneDfdFceHfcMis	A3605696X	Miscellaneous goods and services
GneDfdFceHfcMisOgd	A3605697A	Other goods
GneDfdFceHfcMisIfs	A2302252T	Insurance and other financial services
GneDfdFceHfcMisOsv	A3606485T	Other services