



Department of Econometrics and Business Statistics

<http://business.monash.edu/econometrics-and-business-statistics/research/publications>

# **Probabilistic Forecasts in Hierarchical Time Series**

Puwasala Gamakumara  
Anastasios Panagiotelis  
George Athanasopoulos  
Rob J Hyndman

July 2018

Working Paper ??/??

# Probabilistic Forecasts in Hierarchical Time Series

**Puwasala Gamakumara**

Department of Econometrics and Business Statistics,  
Monash University,  
VIC 3800, Australia.

Email: Puwasala.Gamakumara@monash.edu

**Anastasios Panagiotelis**

Department of Econometrics and Business Statistics,  
Monash University,  
VIC 3800, Australia.

Email: Anastasios.Panagiotelis@monash.edu

**George Athanasopoulos**

Department of Econometrics and Business Statistics,  
Monash University,  
VIC 3800, Australia.

Email: George.Athanasopoulos@monash.edu

**Rob J Hyndman**

Department of Econometrics and Business Statistics,  
Monash University,  
VIC 3800, Australia.

Email: Rob.Hyndman@monash.edu

2 July 2018

**JEL classification:** ??

# Probabilistic Forecasts in Hierarchical Time Series

## Abstract

TBC

## 1 Introduction

Sketch narrative to tell

- Coherence is important in point forecasting to
  - Align decisions
  - Improve forecasts
- We propose a geometric interpretation via linear transformations including special case of projections.
- These nest existing methods but the interpretation has advantages
  - Prove that projections preserve unbiasedness. Previously assumed in special cases.
  - Prove reconciliation (using projection) always improves sum of squared errors (or a scaled/transformed version thereof).
  - Most importantly generalises to probabilistic forecasting.
- For elliptical distributions, true predictive can be recovered via linear reconciliation and in special cases, via projection.
- Discuss evaluation of forecasts
  - Log score improper when coherent v incoherent
  - Energy score can be used instead
  - Discuss scores for coherent v coherent.
- Finally evidence MinT is best from simulation

Many research applications involve a large collection of time series, some of which are aggregates of others. These are called hierarchical time series. For example, electricity demand of a country

can be disaggregated along a geographical hierarchy: the electricity demand of the whole country can be divided into the demand of states, cities, and households.

When forecasting such time series, it is important to have “coherent” forecasts across the hierarchy: aggregates of the forecasts at lower levels should be equal to the forecasts at the upper levels of aggregation. In other words, sums of forecasts should be equal to the forecasts of the sums.

The traditional approaches to produce coherent point forecasts are the bottom-up, top-down and middle-out methods. In the bottom-up approach, forecasts of the lowest level are first generated and they are simply aggregated to forecast upper levels of the hierarchy (Dunn, Williams, and Dechaine, 1976). In contrast, the top-down approach involves forecasting the most aggregated series first and then disaggregating these forecasts down the hierarchy based on the corresponding proportions of observed data (Gross and Sohl, 1990). Many studies have discussed the relative advantages and disadvantages of bottom-up and top-down methods, and situations in which each would provide reliable forecasts (Schwarzkopf, Tersine, and Morris, 1988; Kahn, 1998; Lapide, 1998; Fliedner, 2001). A compromise between these two approaches is the middle-out method which entails forecasting each series of a selected middle level in the hierarchy and then forecasting upper levels by the bottom-up method and lower levels by the top-down method.

It is apparent that these three approaches use only part of the information available when producing coherent forecasts. This might result in inaccurate forecasts. For example, if the bottom-level series are highly volatile or noisy, and hence challenging to forecast, then the resulting forecasts from the bottom-up approach are likely to be inaccurate.

As an alternative to these traditional methods, Hyndman et al. (2011) proposed to utilize the information from all levels of the hierarchy to obtain coherent point forecasts in a two stage process. In the first stage, the forecasts of all series are independently obtained by fitting univariate models for individual series in the hierarchy. It is very unlikely that these forecasts are coherent. Thus in the second stage, these forecasts are optimally combined through a regression model to obtain coherent forecasts. This second step is referred to as “reconciliation” since it takes a set of incoherent forecasts and revises them to be coherent. The approach was further improved by Wickramasuriya, Athanasopoulos, and Hyndman (2018) who proposed the “MinT” algorithm to obtain optimally reconciled point forecasts by minimizing the mean squared coherent forecast errors.

Traditional bottom-up, top-down and middle-out forecasting methods are not strictly reconciliation methods since they use only a part of the information from the hierarchy to produce coherent forecasts.

Previous studies on coherent point forecasting have shown that reconciliation provides better coherent forecasts than the traditional bottom-up and top-down methods (Hyndman et al., 2011; Erven and Cugliari, 2014; Wickramasuriya, Athanasopoulos, and Hyndman, 2018). However, this idea has not been explored in the context of probabilistic forecasting.

Point forecasts are limited because they provide no indication of forecast uncertainty. Providing prediction intervals helps, but a richer description of forecast uncertainty is obtained by estimating the entire forecast distribution. These are often called “probabilistic forecasts” (Gneiting and Katzfuss, 2014). For example, McSharry, Bouwman, and Bloemhof (2005) produced probabilistic forecasts for electricity demand, Ben Taieb et al. (2017) for smart meter data, Pinson et al. (2009) for wind power generation, and Gel, Raftery, and Gneiting (2004), Gneiting et al. (2005) and Gneiting and Raftery (2005) for various weather variables.

Although there is a rich and growing literature on producing coherent point forecasts of hierarchical time series, little attention has been given to coherent probabilistic forecasts. The only relevant paper we are aware of is Ben Taieb et al. (2017), who recently proposed an algorithm to produce coherent probabilistic forecasts and applied it to UK electricity smart meter data. In their approach, a sample from the bottom-level forecast distribution is first generated, and then aggregated to obtain coherent probabilistic forecasts of the upper levels of the hierarchy. Hence this method is a bottom-up approach. They propose to first use the MinT algorithm to reconcile the means of the bottom-level forecast distributions, and then a copula-based approach is employed to model the dependency structure of the hierarchy. The resulting multi-dimensional distribution is used to generating empirical forecast distributions for all bottom-level series. Thus, while Ben Taieb et al. (2017) provide coherent probabilistic forecasts, they do no forecast reconciliation of the distributions. In that sense, their approach is analogous to bottom-up point forecasting rather than forecast reconciliation.

After introducing our notation in Section 2, we define what is meant by probabilistic forecast reconciliation for hierarchical time series in Section 3. First, we provide a new definition for coherency of point forecasts, and the reconciliation of a set of incoherent point forecasts, using concepts related to vector spaces and measure theory. Based on these, we provide a rigorous

definition for probabilistic forecast reconciliation, and how we can reconcile the incoherent forecast densities in practice.

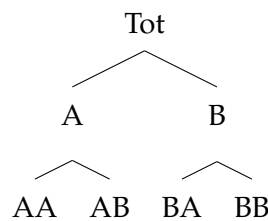
Further, due to the aggregation structure of the hierarchy, the probability distribution is degenerate and hence the forecast distribution should also be degenerate. In Section 4, we discuss in detail how this degeneracy will be taken care of in probabilistic forecast reconciliation, and in Section 5 we consider the evaluation of probabilistic hierarchical forecasts.

Some theoretical results on probabilistic forecast reconciliation in the Gaussian framework are given in Section 6, including a simulation study to show the importance of reconciliation in the probabilistic framework.

We conclude with some thoughts on extensions and limitations in Section 7.

## 2 Hierarchical Time Series

In the section, and throughout the paper, we will try to follow notational conventions used in Wickramasuriya, Athanasopoulos, and Hyndman (2018) as much as possible. A *hierarchical time series* is a collection of  $n$  variables where some variables are aggregates of other variables. For example in the hierarchy depicted in Figure 1 below, the variable labelled *Tot* is the sum of the series *A* and series *B*, the series *A* is the sum of series *AA* and series *AB* and the series *B* is the sum of the series *BA* and *BB*. The *bottom level series* are defined as those  $m$  variables that cannot be formed as aggregates of other variables, in the example in Figure 1 these are the series *AA*, *AB*, *BA* and *BB*.



**Figure 1:** Two level hierarchical diagram.

We let  $\mathbf{y}_t \in \mathbb{R}^n$  be a vector comprised of observations of all variables in the hierarchy at time  $t$ , and  $\mathbf{b}_t \in \mathbb{R}^m$  is a vector comprised of observations of all bottom-level series at time  $t$ . The hierarchical structure of the data imply the following holds for all  $t$

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

where  $S$  is an  $n \times m$  constant matrix that encodes the aggregation constraints. For the hierarchy in Figure 1,  $\mathbf{y}_t = [y_{Tot,t}, y_{A,t}, y_{B,t}, y_{C,t}, y_{AA,t}, y_{AB,t}, y_{BA,t}, y_{BB,t}]'$ ,  $\mathbf{b}_t = [y_{AA,t}, y_{AB,t}, y_{BA,t}, y_{BB,t}]'$ ,  $m = 4$ ,  $n = 7$ , and

$$S = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ & & & I_4 \end{pmatrix},$$

where  $I_4$  is a  $4 \times 4$  identity matrix.

### 3 Coherent forecasts

It is desirable that forecasts, whether point forecasts or probabilistic forecasts, should in some sense respect aggregation constraints. We follow other authors Wickramasuriya, Athanasopoulos, and Hyndman, 2018 in using the nomenclature *coherence* to describe this property.

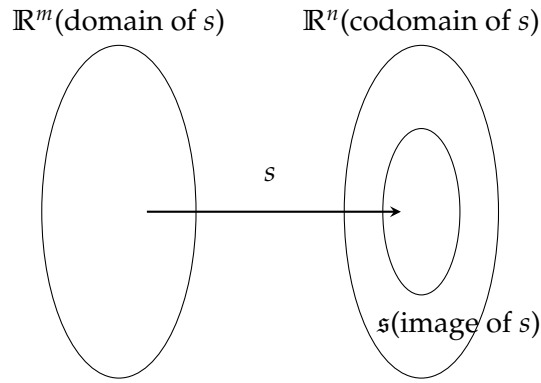
FPP?

We now provide new definitions for coherent forecasts in terms of vector spaces that give a geometric understanding of the problem thus facilitating the development of the probabilistic forecast reconciliation in section 4.

Recall that  $\mathbf{y}_t \in \mathbb{R}^n$  is a  $n$ -dimensional time series subject to the linear aggregation constraint  $\mathbf{y}_t = S\mathbf{b}_t$ , where  $\mathbf{b}_t \in \mathbb{R}^m$  and  $S$  is an  $n \times m$  constant matrix.

**Definition 3.1** (Coherent subspace). The  $m$ -dimensional linear subspace  $\mathfrak{s} \subset \mathbb{R}^n$  that is spanned by the columns of  $S$ , i.e.  $\mathfrak{s} = \text{span}(S)$ , is defined as the *coherent space*.

Also at times it will be useful to think of pre-multiplication by  $S$  as a linear mapping from  $\mathbb{R}^m$  to  $\mathbb{R}^n$  in which case we use the notation  $s(\cdot)$ . Although the codomain of  $s(\cdot)$  is  $\mathbb{R}^n$  its image is the coherent space  $\mathfrak{s}$  as depicted in Figure 2.



**Figure 2:** The domain, codomain and image of the mapping  $s$ .

**Definition 3.2** (Coherent Point Forecasts). Let  $\check{y}_{t+h|t} \in \mathbb{R}^n$  be a point forecast of the values of all series in the hierarchy at time  $t + h$ , made using information up to and including time  $t$ . Then  $\check{y}_{t+h|t}$  is *coherent* if  $\check{y}_{t+h|t} \in \mathfrak{s}$ .

Let  $(\mathbb{R}^m, \mathcal{F}_{\mathbb{R}^m}, \nu)$  be a probability triple, where  $\mathcal{F}_{\mathbb{R}^m}$  is the usual  $\sigma$ -algebra on  $\mathbb{R}^m$ . Let  $\check{\nu}$  be a probability measure on  $\mathfrak{s}$  with  $\sigma$ -algebra  $\mathcal{F}_{\mathfrak{s}}$ . Here  $\mathcal{F}_{\mathfrak{s}}$  is formed as collection of sets  $s(\mathcal{B})$ , where  $s(\mathcal{B})$  denotes the image of the set  $\mathcal{B} \in \mathcal{F}_{\mathbb{R}^m}$  under the mapping  $s(\cdot)$ .

**Definition 3.3** (Coherent Probabilistic Forecasts). The measure  $\check{\nu}$  is coherent if it has the property

$$\check{\nu}(s(\mathcal{B})) = \nu(\mathcal{B}) \quad \forall \mathcal{B} \in \mathcal{F}_{\mathbb{R}^m},$$

A probabilistic forecast for time  $t + h$  is coherent if uncertainty in  $y_{t+h|t}$  conditional on all information up to time  $t$  is characterised by the probability triple  $(\mathfrak{s}, \mathcal{F}_{\mathfrak{s}}, \check{\nu})$ .

These definitions of the coherent space  $\mathfrak{s}$  and coherent point and probabilistic forecasts are defined in terms of the mapping  $s(\cdot)$  and may give the impression that the bottom level series play an important role in the definition. However, alternative definitions could be formed using any set of basis vectors that spans  $\mathfrak{s}$ . For example, consider the most simple three variable hierarchy where  $y_{1,t} = y_{2,t} + y_{3,t}$ . In this case the matrix  $S$  has columns  $(1, 1, 0)'$  and  $(1, 0, 1)'$  spanning  $\mathfrak{s}$  and premultiplying by  $S$  transforms arbitrary values of  $y_{2,t}$  and  $y_{3,t}$  into a coherent vector for the full hierarchy. However the columns  $(1, 0, 1)'$  and  $(0, 1, -1)'$  also span  $\mathfrak{s}$  and define a mapping that transforms arbitrary values of  $y_{1,t}$  and  $y_{2,t}$  into a coherent vector for the full hierarchy. The definitions above could be made in terms of any series and not just the bottom level series. In general, we call the series (or linear combinations thereof) used in the definitions



of coherence *basis series*. Unless stated otherwise, we will always assume that the basis series are the bottom level series as in Definition 3.2 and Defintion 3.3, since this facilitates comparison with existing approaches in the literature.

To the best of our knowledge, the only other definition of coherent probabilistic forecasts is given by Ben Taieb et al. (2017) who define coherent probabilistic forecasts in terms of convolutions. According to their definition, probabilistic forecasts are coherent when a convolution of forecast distributions of disaggregate series is identical to the forecast distribution of the corresponding aggregate series. Their definition is consistent with our definition, our reason for providing a different definition is that the geometric understanding of coherence will facilitate our definitions of point and probabilistic forecast reconciliation to which we now turn our attention.

## 4 Forecast reconciliation

Initially we define point forecast reconciliation, before extending the idea to the probabilistic setting.

### 4.1 Point forecast reconciliation

Let  $\hat{\mathbf{y}}_{t+h|t} \in \mathbb{R}^n$  be any set of incoherent point forecasts at time  $t + h$  using information up to and including time  $t$ . Let  $\mathbf{G}$  and  $\mathbf{d}$  be an  $m \times n$  matrix and  $m \times 1$  vector respectively and  $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be the mapping  $g(\mathbf{y}) = \mathbf{G}\mathbf{y} + \mathbf{d}$ .

**Definition 4.1.** The point forecast  $\tilde{\mathbf{y}}_{t+h|t}$  “reconciles”  $\hat{\mathbf{y}}_{t+h|t}$  with respect to the mapping  $g(\cdot)$  iff

$$\tilde{\mathbf{y}}_{t+h|t} = \mathbf{S} (\mathbf{G}\hat{\mathbf{y}}_{t+h|t} + \mathbf{d}) .$$

Many choices of  $g(\cdot)$  currently extant in the literature including the so called OLS Hyndman et al., 2011, WLS and MinTWickramasuriya, Athanasopoulos, and Hyndman, 2018 methods are special cases where  $s \circ g$  is a projection and are summarised in Table 1. These can be defined so that  $\mathbf{G} = (\mathbf{R}'_{\perp} \mathbf{S})^{-1} \mathbf{S}'$  and  $\mathbf{d} = \mathbf{0}$ . Here,  $\mathbf{R}_{\perp}$  is a  $n \times m$  orthogonal complement to the an  $n \times (n - m)$  matrix  $\mathbf{R}$  where the columns of the latter span the null space of  $\mathbf{s}_{\perp}$ . For example, a straightforward choice of  $\mathbf{R}$  for the most simple three variable hierarchy where  $y_{1,t} = y_{2,t} + y_{3,t}$ , is the vector  $(1, -1, -1)$  which is orthogonal (in the Euclidean sense) to the columns of  $\mathbf{S}$ . In this case, the matrix  $\mathbf{R}$  can be interpreted as a ‘restrictions’ matrix since it has the property that  $\mathbf{R}'\mathbf{y} = \mathbf{0}$  for coherent  $\mathbf{y}$ . For the example provided,  $\mathbf{R}'_{\perp} = \mathbf{S}$  and reconciliation corresponds to

First  
WLS  
refer-  
ence?

**Table 1:** Several possible estimates of  $R'_\perp$ . For  $n < T$ ,  $\hat{W}_{T+1}^{sam}$  is an unbiased and consistent estimator for  $W_{T+1}$ .  $\hat{W}_{T+1}^{shr}$  is a shrinkage estimator which is more suitable for large dimensions.  $\hat{W}_{T+1}^{shr}$  was proposed by Schäfer and Strimmer (2005) and also used by Wickramasuriya, Athanasopoulos, and Hyndman (2018), where  $\text{Diag}(A)$  denotes the diagonal matrix of  $A$ ,  $\tau = \frac{\sum_{i \neq j} \text{Var}(\hat{r}_{ij})}{\sum_{i \neq j} \hat{r}_{ij}^2}$ , and  $\hat{r}_{ij}$  is the  $ij$ th element of the sample correlation matrix.

Method	Estimate of $W_h$	Estimate of $R'_\perp$
OLS	$I$	$S'$
MinT(Sample)	$\hat{W}_{T+1}^{sam}$	$S'(\hat{W}_{T+1}^{sam})^{-1}$
MinT(Shrink)	$\hat{W}_{T+1}^{shr} = \tau \text{Diag}(\hat{W}_{T+1}^{sam}) + (1 - \tau) \hat{W}_{T+1}^{sam}$	$S'(\hat{W}_{T+1}^{shr})^{-1}$
WLS	$\hat{W}_{T+1}^{wls} = \text{Diag}(\hat{W}_{T+1}^{shr})$	$S'(\hat{W}_{T+1}^{wls})^{-1}$

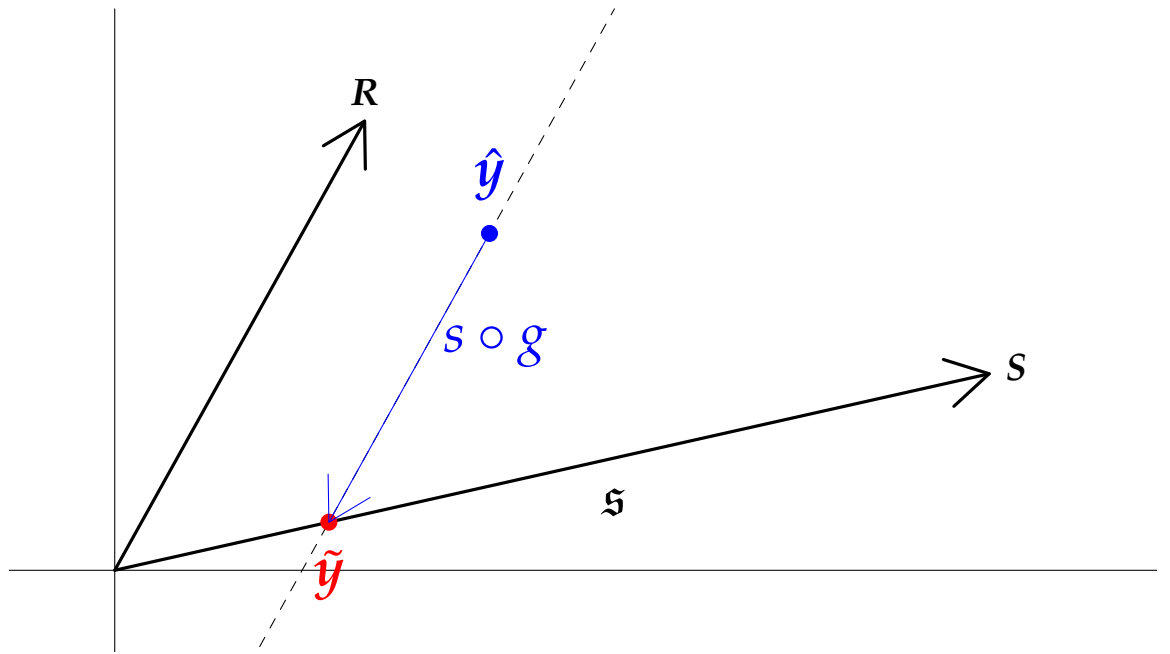
the so called ‘OLS’ method Hyndman et al., 2011. For the case where  $R'_\perp \neq S$ , for example WLS and MinT, there are two possible interpretations. One is that these are oblique projections in Euclidean space where the columns of  $R$  are ‘directions’ along which incoherent point forecasts are projected onto the coherent space  $\mathfrak{s}$ . Alternatively, since  $R'_\perp$  is usually written in the form  $S'W^{-1}$ , these projections can be thought of as orthogonal projections after applying the transformation  $W^{-1/2}$ . A schematic providing a geometric interpretation of point reconciliation is given in Figure 3.

To illustrate further note that the columns of  $S$  and  $R$  provide a basis for  $\mathbb{R}^n$ . As such any incoherent set of point forecasts  $\hat{y}_{t+h|t} \in \mathbb{R}^n$ , can be expressed in terms of coordinates in the basis defined by  $S$  and  $R$ . Let  $\tilde{b}_{t+h|t}$  and  $\tilde{a}_{t+h|t}$  be the coordinates corresponding to  $S$  and  $R$  respectively after a change of basis. The process of reconciliation involves setting  $\tilde{b}_{t+h|t}$  to be the values of the reconciled bottom-level series and setting  $\tilde{a}_{t+h} = \mathbf{0}$  to ensure coherence. From properties of linear algebra it follows that

$$\hat{y}_{t+h|t} = (S \ R) \begin{pmatrix} \tilde{b}_{t+h|t} \\ \tilde{a}_{t+h|t} \end{pmatrix} = S\tilde{b}_{t+h|t} + R\tilde{a}_{t+h|t},$$

while setting  $\tilde{a}_{t+h|t} = \mathbf{0}$  gives the reconciled point forecast

$$\tilde{y}_{t+h|t} = S\tilde{b}_{t+h|t}$$



**Figure 3:** Summary of probabilistic point reconciliation. The mapping  $G(\cdot)$  projects the unreconciled forecast  $\hat{\mathbf{y}}$  onto  $\mathbf{S}$ . Note that since the smallest hierarchy involves three dimensions, this figure is only a schematic

In order to find  $\tilde{\mathbf{b}}_{t+h|t}$  we require the inverse  $(\mathbf{S} \ \mathbf{R})^{-1}$  which is given by

$$(\mathbf{S} \ \mathbf{R})^{-1} = \begin{pmatrix} (\mathbf{R}'_{\perp} \mathbf{S})^{-1} \mathbf{R}'_{\perp} \\ (\mathbf{S}'_{\perp} \mathbf{R})^{-1} \mathbf{S}'_{\perp} \end{pmatrix},$$

where  $\mathbf{S}_{\perp}$  is the orthogonal complements of  $\mathbf{S}$ . Thus it follows that  $\tilde{\mathbf{b}}_{t+h} = (\mathbf{R}'_{\perp} \mathbf{S})^{-1} \mathbf{R}'_{\perp} \hat{\mathbf{y}}_{t+h}$  and  $\tilde{\mathbf{y}}_{t+h|t} = \mathbf{S}(\mathbf{R}'_{\perp} \mathbf{S})^{-1} \mathbf{R}'_{\perp} \hat{\mathbf{y}}_{t+h|t}$ . Here  $(\mathbf{R}'_{\perp} \mathbf{S})^{-1} \mathbf{R}'_{\perp}$  corresponds to  $\mathbf{G}$  as defined previously.

Point reconciliation methods based on projections will always minimise the distance between unreconciled and reconciled forecasts, however the specific distance will depend on the choice of  $\mathbf{R}$ . For example Hyndman et al., 2011 consider  $\hat{\mathbf{y}}_{t+h}^{OLS} = \mathbf{S}(\mathbf{S}'\mathbf{S})^{-1}\mathbf{S}'\hat{\mathbf{y}}_{t+h}$  which minimises the Euclidean distance between  $\hat{\mathbf{y}}$  and  $\tilde{\mathbf{y}}$ . Wickramasuriya, Athanasopoulos, and Hyndman, 2018 consider  $\hat{\mathbf{y}}_{t+h}^{MinT} = \mathbf{S}(\mathbf{S}'\mathbf{W}_h^{-1}\mathbf{S})^{-1}\mathbf{S}'\mathbf{W}_h^{-1}\hat{\mathbf{y}}_{t+h}$ , where  $\mathbf{W}$  is an estimate of the variance covariance

matrix of the unreconciled errors. This minimises the Mahalanobis distance between  $\hat{\mathbf{y}}$  and  $\tilde{\mathbf{y}}$ . Bottom up methods minimise distance between reconciled and unreconciled forecasts only along dimensions corresponding to the bottom level series. As such, bottom up methods should be thought of as a boundary case of reconciliation methods, since they ultimately do not use information at all levels of the hierarchy.

Before generalising the concept of point reconciliation to probabilistic forecasts we state two theorems that motivate the use of projections for point forecast reconciliation. First, let  $\boldsymbol{\mu}_{t+h|t} := E(\mathbf{y}_{t+h} | \mathbf{y}_1, \dots, \mathbf{y}_t)$  and assume  $\hat{\mathbf{y}}_{t+h|t}$  is an unbiased prediction, that is  $E_{1:t}(\hat{\mathbf{y}}_{t+h|t}) = \boldsymbol{\mu}_{t+h|t}$  where the subscript  $1:t$  denotes an expectation taken over the training sample.

mention  
reduc-  
ing dis-  
tance to  
truth?

**Theorem 4.1** (Unbiasedness preserving property). *The reconciled point forecast will also be an unbiased prediction as long as  $s \circ g$  is a projection*

*Proof.* The expected value of the reconciled forecast is given by

$$\begin{aligned} E_{1:t}(\tilde{\mathbf{y}}_{t+h|t}) &= E_{1:t}(\mathbf{S}\mathbf{G}\hat{\mathbf{y}}_{t+h|t}) \\ &= \mathbf{S}\mathbf{G}E_{1:t}(\hat{\mathbf{y}}_{t+h|t}) \\ &= \mathbf{S}\mathbf{G}\boldsymbol{\mu}_{t+h|t} \end{aligned}$$

Since the aggregation constraints hold for the true data generating process  $\boldsymbol{\mu}_{t+h|t}$  must lie in  $\mathfrak{s}$ . If  $\mathbf{S}\mathbf{G}$  is a projection then it is equivalent to the identity map for all vectors that lie in its range. Therefore  $\mathbf{S}\mathbf{G}\boldsymbol{\mu}_{t+h|t} = \boldsymbol{\mu}_{t+h|t}$  when  $\mathbf{S}\mathbf{G}$  is a projection matrix.  $\square$

We note the same result does not hold for general  $\mathbf{G}$  even though the range of  $s \circ g$  is  $\mathfrak{s}$ . Now let  $\mathbf{y}_{t+h}$  be the realisation of the data generating process at time  $t+h$  and let  $\|\mathbf{v}\|_2$  be the L2 norm of vector  $\mathbf{v}$ . The following theorem shows that reconciliation never increases, and in most cases reduces the sum of squared errors of point forecasts.

**Theorem 4.2** (Distance reducing property). *If  $\tilde{\mathbf{y}}_{t+h|t} = \mathbf{S}\mathbf{G}\hat{\mathbf{y}}_{t+h|t}$  where  $\mathbf{G}$  is such that  $\mathbf{S}\mathbf{G}$  is an orthogonal projection then the following inequality holds*

$$\|(\tilde{\mathbf{y}}_{t+h|t} - \mathbf{y}_{t+h})\|_2^2 \leq \|(\hat{\mathbf{y}}_{t+h|t} - \mathbf{y}_{t+h})\|_2^2$$

*Proof.* Since the aggregation constraints must hold for all realisations,  $\mathbf{y}_{t+h} \in \mathfrak{s}$  and  $\mathbf{y}_{t+h} = \mathbf{S}\mathbf{G}\mathbf{y}_{t+h}$  whenever  $\mathbf{S}\mathbf{G}$  is a projection. Therefore

$$\begin{aligned} \|(\hat{\mathbf{y}}_{t+h|t} - \mathbf{y}_{t+h})\|_2 &= \|(\mathbf{S}\mathbf{G}\hat{\mathbf{y}}_{t+h|t} - \mathbf{S}\mathbf{G}\mathbf{y}_{t+h})\|_2 \\ &= \|\mathbf{S}\mathbf{G}(\hat{\mathbf{y}}_{t+h|t} - \mathbf{y}_{t+h})\|_2 \end{aligned}$$

The Cauchy-Schwarz inequality can be used to show that orthogonal projections are bounded operators, therefore

$$\|\mathbf{S}\mathbf{G}(\hat{\mathbf{y}}_{t+h|t} - \mathbf{y}_{t+h})\|_2 \leq \|(\hat{\mathbf{y}}_{t+h|t} - \mathbf{y}_{t+h})\|_2$$

□

 Find  
refer-  
ence or  
prove

The inequality is strict whenever  $\hat{\mathbf{y}}_{t+h|t} \notin \mathfrak{s}$ .

## 4.2 Probabilistic forecast reconciliation

We now extend the methodology of point forecast reconciliation to probabilistic forecasts

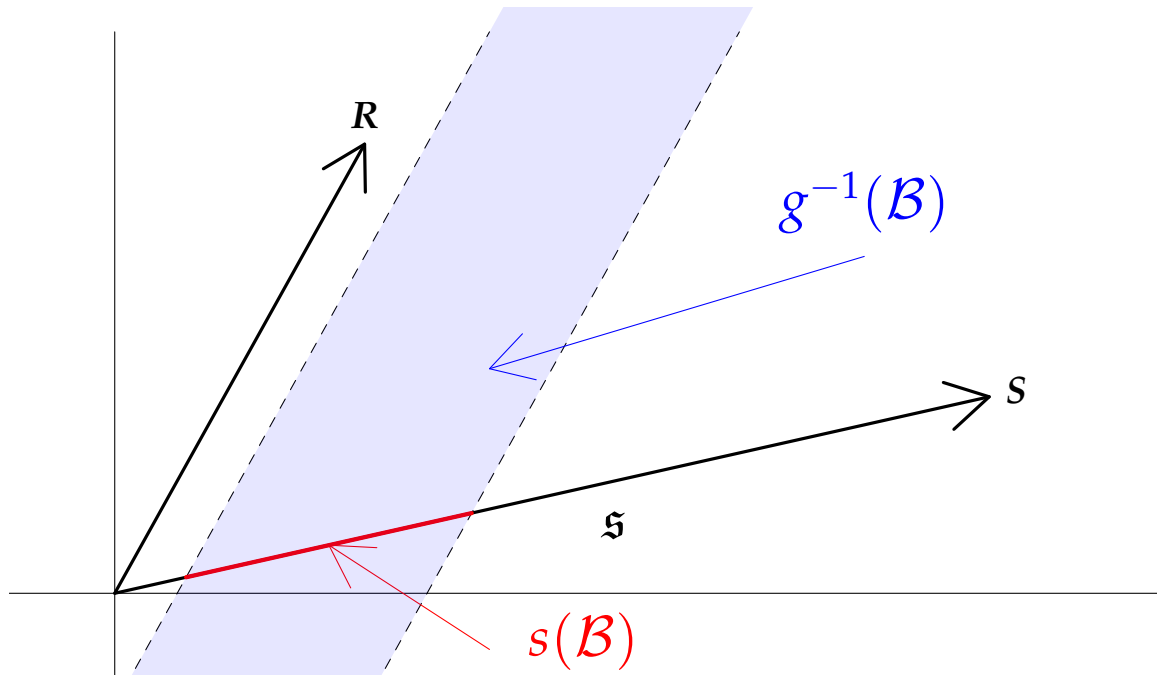
Let  $(\mathbb{R}^n, \mathcal{F}_{\mathbb{R}^n}, \hat{\nu})$  be an probability triple, that is incoherent and that characterises forecast uncertainty for all variables in the hierarchy at time  $t + h$  conditional on all information up to time  $t$ . This may be obtained from the first stage of the forecasting process e.g. by modelling and forecasting each series individually. Let  $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be a linear mapping. Let  $(\mathbb{R}^m, \mathcal{F}_{\mathbb{R}^m}, \nu)$  be a probability triple defined on  $\mathbb{R}^m$ .

**Definition 4.2.** We define the reconciled probability measure of  $\hat{\nu}$  with respect to the mapping  $g(\cdot)$  as a probability measure  $\tilde{\nu}$  on  $\mathfrak{s}$  with  $\sigma$ -algebra  $\mathcal{F}_{\mathfrak{s}}$  where the following holds

$$\tilde{\nu}(g(\mathcal{B})) = \nu(\mathcal{B}) = \hat{\nu}(g^{-1}(\mathcal{B})) \quad \forall \mathcal{B} \in \mathcal{F}_{\mathbb{R}^m},$$

where  $g^{-1}(\mathcal{B}) := \{\tilde{\mathbf{y}} \in \mathbb{R}^n : g(\tilde{\mathbf{y}}) \in \mathcal{B}\}$  is the pre-image of  $\mathcal{B}$ , that is the set of all points in  $\mathbb{R}^n$  that  $g(\cdot)$  maps to a point in  $\mathcal{B}$ .

This definition extends the notion of forecast reconciliation to the probabilistic setting. Under point reconciliation methods, the reconciled point forecast is equal to the unreconciled point forecast after the latter is passed through two linear mappings. Similarly, probabilistic forecast reconciliation assigns the same probability to two sets where the points in one set are obtained by



**Figure 4:** Summary of probabilistic forecast reconciliation. The probability that  $\mathbf{y}_{t+h|t}$  lies in the red line segment under the reconciled probabilistic forecast is defined to be equal to the probability that  $\mathbf{y}_{t+h|t}$  lies in the shaded blue area under the unreconciled probabilistic forecast. Note that since the smallest hierarchy involves three dimensions, this figure is only a schematic

passing all points in the other set through two linear mappings. This is depicted schematically when  $s \circ g$  is a projection in Figure 4.

Recall that when  $s \circ g$  is a projection, the case of point forecast reconciliation could be broken down into three steps. In the first,  $\hat{\mathbf{y}}_{t+h|t}$  is transformed into coordinates  $\tilde{\mathbf{b}}_{t+h|t}$  and  $\tilde{\mathbf{a}}_{t+h|t}$  via a change of basis. In the second,  $\tilde{\mathbf{a}}_{t+h|t}$  is discarded and  $\tilde{\mathbf{b}}_{t+h|t}$  are kept as the bottom level reconciled forecasts. In the third, reconciled forecasts for the entire hierarchy are recovered via  $\tilde{\mathbf{y}}_{t+h|t} = \mathbf{S}\tilde{\mathbf{b}}_{t+h|t}$ . We now outline the analogues to these three steps for probabilistic forecasts when predictive densities are available.

While  $\hat{\nu}$  is a probability measure for an  $n$ -vector  $\hat{\mathbf{y}}_{t+h|t}$ , probability statements in terms of a different coordinate system can be made via an appropriate change of basis. Letting  $f(\cdot)$

be generic notation for a probability density functions and following the notation from our definition of point forecast reconciliation where  $\hat{\mathbf{y}}_{t+h|t} = \mathbf{S}\tilde{\mathbf{b}}_{t+h|t} + \mathbf{R}\tilde{\mathbf{a}}_{t+h|t}$  we obtain

$$f(\hat{\mathbf{y}}_{t+h|t}) = f(\mathbf{S}\tilde{\mathbf{b}}_{t+h|t} + \mathbf{R}\tilde{\mathbf{a}}_{t+h|t}) |(\mathbf{S} \ \mathbf{R})|$$

The expression  $\hat{\nu}(g^{-1}(\mathcal{B}))$  in Definition 4.2 is equivalent to the probability statement  $\Pr(\hat{\mathbf{y}}_{t+h|t} \in g^{-1}(\mathcal{B}))$ . After the change of basis this is equivalent to  $\Pr(\tilde{\mathbf{b}} \in \mathcal{B})$  which implies

$$\begin{aligned} \Pr(\hat{\mathbf{y}}_{t+h|t} \in g^{-1}(\mathcal{B})) &= \int_{g^{-1}(\mathcal{B})} f(\hat{\mathbf{y}}_{t+h|t}) d\hat{\mathbf{y}}_{t+h|t} \\ &= \int_{\mathcal{B}} \int f(\mathbf{S}\tilde{\mathbf{b}}_{t+h|t} + \mathbf{R}\tilde{\mathbf{a}}_{t+h|t}) |(\mathbf{S} \ \mathbf{R})| d\tilde{\mathbf{a}}_{t+h|t} d\tilde{\mathbf{b}}_{t+h|t} \end{aligned}$$

After integrating out over  $\tilde{\mathbf{a}}_{t+h|t}$ , a step analogous to setting  $\tilde{\mathbf{a}}_{t+h|t} = 0$  for point forecasting, we obtain an expression that gives the probability the reconciled bottom level series lies in the region  $\mathcal{B}$ . This corresponds to  $\nu(\mathcal{B})$  in Definition 4.2. To make a valid probability statement about the entire hierarchy we simply use the bottom level probabilistic forecasts together with Definition 3.3.

### Example: Gaussian Distributions

Suppose an unreconciled probabilistic forecast is Gaussian with mean  $\hat{\boldsymbol{\mu}}$  and variance-covariance matrix  $\hat{\boldsymbol{\Sigma}}$ . The subscripts  $t + h|t$  are suppressed for brevity. The unreconciled density

$$f(\hat{\mathbf{y}}) = (2\pi)^{-n/2} |\hat{\boldsymbol{\Sigma}}|^{-1/2} \exp \left\{ -\frac{1}{2} [(\hat{\mathbf{y}} - \hat{\boldsymbol{\mu}})' \hat{\boldsymbol{\Sigma}}^{-1} (\hat{\mathbf{y}} - \hat{\boldsymbol{\mu}})] \right\}$$

After a change in basis

$$f(\tilde{\mathbf{b}}, \tilde{\mathbf{a}}) = (2\pi)^{-\frac{n}{2}} |\hat{\boldsymbol{\Sigma}}_{t+h}|^{-\frac{1}{2}} |(\mathbf{S} \ \mathbf{R})| \exp \left\{ -\frac{1}{2} q \right\},$$

where

$$q = (\mathbf{S}\tilde{\mathbf{b}} + \mathbf{R}\tilde{\mathbf{a}} - \hat{\boldsymbol{\mu}})' \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{S}\tilde{\mathbf{b}} + \mathbf{R}\tilde{\mathbf{a}} - \hat{\boldsymbol{\mu}})$$

The quadratic form  $q$  can be rearranged as

$$\begin{aligned} q &= \left( (S \ R) \begin{pmatrix} \tilde{\mathbf{b}} \\ \tilde{\mathbf{a}} \end{pmatrix} - \hat{\boldsymbol{\mu}} \right)' \hat{\boldsymbol{\Sigma}}^{-1} \left( (S \ R) \begin{pmatrix} \tilde{\mathbf{b}} \\ \tilde{\mathbf{a}} \end{pmatrix} - \hat{\boldsymbol{\mu}} \right), \\ &= \left( \begin{pmatrix} \tilde{\mathbf{b}} \\ \tilde{\mathbf{a}} \end{pmatrix} - (S \ R)^{-1} \hat{\boldsymbol{\mu}}_{t+h} \right)' \left[ (S \ R)^{-1} \hat{\boldsymbol{\Sigma}}_{t+h} \left( (S \ R)^{-1} \right)' \right]^{-1} \left( \begin{pmatrix} \tilde{\mathbf{b}} \\ \tilde{\mathbf{a}} \end{pmatrix} - (S \ R)^{-1} \hat{\boldsymbol{\mu}}_{t+h} \right). \end{aligned}$$

Recall that

$$(S \ R)^{-1} = \begin{pmatrix} (R'_{\perp} S)^{-1} R'_{\perp} \\ (S'_{\perp} R)^{-1} S'_{\perp} \end{pmatrix} := \begin{pmatrix} G \\ H \end{pmatrix}.$$

Then  $q$  can be rearranged further as

$$\begin{aligned} q &= \left[ \begin{pmatrix} \tilde{\mathbf{b}} \\ \tilde{\mathbf{a}} \end{pmatrix} - \begin{pmatrix} G \\ H \end{pmatrix} \hat{\boldsymbol{\mu}}_{t+h} \right]' \left[ \begin{pmatrix} G \\ H \end{pmatrix} \hat{\boldsymbol{\Sigma}}_{t+h} \begin{pmatrix} G \\ H \end{pmatrix}' \right]^{-1} \left[ \begin{pmatrix} \tilde{\mathbf{b}} \\ \tilde{\mathbf{a}} \end{pmatrix} - \begin{pmatrix} G \\ H \end{pmatrix} \hat{\boldsymbol{\mu}}_{t+h} \right] \\ &= \begin{pmatrix} \tilde{\mathbf{b}} - G\hat{\boldsymbol{\mu}} \\ \tilde{\mathbf{a}} - H\hat{\boldsymbol{\mu}} \end{pmatrix}' \left[ \begin{pmatrix} G \\ H \end{pmatrix} \hat{\boldsymbol{\Sigma}}_{t+h} \begin{pmatrix} G \\ H \end{pmatrix}' \right]^{-1} \begin{pmatrix} \tilde{\mathbf{b}} - G\hat{\boldsymbol{\mu}} \\ \tilde{\mathbf{a}} - H\hat{\boldsymbol{\mu}} \end{pmatrix} \end{aligned}$$

Similar manipulations on determinant of the covariance matrix lead to the following expression for the density

$$\begin{aligned} f(\tilde{\mathbf{b}}, \tilde{\mathbf{a}}) &= (2\pi)^{-\frac{n}{2}} \left| \begin{pmatrix} G\hat{\boldsymbol{\Sigma}}G' & G\hat{\boldsymbol{\Sigma}}H' \\ H\hat{\boldsymbol{\Sigma}}G' & H\hat{\boldsymbol{\Sigma}}H' \end{pmatrix} \right|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \begin{pmatrix} \tilde{\mathbf{b}} - G\hat{\boldsymbol{\mu}} \\ \tilde{\mathbf{a}} - H\hat{\boldsymbol{\mu}} \end{pmatrix}' \right. \\ &\quad \left. \begin{pmatrix} G\hat{\boldsymbol{\Sigma}}G' & G\hat{\boldsymbol{\Sigma}}H' \\ H\hat{\boldsymbol{\Sigma}}G' & H\hat{\boldsymbol{\Sigma}}H' \end{pmatrix}^{-1} \begin{pmatrix} \tilde{\mathbf{b}} - G\hat{\boldsymbol{\mu}} \\ \tilde{\mathbf{a}} - H\hat{\boldsymbol{\mu}} \end{pmatrix} \right\}. \end{aligned}$$

Marginalising out  $\tilde{\mathbf{a}}$ , leads to the following bottom level reconciled forecasts.

$$f(\tilde{\mathbf{b}}) = (2\pi)^{-\frac{m}{2}} \left| G\hat{\boldsymbol{\Sigma}}G' \right|^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\tilde{\mathbf{b}} - G\hat{\boldsymbol{\mu}})' (G\hat{\boldsymbol{\Sigma}}G')^{-1} (\tilde{\mathbf{b}} - G\hat{\boldsymbol{\mu}}) \right\}.$$

Which implies that the reconciled probabilistic forecast for the bottom level series is  $\tilde{\mathbf{b}}_{t+h} \sim \mathcal{N}(G\hat{\boldsymbol{\mu}}_{t+h}, G\hat{\boldsymbol{\Sigma}}_{t+h}G')$ . The reconciled probabilistic forecasts for the whole hierarchy follow



a degenerate Gaussian distribution with mean  $SG\hat{\mu}$  and rank deficient covariance matrix  $SG\hat{\Sigma}_{t+h}G'S'$ .

### 4.3 Elliptical Distributions

We now show that the true predictive distribution can be recovered for elliptical distributions via linear reconciliation. Here, for any square matrix  $C$ ,  $C^{1/2}$  and  $C^{-1/2}$  are defined to satisfy  $C^{1/2} (C^{1/2})' = C$  and  $C^{-1/2} (C^{-1/2})' = C^{-1}$ , for example  $C^{1/2}$  may be obtained via the Cholesky or eigenvalue decompositions.

**Theorem 4.3** (Reconciliation for Elliptical Distributions). *Let an unreconciled probabilistic forecast come from the elliptical class with location parameter  $\hat{\mu}$  and scale matrix  $\hat{\Sigma}$ . Let the true predictive distribution of  $\mathbf{y}_{t+h|t}$  also belong to the elliptical class with location parameter  $\mu$  and scale matrix  $\Sigma$ . Then the affine reconciliation mapping  $g(\tilde{\mathbf{y}}) = \mathbf{G}_{opt}\tilde{\mathbf{y}} + \mathbf{d}_{opt}$  with  $\mathbf{G}_{opt} = \mathbf{A}\Sigma^{-1/2}$  and  $\mathbf{d}_{opt} = \mu - \mathbf{S}\mathbf{G}_{opt}\hat{\mu}$  recovers the true predictive density where  $\mathbf{A}$  is any  $m \times n$  matrix such that  $\mathbf{A}\mathbf{A}' = \Omega$  and  $\Omega$  is the true variance covariance matrix of the predictive distribution for the bottom level.*

*Proof.* Since elliptical distributions are closed under affine transformations, and are closed under marginalisation, reconciliation of an elliptical distribution yields an elliptical distribution (although the unreconciled and unreconciled distributions may be different members of the class of elliptical distributions). The scale matrix of the reconciled forecast is given by  $\mathbf{S}\mathbf{G}_{opt}\Sigma\mathbf{G}_{opt}'\mathbf{S}'$  while the location matrix is given by  $\mathbf{S}\mathbf{G}_{opt}\hat{\mu} + \mathbf{d}_{opt}$ . The reconciled scale matrix is

$$\begin{aligned}\tilde{\Sigma}_{opt} &= \mathbf{S}\mathbf{A}\Sigma^{-1/2}\Sigma\left(\Sigma^{-1/2}\right)' \mathbf{A}'\mathbf{S}' \\ &= \mathbf{S}\Omega\mathbf{S}' \\ &= \Sigma\end{aligned}$$

For the choices of  $\mathbf{G}_{opt}$  and  $\mathbf{d}_{opt}$  given above, the reconciled location vector is

$$\begin{aligned}\tilde{\mu}_{opt} &= \mathbf{S}\mathbf{G}_{opt}\hat{\mu} + \mu - \mathbf{S}\mathbf{G}_{opt}\hat{\mu} \\ &= \mu\end{aligned}$$

□

A number of insights can be drawn from this theorem. First, although a linear mapping  $g(\cdot)$  can be used to recover the true density in the elliptical case, the same does not hold in general. Second,  $g(\cdot)$  is not, in general, a projection matrix. The conditions for which the true predictive density can be recovered by reconciliation are given below.

**Theorem 4.4** (True predictive via projection). *Assume that the true predictive distribution is elliptical with location  $\mu$  and scale  $\Sigma$ . Consider reconciliation via a projection  $g(y) = (R'_{\perp} S)^{-1} R'_{\perp} y$ . The true predictive distribution can be recovered via reconciliation of an elliptical distribution with location  $\hat{\mu}$  and scale  $\hat{\Sigma}$  when the following conditions hold.*

$$\begin{aligned} sp(\hat{\mu} - \mu) &\subset sp(R) \\ sp(\hat{\Sigma}^{1/2} - \Sigma^{1/2}) &\subset sp(R) \end{aligned}$$

*Proof.* The reconciled location vector will be given by

$$\begin{aligned} \tilde{\mu} &= S(R'_{\perp} S)^{-1} R'_{\perp} \hat{\mu} \\ &= S(R'_{\perp} S)^{-1} R'_{\perp} (\hat{\mu} + \mu - \mu) \\ &= S(R'_{\perp} S)^{-1} R'_{\perp} \mu + S(R'_{\perp} S)^{-1} R'_{\perp} (\hat{\mu} - \mu) \end{aligned}$$

Since  $S(R'_{\perp} S)^{-1} R'_{\perp}$  is a projection onto  $\mathfrak{s}$  and  $\mu \in \mathfrak{s}$  the first term simplified to  $\mu$ . If  $\mu - \hat{\mu}$  lies in the span of  $R$  then multiplication by  $R'_{\perp}$  reduced the second term to  $0$ . By a similar argument it can be shown that  $\hat{\Sigma}^{1/2} = \Sigma^{1/2}$ . The closure property of elliptical distributions under affine transformations ensures that the full true predictive distribution can be recovered.  $\square$

Although these conditions will rarely hold in practice and only apply to a limited class of distributions they do provide some insight into selecting a projection for reconciliation. If the value of  $\hat{\mu}$  were equi-probable in all directions then a projection orthogonal to  $\mathfrak{s}$  would be a sensible choice for  $R$  since it would in some sense represent a ‘median’ direction for  $\mu - \hat{\mu}$ . However, the one step ahead in-sample errors are usually correlated suggesting that  $\hat{\mu}$  is more likely to fall in some directions than others. As such an orthogonal projection after transformation by the inverse of the one step ahead in-sample errors may be more intuitively appealing. This is exactly what the MinT projection estimates, and as simulations will show in Section 6, this projection leads to the best empirical results.

## 5 Evaluation of hierarchical probabilistic forecasts

The necessary final step in hierarchical forecasting is to make sure that our forecast distributions are accurate. In general, forecasters prefer to maximize the sharpness of the forecast distribution subject to calibration (Gneiting and Katzfuss, 2014). Therefore the probabilistic forecasts should be evaluated with respect to these two properties.

Calibration refers to the statistical compatibility between probabilistic forecasts and realizations. In other words, random draws from a perfectly calibrated forecast distribution should be equivalent in distribution to the realizations. On the other hand, sharpness refers to the spread or the concentration of the predictive distributions and it is a property of the forecasts only. The more concentrated the forecast distributions, the sharper the forecasts (Gneiting et al., 2008). However, independently assessing the calibration and sharpness will not help to properly evaluate the probabilistic forecasts. Therefore we need to assess these properties simultaneously using scoring rules.

Scoring rules are summary measures obtained based on the relationship between the forecast distributions and the realizations. In some studies, researchers take the scoring rules to be positively oriented, in which case the scores should be maximized (Gneiting and Raftery, 2007). However, scoring rules have also been defined to be negatively oriented, and then the scores should be minimized (Gneiting and Katzfuss, 2014). We follow the latter convention here.

Let  $P$  be a forecast distribution and let  $Q$  be the true data generating process respectively. Furthermore let  $\omega$  be a realization from  $Q$ . Then a scoring rule is a function  $S(P, \omega)$  that maps  $P, \omega$  to  $\mathbb{R}$ . It is a “proper” scoring rule if

$$E_Q[S(P, \omega)] \leq E_Q[S(Q, \omega)], \quad (1)$$

where  $E_Q[S(P, \omega)]$  is the expected score under the true distribution  $Q$  (Gneiting et al., 2008; Gneiting and Katzfuss, 2014). When this inequality is strict, the scoring rule is said to be strictly proper.

In the context of probabilistic forecast reconciliation there could be two motivations for using scoring rules. The first is to compare unreconciled densities to reconciled densities. Although reconciliation is a valuable goal in and of itself since it can be important in aligning decision making across, for example, different units of an enterprise, in the point forecasting literature, forecast reconciliation has also been shown to improve forecast performance . It will be worth-

while to see whether the same holds in the probabilistic forecasting case. The second motivation for using scoring rules is to compare two or more sets of reconciled probabilistic forecasts to one another. The objective here is to evaluate which reconciliation mapping  $g(\cdot)$  works best in practice.

## 5.1 Univariate Scoring rules

One way to evaluate probabilistic forecasts is via the application of univariate scoring rules to each variable in a hierarchy. A summary can be taken of the expected scores across each margin for example a mean or median. In the simulations of section 6 we consider two scoring rules. The log score is given by the log of the marginal density of each variable. The cumulative rank probability score generalises mean square error and is given by

$$\begin{aligned} \text{CRPS}(\check{F}_i, y_{T+h,i}) &= \int (\check{F}_i(\check{y}_i) - \mathbb{1}(y_i < y_{T+h,i})) dy_i \\ &= E_{\check{Y}_i} |\check{Y}_{T+h,i} - y_{T+h,i}| - \frac{1}{2} E_{\check{Y}_i} |\check{Y}_{T+h,i} - \check{Y}_{T+h,i}^*|, \end{aligned}$$

where  $\check{Y}_i$  and  $\check{Y}_{T+h,i}^*$  are independent copies of a random variable with distribution  $\check{F}_i$  and the latter is the predictive distribution for the  $i^{\text{th}}$  margin of a forecast for time  $t + h$  made at time  $t$ . The expectations in the second line can be approximated by Monte Carlo when a sample from the predictive distribution is available.

An advantage to this approach is that it allows the forecaster to evaluate the levels and individual series of the hierarchy where the gains to reconciliation are greatest. For this reason this approach has been used in the limited literature on probabilistic forecasting for hierarchies Ben Taieb et al., 2017 and Jeon et al to date. A major shortcomings to this approach however is that, evaluating univariate scores on the margins do not account for the dependence in the hierarchy.

## 5.2 Multivariate Scoring rules

While the a number of alternative proper scoring rules are available for univariate forecasts, the multivariate case is somewhat more limited. Here we focus on three scoring rules: the log score, the energy score and the variogram score. These are summarized in table 2.

The log score can be approximated using a sample of values from the probabilistic forecast density (Jordan, Krüger, and Lerch, 2017) however it is more commonly used when a parametric form for the density is available for the probabilistic forecast. So far, we have mainly defined probabilistic forecasts in terms of probability measures. Although densities can be obtained for

Jeon  
Pana-  
giotelis  
and  
Petropou-  
los  
refer-  
ence

both reconciled and unreconciled forecasts, the degeneracy of reconciled forecasts is problematic when using log scores. We will discuss this further in the next subsection.

The energy score on the other hand can be defined in terms of the characteristic function of the probabilistic forecast, but the representation in Table 2 in terms of expectations leads itself to easy computation when samples from the probabilistic forecast are available. An interesting case is where  $\alpha = 2$ , where it can be easily shown that

$$\text{ES}(\mathbf{Y}_{T+h}, \check{\mathbf{y}}_{T+h}) = \|\mathbf{y}_{T+h} - \check{\boldsymbol{\mu}}_{T+h}\|^2,$$

where  $\check{\boldsymbol{\mu}}_{T+h} = \mathbb{E}_F(\check{\mathbf{Y}}_{T+h})$ . In this limiting case, the energy score only measures the accuracy of the forecast mean, and not the entire distribution and the energy score is proper, but not strictly proper. Pinson and Tastu (2013) also argues that the energy score has very low discriminative ability for incorrectly specified covariances, even though it discriminates the misspecified means well.

In contrast, Scheuerer and Hamill (2015) have shown that the variogram score has a higher discrimination ability of misspecified means, variances and correlation structures than the energy score. For a finite sample of size  $B$  from the multivariate forecast density  $\check{\mathbf{F}}$ , the empirical variogram score is defined as

$$\text{VS}(\check{\mathbf{F}}, \mathbf{y}_{T+h}) = \sum_{i=1}^n \sum_{j=1}^n w_{ij} \left( |\mathbf{y}_{T+h,i} - \mathbf{y}_{T+h,j}|^p - \frac{1}{B} \sum_{k=1}^B |\check{\mathbf{Y}}_{T+h,i}^k - \check{\mathbf{Y}}_{T+h,j}^k|^p \right)^2.$$

**Table 2:** Scoring rules to evaluate multivariate forecast densities. Here,  $\check{\mathbf{y}}_{T+h}$  and  $\check{\mathbf{y}}_{T+h}^*$  are two independent random vectors from the coherent forecast distribution  $\check{\mathbf{F}}$  with density function  $\check{f}(\cdot)$  at time  $T + h$ , and  $\mathbf{y}_{T+h}$  is the vector of realizations. Further,  $\check{Y}_{T+h,i}$  and  $\check{Y}_{T+h,j}$  are the  $i$ th and  $j$ th components of the vector  $\check{\mathbf{Y}}_{T+h}$ . The variogram score is given for order  $p$ , where  $w_{ij}$  denote non-negative weights.

Scoring rule	Expression	Reference
Log score	$\text{LS}(\check{\mathbf{F}}, \mathbf{y}_{T+h}) = -\log \check{f}(\mathbf{y}_{T+h})$	Gneiting and Raftery (2007)
Energy score	$\text{ES}(\check{\mathbf{Y}}_{T+h}, \mathbf{y}_{T+h}) = \mathbb{E}_{\check{\mathbf{F}}} \ \check{\mathbf{Y}}_{T+h} - \mathbf{y}_{T+h}\ ^\alpha - \frac{1}{2} \mathbb{E}_{\check{\mathbf{F}}} \ \check{\mathbf{Y}}_{T+h} - \check{\mathbf{Y}}_{T+h}^*\ ^\alpha, \quad \alpha \in (0, 2]$	Gneiting et al. (2008)
Variogram score	$\text{VS}(\check{\mathbf{F}}, \mathbf{y}_{T+h}) = \sum_{i=1}^n \sum_{j=1}^n w_{ij} \left(  \mathbf{y}_{T+h,i} - \mathbf{y}_{T+h,j} ^p - \mathbb{E}_{\check{\mathbf{F}}}  \check{Y}_{T+h,i} - \check{Y}_{T+h,j} ^p \right)^2$	Scheuerer and Hamill (2015)

Scheuerer and Hamill (2015) recommend using  $p = 0.5$ .

### 5.2.1 Comparing Unreconciled Forecasts to Reconciled Forecasts

For both reconciled and unreconciled densities it is possible to obtain a density from the probability measures defined in 3. As such it may seem sensible to compare unreconciled densities to reconciled densities on the basis of log score. The following theorem shows that using the log score may fail in the case of multivariate distributions with a degeneracy.

**Theorem 5.1** (Impropriety of log score). *When the true data generating process is a coherent measure, then the log score is improper with respect to the class of incoherent measures.*

*Proof.* Consider a rotated version of hierarchical time series  $\mathbf{z}_t = \mathbf{U}\mathbf{y}_t$  so that the first  $m$  elements of  $\mathbf{z}_t$  denoted  $\mathbf{z}_t^{(1)}$  are unconstrained, while the remaining  $n - m$  elements denoted  $\mathbf{z}_t^{(2)}$  equal 0 when the aggregation constraints hold. An example of the  $n \times n$   $\mathbf{U}$  could be the matrix of left singular vectors of  $\mathbf{S}$ . For a non-degenerate probability measure on  $\mathbb{R}^n$ , the density is the Radon-Nikodym derivative with respect to the usual Lebesgue measure on  $\mathbb{R}^n$ .

Consider the case where the true density is  $f_1(\mathbf{z}_t^{(1)})\mathbb{1}(\mathbf{z}_t^{(2)})$ , and is compared to an incoherent density is given by  $f_1(\mathbf{z}_t^{(1)})f_2(\mathbf{z}_t^{(2)})$ , where  $f_2$  is highly concentrated around 0 but still a proper density. For example  $f_2$  may be Gaussian with variance  $\sigma^2\mathbf{I}$  with  $\sigma^2 < (2\pi)^{-1}$ . The log score under the true DGP is

$$S(\tilde{f}, \mathbf{z}_t^{(1)}) = -\log f_1(\mathbf{z}_t^{(1)}),$$

while that of the unreconciled density is

$$\begin{aligned} S(\hat{f}, \mathbf{z}_t^{(1)}) &= -\log f_1(\mathbf{z}_t^{(1)}) - f_2(\mathbf{z}_t^{(1)}) \\ &= -\log f_1(\mathbf{z}_t^{(1)}) + \frac{n-m}{2} \log(2\pi\sigma^2) \\ &< -\log f_1(\mathbf{z}_t^{(1)}). \end{aligned}$$

After taking expectations  $ES(f, f) > ES(\hat{f}, f)$ , violating the condition (1) for a proper scoring rule. □

A similar issue also arises when discrete random variables are modelled as if they were continuous, an issue discussed in Section 4.1, page 366 of Gneiting and Raftery, 2007. This implies

that the log score should not be used to evaluate multivariate densities with degeneracies and should be avoided when comparing reconciled and unreconciled probabilistic forecasts.

### 5.2.2 Comparing Reconciled Forecasts to one another

Coherent probabilistic forecasts can be completely characterised in terms of basis series; if a probabilistic forecast is available for the basis series then a probabilistic forecast can be recovered for the entire hierarchy via Definition 3.3. This may suggest that it is adequate to merely compare two coherent forecasts to one another using the basis series only. We now show how this is dependent on the specific scoring rule used.

For the log score, suppose the coherent probabilistic forecast has density  $f(\mathbf{b})$ . The density for the full hierarchy is given by  $f(\mathbf{y}) = f(\mathbf{S}\mathbf{b}) = f(\mathbf{b})J^{-1}$  where  $J = \prod_{j=1}^m \lambda_j$  is a pseudo-determinant of the non-square matrix  $\mathbf{S}$  and  $\lambda_j$  are the non-zero singular values of  $\mathbf{S}$ . Therefore for any coherent density the log score of the full hierarchy differs from the log score for the bottom level series by the term  $\log(J)$ . This term depends only on the structure of the hierarchy and is fixed across different reconciliation methods. Therefore if one method achieves a lower expected log score compared to an alternative method when assessed using the bottom level series, the same ordering is preserved when an assessment is made on the basis of the full hierarchy.

The same property does not hold for all scores in general. For example, energy score can be expressed in terms of expectations of norms. In general, since norms are invariant under orthogonal rotations the energy score is also invariant under orthogonal transformations (Székely and Rizzo, 2013; Gneiting and Raftery, 2007). In the context of two coherent forecasts the same is true of a semi-orthogonal transformation from a lower dimensional basis series to the full hierarchy. However, when  $\mathbf{S}$  is the usual summing matrix, it is not semi-orthogonal. As such the energy score computed on the bottom level series will differ from the energy score computed using the full hierarchy and the ordering of different reconciliation methods may change depending on the basis series used. In this case we recommend computing energy score using the full hierarchy. Although the discussion here is related to energy score, the same logic holds for other multivariate scores that are not invariant to orthogonal rotations, for example the variogram score.

The properties of multivariate scoring rules in the context of evaluating reconciled probabilistic forecasts in Table 3.

	Coherent v Incoherent	Coherent v Coherent
Log Score	Not proper	Ordering preserved if compared using bottom level only
Energy/ Variogram Score	Proper	Full hierarchy should be used

**Table 3:** Summary of properties of scoring rules in the context of reconciled probabilistic forecasts.

## 6 Simulation Study

We now turn our attention to comparing different reconciliation methods in a simulation study where the data is conditionally Gaussian. We choose the Gaussian case due to its analytical tractability which allows for evaluation on the basis of all scoring rules (including the log score). The non-Gaussian case lies beyond the scope of this simulation study, but can be handled by a bootstrapping approach proposed in separate work.

For the data generating process, we consider the hierarchy given in Figure 1, comprising two aggregation levels with four bottom-level series. Each bottom-level series will be generated first, and then summed to obtain the data for the upper-level series. In practice, hierarchical time series tend to contain much noisier series at lower levels of aggregation. In order to replicate this feature in our simulations, we follow the data generating process proposed by Wickramasuriya, Athanasopoulos, and Hyndman (2018).

First  $\{w_{AA,t}, w_{AB,t}, w_{BA,t}, w_{BB,t}\}$  are generated from  $\text{ARIMA}(p, d, q)$  processes, where  $(p, q)$  and  $d$  take integers from  $\{1, 2\}$  and  $\{0, 1\}$  respectively with equal probability. The errors driving these ARIMA processes denoted  $\varepsilon$  are jointly normal  $\{\varepsilon_{AA,t}, \varepsilon_{AB,t}, \varepsilon_{BA,t}, \varepsilon_{BB,t}\} \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}) \forall t$ . The parameters for the AR and MA components are randomly and uniformly generated from  $[0.3, 0.5]$  and  $[0.3, 0.7]$  respectively. Then the bottom-level series  $\{y_{AA,t}, y_{AB,t}, y_{BA,t}, y_{BB,t}\}$  are given by:

$$y_{AA,t} = w_{AA,t} + u_t - 0.5v_t,$$

$$y_{AB,t} = w_{AB,t} - u_t - 0.5v_t,$$

$$y_{BA,t} = w_{BA,t} + u_t + 0.5v_t,$$

$$y_{BB,t} = w_{BB,t} - u_t + 0.5v_t,$$



where  $u_t \sim \mathcal{N}(0, \sigma_u^2)$  and  $v_t \sim \mathcal{N}(0, \sigma_v^2)$ . The aggregate series at in the middle level level, are given by:

$$y_{A,t} = w_{AA,t} + w_{AB,t} - v_t,$$

$$y_{B,t} = w_{BA,t} + w_{BB,t} + v_t,$$

and the total series is given by

$$y_{Tot,t} = w_{AA,t} + w_{AB,t} + w_{BA,t} + w_{BB,t}.$$

To ensure noisier disaggregate series than aggregate series, we choose  $\Sigma, \sigma_u^2$  and  $\sigma_v^2$  such that

$$\text{Var}(\varepsilon_{AA,t} + \varepsilon_{AB,t} + \varepsilon_{BA,t} + \varepsilon_{BB,t}) \leq \text{Var}(\varepsilon_{AA,t} + \varepsilon_{AB,t} - v_t) \leq \text{Var}(\varepsilon_{AA,t} + u_t - 0.5v_t),$$

and similar inequalities hold when  $\varepsilon_{AA,t}$  is replaced by  $\varepsilon_{AB,t}$ ,  $\varepsilon_{BA,t}$  and  $\varepsilon_{BB,t}$  in the second and third terms. The values of  $\Sigma$ ,  $\sigma_u^2$  and  $\sigma_v^2$ , that we use and which satisfy these constraints are:

$$\Sigma = \begin{pmatrix} 5.0 & 3.1 & 0.6 & 0.4 \\ 3.1 & 4.0 & 0.9 & 1.4 \\ 0.6 & 0.9 & 2.0 & 1.8 \\ 0.4 & 1.4 & 1.8 & 3.0 \end{pmatrix}, \sigma_u^2 = 19 \text{ and } \sigma_v^2 = 18 \text{ in our simulation setting.}$$

We generate data a sample size of  $T = 501$ . Univariate ARIMA models are selected for each series using the *auto.arima* function in the *forecast* package (Hyndman, 2017) in R (R Core Team, 2018). The same package was used to fit each series independently using the first 500 observations, and evaluate 1-step ahead base (incoherent) probabilistic forecasts. These were then reconciled using different projections summarized in Table 1. This process was replicated using 1000 different data sets from the same data generating processes.

To assess the predictive performance of different forecasting methods, we use scoring rules as discussed in Section 5. To facilitate comparisons, we report skill scores (Gneiting and Raftery, 2007). For a given forecasting method, evaluated by a particular scoring rule  $S(\cdot)$ , the skill score gives the percentage improvement of the preferred forecasting method relative to a reference method. A negative valued skill score indicates that a method is worse than the reference method, whereas any positive value indicates that method is superior to the reference method.

Table 4 summarizes the forecasting performance of unreconciled, bottom-up, OLS, WLS and two MinT reconciliation methods using log score, energy score and variogram score. In all cases

Are the DGPs the same in all cases?

**Table 4:** Comparison of coherent forecasts. “Energy score” and “Variogram score” columns give scores based on the joint forecast distribution of whole hierarchy. “Log score” column gives the log scores of the joint forecast distribution of bottom level. “Skill score” columns give the percentage skill score with reference to the bottom-up method. Entries in these columns show the percentage increase of score for different reconciliation methods relative to the bottom-up method.

Forecasting method	Energy score		Variogram score		Log score	
	Mean score	Skill score %	Mean score	Skill score %	Mean score	Skill score %
MinT(Shrink)	10.03	18.79	8.44	8.46	11.30	6.22
MinT(Sample)	10.01	18.95	8.41	8.79	11.29	6.31
MinT(WLS)	10.53	14.74	9.02	2.17	12.61	−4.65
OLS	10.53	14.74	8.86	3.09	11.54	4.23
Bottom-up	12.35		9.22		12.05	
Incoherent	11.12		9.53			

skill scores are calculated with the bottom-up method as reference. All log scores are evaluated on the basis of bottom level series only, however these only differ from the log scores for the full hierarchy by a fixed constant. The cell for log score of unreconciled forecasts is left blank since the log score is not proper in this context. Overall, the MinT methods provide the best performance irrespective of the scoring rule, and all methods that reconcile using information at all levels of the forecast improve upon unreconciled forecasts. Bottom up forecasts perform even worse than unreconciled forecasts.

Tables 5 and 6 break down the forecasting performance of different reconciliation methods by considering univariate scores on each individual margin. The log score and CRPS are considered while skill scores are computed with the unreconciled forecast as a reference. When broken down in this fashion, the methods based on MinT perform best for all series and always outperform bottom up and unreconciled forecasts. The same cannot be said for OLS which performs worse than bottom up and incoherent forecasts on every individual series.

Puzzling  
it does  
so bad,  
but  
outper-  
forms  
unrec-  
onciled  
on joint  
evalua-  
tion

**Table 5:** Comparison of incoherent vs coherent forecasts based on the univariate forecast distribution of aggregate series. The “Incoherent” row shows the average scores for incoherent forecasts. Each entry above this row represents the percentage skill score with reference to the incoherent forecasts. These entries show the percentage increase in score for different forecasting methods relative to the incoherent forecasts.

Forecasting method	Total		Series - A		Series - B	
	CRPS	LogS	CRPS	LogS	CRPS	LogS
MinT(Shrink)	0.74	0.00	10.49	3.24	9.16	2.73
MinT(Sample)	0.74	0.00	10.49	3.24	9.16	2.73
MinT(WLS)	−2.96	−2.36	6.10	−4.12	5.66	−3.03
OLS	−9.26	−3.36	7.07	2.06	7.01	1.82
Bottom-up	−91.48	−22.22	−8.05	−2.06	−6.20	−1.82
<i>Incoherent</i>	2.70	2.97	4.10	3.40	3.71	3.30

**Table 6:** Comparison of incoherent vs coherent forecasts based univariate forecast distribution of bottom-level series. The “Incoherent” row shows the average scores for incoherent forecasts.

Forecasting method	Series - AA		Series - AB		Series - BA		Series - BB	
	CRPS	LogS	CRPS	LogS	CRPS	LogS	CRPS	LogS
MinT(Shrink)	7.61	2.43	10.82	3.02	5.93	1.86	7.76	2.47
MinT(Sample)	7.88	2.43	11.08	3.02	6.20	1.86	8.05	2.47
MinT(WLS)	3.53	0.00	6.33	0.60	2.43	−0.62	4.89	0.62
OLS	2.99	0.91	5.28	1.51	2.90	0.62	4.31	1.23
<i>Incoherent</i>	3.68	3.29	3.79	3.31	3.45	3.22	3.48	3.24

## 7 Conclusions

Points to discuss

- Optimal reconciliation likely to depend on
  - Properties of base probabilistic forecast
  - Scoring rule
- Case where no parametric form available
- Cases where  $g$  and maybe even  $s$  are non linear

Although the problem of hierarchical point forecasts is well studied in the literature, there is a lack of attention in the context of probabilistic forecasts. Thus we attempted to fill this gap in the literature by providing substantial theoretical background to the problem. We initially provided rigorous definitions for the coherent point and probabilistic forecasts using the principles of measure theory. Due to the aggregation nature of hierarchy, the probability

density is a degenerate density. Thus the forecast distribution that we opt to find should also lie in a lower dimensional subspace of  $\mathbb{R}^n$ .

As it was well established that the reconciliation outperforms other conventional point forecasting methods in the hierarchical literature, we proposed to use reconciliation in probabilistic framework to obtain coherent degenerate densities. We provided a distinct definition for density forecast reconciliation and how it can be used to reconcile incoherent densities in practice.

Assuming a multivariate Gaussian distribution for the hierarchy, we showed how to obtain reconciled Gaussian forecast densities, utilizing available information in the hierarchy. An extensive Monte Carlo simulation study further showed that the MinT reconciliation method (Wickramasuriya, Athanasopoulos, and Hyndman, [2018](#)) is useful in producing improved coherent probabilistic forecasts at least in the Gaussian framework.

## References

- Ben Taieb, S, Huser, R, Hyndman, RJ, and Genton, MG (2017). Forecasting uncertainty in electricity smart meter data by boosting additive quantile regression. *IEEE Transactions on Smart Grid* **7**(5), 2448–2455.
- Dunn, DM, Williams, WH, and Dechaine, TL (1976). Aggregate Versus Subaggregate Models in Local Area Forecasting. *Journal of American Statistical Association* **71**(353), 68–71.
- Erven, T van and Cugliari, J (2014). *Game-Theoretically Optimal reconciliation of contemporaneous hierarchical time series forecasts*. Ed. by A Antoniadis, X Brossat, and J Poggi, pp. 297–317.
- Fliedner, G (2001). Hierarchical forecasting: issues and use guidelines. *Industrial Management & Data Systems* **101**(1), 5–12.
- Gel, Y, Raftery, AE, and Gneiting, T (2004). Calibrated Probabilistic Mesoscale Weather Field Forecasting. *Journal of the American Statistical Association* **99**(July), 575–583.
- Gneiting, T and Katzfuss, M (2014). Probabilistic Forecasting. *Annual Review of Statistics and Its Application* **1**, 125–151.
- Gneiting, T and Raftery, AE (2007). Strictly Proper Scoring Rules, Prediction, and Estimation. *Journal of the American Statistical Association* **102**(477), 359–378.
- Gneiting, T, Raftery, AE, Westveld, AH, and Goldman, T (2005). Calibrated Probabilistic Forecasting Using Ensemble Model Output Statistics and Minimum CRPS Estimation. *Monthly Weather Review* **133**(5), 1098–1118.
- Gneiting, T and Raftery, AE (2005). Weather\_forecasting\_with\_ensem.PDF. *Science* **310**.5746, 248–249.
- Gneiting, T, Stanberry, LI, Grimit, EP, Held, L, and Johnson, NA (2008). “Assessing probabilistic forecasts of multivariate quantities, with an application to ensemble predictions of surface winds”.
- Gross, CW and Sohl, JE (1990). Disaggregation methods to expedite product line forecasting. *Journal of Forecasting* **9**(3), 233–254.
- Hyndman, R (2017). forecast: Forecasting Functions for Time Series and Linear Models, R package version 8.0. URL: <http://github.com/robjhyndman/forecast>.
- Hyndman, RJ, Ahmed, RA, Athanasopoulos, G, and Shang, HL (2011). Optimal combination forecasts for hierarchical time series. *Computational Statistics and Data Analysis* **55**(9), 2579–2589.
- Hyndman, RJ, Lee, AJ, and Wang, E (2016). Fast computation of reconciled forecasts for hierarchical and grouped time series. *Computational Statistics and Data Analysis* **97**, 16–32.

- Jordan, A, Krüger, F, and Lerch, S (2017). Evaluating probabilistic forecasts with the R package `scoringRules`. arXiv: [1709.04743](https://arxiv.org/abs/1709.04743).
- Kahn, KB (1998). *Revisiting top-down versus bottom-up forecasting*. <http://search.ebscohost.com/login.aspx?direct=true%7B%5C%7Ddb=bth%7B%5C%7DAN=985713%7B%5C%7Dlang=pt-br%7B%5C%7Dsite=ehost-live>.
- Lapide, L (1998). A simple view of top-down vs bottom-up forecasting.pdf. *Journal of Business Forecasting Methods & Systems* **17**, 28–31.
- McSharry, PE, Bouwman, S, and Bloemhof, G (2005). Probabilistic forecasts of the magnitude and timing of peak electricity demand. *IEEE Transactions on Power Systems* **20**(2), 1166–1172.
- Pinson, P and Tastu, J (2013). *Discrimination ability of the Energy score*. Tech. rep. Technical University of Denmark.
- Pinson, P, Madsen, H, Papaefthymiou, G, and Klöckl, B (2009). From Probabilistic Forecasts to Wind Power Production. *Wind Energy* **12**(1), 51–62.
- R Core Team (2018). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. Vienna, Austria. <https://www.R-project.org/>.
- Schäfer, J and 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).
- Scheuerer, M and Hamill, TM (2015). Variogram-Based Proper Scoring Rules for Probabilistic Forecasts of Multivariate Quantities \*. *Monthly Weather Review* **143**(4), 1321–1334.
- Schwarzkopf, AB, Tersine, RJ, and Morris, JS (1988). Top-down versus bottom-up forecasting strategies. *International Journal of Production Research* **26**(11), 1833.
- Székely, GJ and Rizzo, ML (2013). Energy statistics: A class of statistics based on distances. *Journal of Statistical Planning and Inference* **143**(8), 1249–1272.
- Wickramasuriya, SL, Athanasopoulos, G, and Hyndman, RJ (2018). Optimal forecast reconciliation for hierarchical and grouped time series through trace minimization. *J American Statistical Association*. to appear.
- Yao, Q and Brockwell, PJ (2006). Gaussian maximum likelihood estimation for ARMA models. I. Time series. *Journal of Time Series Analysis* **27**(6), 857–875.