

Probabilistic Forecasts for Hierarchical Time Series

Anastasios Panagiotelis*

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

Email: anastasios.panagiotelis@monash.edu
and

Puwasala Gamakumara

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

Email: puwasala.gamakumara@monash.edu
and

George Athanasopoulos

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

Email: george.athanasopoulos@monash.edu
and

Rob J Hyndman

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

Email: rob.hyndman@monash.edu

August 3, 2020

*The authors gratefully acknowledge the support of Australian Research Council Grant DP140103220.
We also thank Professor Mervyn Silvapulle for valuable comments.

Abstract

We develop a framework for prediction of multivariate data that follow some known linear constraints, such as the example where some variables are aggregates of others. This is particularly common when forecasting time series (predicting the future), but also arises in other types of prediction. For point prediction, an increasingly popular technique is reconciliation, whereby predictions are made for all series (so-called ‘base’ predictions) and subsequently adjusted to ensure coherence with the constraints. This paper extends reconciliation from the setting of point prediction to probabilistic prediction. A novel definition of reconciliation is developed and used to construct densities and draw samples from a reconciled probabilistic prediction. In the elliptical case, it is proven that the true predictive distribution can be recovered from reconciliation even when the location and scale matrix of the base prediction are chosen arbitrarily. To find reconciliation weights, an objective function based on scoring rules is optimised. The energy and variogram scores are considered since the log score is improper in the context of comparing unreconciled to reconciled predictions, a result also proved in this paper. To account for the stochastic nature of the energy and variogram scores, optimisation is achieved using stochastic gradient descent. This method is shown to improve base predictions in simulation studies and in an empirical application, particularly when the base prediction models are severely misspecified. When misspecification is not too severe, extending popular reconciliation methods for point prediction can result in a similar performance to score optimisation via stochastic gradient descent. The methods described here are implemented in the *ProbReco* package for R.

Keywords: Scoring Rules, Probabilistic Forecasting, Hierarchical Time Series, Stochastic Gradient Descent.

1 Introduction

Many multivariate prediction problems involve data that follow some linear constraints. For instance, in retail or tourism it is important to forecast demand in individual regions as well as aggregate demand of a whole country. In recent years reconciliation has become an increasingly popular method for handling such problems (see Hyndman & Athanasopoulos 2018, for an overview). Reconciliation involves producing predictions for all variables and making a subsequent adjustment to ensure these adhere to known linear constraints. While this methodology has been extensively developed for point prediction, there is a paucity of literature dealing with probabilistic predictions. This paper develops a formal framework for probabilistic reconciliation, derives theoretical results that allow reconciled probabilistic forecasts to be constructed and evaluated, and proposes an algorithm for optimally reconciling probabilistic forecasts with respect to a proper scoring rule.

Before describing the need for probabilistic reconciliation we briefly review the literature on point forecast¹ reconciliation. Prior to the development of forecast reconciliation, the focus was on finding a subset of variables that could be subsequently aggregated or disaggregated to find forecasts for all series (see Dunn et al. 1976, Gross & Sohl 1990, and references therein). An alternative approach emerged with Athanasopoulos et al. (2009) and Hyndman et al. (2011) who recommended producing forecasts of all series and then adjusting, or ‘reconciling’, these forecasts to be ‘coherent’, i.e. adhere to the aggregation constraints. These papers formulated reconciliation as a regression model, however subsequent work has formulated reconciliation as an optimisation problem where weights are chosen to minimise a loss, such as a weighted squared error (Van Erven & Cugliari 2015, Nystrup et al. 2020), a penalised version thereof (Ben Taieb & Koo 2019), or the trace of the forecast error covariance (Wickramasuriya et al. 2019).

In contrast to the point forecasts, the entire probability distribution of future values provides a full description of the uncertainty associated with the predictions (Abramson & Clemen 1995, Gneiting & Katzfuss 2014). Therefore probabilistic forecasting has become of great interest in many disciplines such as, economics (Zarnowitz & Lambros 1987, Rossi

¹Such has been the dominance of forecasting in the literature on reconciliation, that we will refer to forecasting throughout the remainder of the paper. However, we note that the techniques discussed throughout the paper generalise to prediction problems in general and are not limited to time series.

2014), meteorological studies (Pinson et al. 2009, McLean Sloughter et al. 2013), energy forecasting (Wytock & Kolter 2013, Ben Taieb et al. 2017) and retail forecasting (Böse et al. 2017). An early attempt towards probabilistic forecast reconciliation came from Shang & Hyndman (2017) who applied reconciliation to forecast quantiles, rather than to the point forecasts, in order to construct prediction intervals. This idea was extended to constructing a full probabilistic forecast by Jeon et al. (2019) who propose a number of algorithms, one of which is equivalent to reconciling a large number of forecast quantiles. Ben Taieb et al. (2020) also propose an algorithm to obtain probabilistic forecasts that cohere to linear constraints. In particular, Ben Taieb et al. (2020) draw a sample from the probabilistic forecasts of univariate models for the bottom level data, reorder these to match the empirical copula of residuals, and aggregate these in a bottom-up fashion. The only sense in which top level forecasts are used is in the mean, which is adjusted to match that obtained using the MinT reconciliation method (Wickramasuriya et al. 2019).

There are a number of shortcomings to Jeon et al. (2019) and Ben Taieb et al. (2020) which to the best of our knowledge represent the only attempts to develop algorithms for probabilistic forecast reconciliation. First, little formal justification is provided for the algorithms, or for the sense in which they generalise forecast reconciliation to the probabilistic domain. As such, both algorithms are based on sampling and neither can be used to obtain a reconciled density analytically. Both algorithms are tailored towards specific applications and conflate reconciliation with steps that involve reordering the base forecasts. For example, while Jeon et al. (2019) show that ranking draws from independent base probabilistic forecasts before reconciliation is effective, this may only be true due to the highly dependent time series considered in their application. A limitation of Ben Taieb et al. (2020) is that to ensure their sample from the base probabilistic forecast has the same empirical copula as the data, it must be of the same size as the training data. This will be problematic in applications with fewer observations than the smart meter data they consider. Further, Ben Taieb et al. (2020) only incorporate information from the forecast mean of aggregate variables, missing out on potentially valuable information in the probabilistic forecasts of aggregate data.

In this paper we seek to address a number of open issues in probabilistic forecast reconciliation. First, we develop in a formal way, definitions and a framework that generalise

reconciliation from the point setting to the probabilistic setting. This is achieved by extending the geometric framework proposed by Panagiotelis et al. (2020) for point forecast reconciliation. Second, we utilise these definitions to show how a reconciled forecast can be constructed from an arbitrary base forecast. Solutions are provided in the case where a density of the base probabilistic forecast is available and in the case where it is only possible to draw a sample from the base forecasting distribution. Third, we show that in the elliptical case, the correct predictive distribution can be recovered via linear reconciliation irrespective of the location and scale parameters of the base forecasts. We also derive conditions for when this also holds for the special case of reconciliation via projection. Fourth, we derive theoretical results on the evaluation of reconciled probabilistic forecasts using multivariate scoring rules, including showing that the log score is improper when used to compare reconciled to unreconciled forecasts. Fifth, we propose an algorithm for choosing reconciliation weights by optimising a scoring rule. This algorithm takes advantages of advances in stochastic gradient descent and is thus suited to scoring rules that are themselves often only known up to an approximation. The algorithm and other methodological contributions described in this paper are implemented in the *ProbReco* package (Panagiotelis 2020).

The remainder of the paper is structured as follows. In Section 2, after a brief review of point forecast reconciliation, novel definitions are provided for coherent forecasts and reconciliation in the probabilistic setting. In Section 3, we outline how reconciliation can be achieved in the both the case where the density of the base probabilistic forecast is available, and in the case where a sample has been generated from the base probabilistic forecast. In Section 4, we consider the evaluation of probabilistic hierarchical forecasts via scoring rules, including theoretical results on the impropriety of the log score in the context of forecast reconciliation. The use of scoring rules motivates our algorithm for finding optimal reconciliation weights using stochastic gradient descent, which is described in Section 5 and evaluated in an extensive simulation study in Section 6. An empirical application on forecasting electricity generation from different sources is contained in Section 7. Finally Section 8 concludes with some discussion and thoughts on future research.

2 Hierarchical probabilistic forecasts

Before introducing coherence and reconciliation to the probabilistic setting, we first briefly refresh these concepts in the case of point forecasts. In doing so, we follow the geometric interpretation introduced by Panagiotelis et al. (2020), since this formulation naturally generalises to probabilistic forecasting.

2.1 Point Forecasting

A *hierarchical time series* is a collection of time series adhering to some known linear constraints. Stacking the value of each series at time t into an n -vector \mathbf{y}_t , the constraints imply that \mathbf{y}_t lies in an m -dimensional linear subspace of \mathbb{R}^n for all t . This subspace is referred to as the *coherent subspace* and is denoted as \mathfrak{s} . A typical (and the original) motivating example is a collection of time series some of which are aggregates of other series. In this case $\mathbf{b}_t \in \mathbb{R}^m$ can be defined as the values of the most disaggregated or *bottom-level series* at time t and the aggregation constraints can be formulated as,

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

where \mathbf{S} is an $n \times m$ constant matrix for a given hierarchical structure.

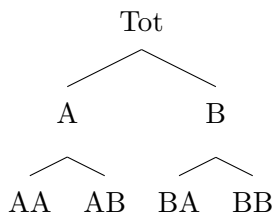


Figure 1: An example of a two-level hierarchical structure.

An example of a hierarchy is shown in Figure 1. There are $n = 7$ series of which $m = 4$ are bottom-level series. Also, $\mathbf{b}_t = [y_{AA,t}, y_{AB,t}, y_{BA,t}, y_{BB,t}]'$, $\mathbf{y}_t = [y_{Tot,t}, y_{A,t}, y_{B,t}, \mathbf{b}_t']'$, and

$$\mathbf{S} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ \mathbf{I}_4 \end{pmatrix},$$

where \mathbf{I}_4 is the 4×4 identity matrix.

The connection between this characterisation and the coherent subspace is that the columns of \mathbf{S} span \mathfrak{s} . Below, the notation $s : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is used when premultiplication by \mathbf{S} is thought of as a mapping. Finally, while \mathbf{S} is defined in terms of m bottom-level series here, in general any m series can be chosen with the \mathbf{S} matrix redefined accordingly. The columns of all appropriately defined \mathbf{S} matrices span the same coherent subspace \mathfrak{s} .

When forecasts of all n series are produced, they may not adhere to constraints. In this case forecasts are called *incoherent base* forecasts and are denoted $\hat{\mathbf{y}}_{t+h}$, with the subscript $t+h$ implying a h -step-ahead forecast at time t . To exploit the fact that the target of the forecast adheres to known linear constraints, these forecasts can be adjusted in a process known as *forecast reconciliation*. At its most general, this involves selecting a mapping $\psi : \mathbb{R}^n \rightarrow \mathfrak{s}$ and then setting $\tilde{\mathbf{y}}_{t+h} = \psi(\hat{\mathbf{y}}_{t+h})$, where $\tilde{\mathbf{y}}_{t+h} \in \mathfrak{s}$ is called the *reconciled* forecast. The mapping ψ may be considered as the composition of two mappings $\psi = s \circ g$. Here, $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ combines incoherent base forecasts of all series to produce new bottom-level forecasts, which are then aggregated via s . Many existing point forecasting approaches including the bottom-up (Dunn et al. 1976), OLS (Hyndman et al. 2011), WLS (Hyndman et al. 2016, Athanasopoulos et al. 2017) and MinT (Wickramasuriya et al. 2019) methods, are special cases where g involves premultiplication by a matrix \mathbf{G} and where \mathbf{SG} is a projection matrix. These are summarised in Table 1.

2.2 Coherent probabilistic forecasts

We now turn our attention towards a novel definition of coherence in a probabilistic setting. First let $(\mathbb{R}^m, \mathcal{F}_{\mathbb{R}^m}, \nu)$ be a probability triple, where $\mathcal{F}_{\mathbb{R}^m}$ is the usual Borel σ -algebra on \mathbb{R}^m . This triple can be thought of as a probabilistic forecast for the bottom-level series. A σ -algebra $\mathcal{F}_{\mathfrak{s}}$ can then be constructed as the collection of sets $s(\mathcal{B})$ for all $\mathcal{B} \in \mathcal{F}_{\mathbb{R}^m}$, where $s(\mathcal{B})$ denotes the image of \mathcal{B} under the mapping s .

Definition 2.1 (Coherent Probabilistic Forecasts). Given the triple, $(\mathbb{R}^m, \mathcal{F}_{\mathbb{R}^m}, \nu)$, a coherent probability triple $(\mathfrak{s}, \mathcal{F}_{\mathfrak{s}}, \check{\nu})$, is given by \mathfrak{s} , the σ -algebra $\mathcal{F}_{\mathfrak{s}}$ and a measure $\check{\nu}$, such that

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

Table 1: Summary of reconciliation methods for which \mathbf{SG} is a projection matrix. Here \mathbf{W} is some diagonal matrix, $\hat{\Sigma}_{sam}$ is a sample estimate of the residual covariance matrix and $\hat{\Sigma}_{shr}$ is a shrinkage estimator proposed by Schäfer & Strimmer (2005), given by $\tau \text{diag}(\hat{\Sigma}_{sam}) + (1 - \tau)\hat{\Sigma}_{sam}$ where $\tau = \frac{\sum_{i \neq j} \text{Var}(\hat{\sigma}_{ij})}{\sum_{i \neq j} \hat{\sigma}_{ij}^2}$ and σ_{ij} denotes the (i, j) th element of $\hat{\Sigma}_{sam}$.

Reconciliation method	\mathbf{G}
OLS	$(\mathbf{S}'\mathbf{S})^{-1}\mathbf{S}'$
WLS	$(\mathbf{S}'\mathbf{W}\mathbf{S})^{-1}\mathbf{S}'\mathbf{W}$
MinT(Sample)	$(\mathbf{S}'\hat{\Sigma}_{sam}^{-1}\mathbf{S})^{-1}\mathbf{S}'\hat{\Sigma}_{sam}^{-1}$
MinT(Shrink)	$(\mathbf{S}'\hat{\Sigma}_{shr}^{-1}\mathbf{S})^{-1}\mathbf{S}'\hat{\Sigma}_{shr}^{-1}$

To the best of our knowledge, the only other definition of coherent probabilistic forecasts is given by Ben Taieb et al. (2020) who define them in terms of convolutions. While these definitions do not contradict one another, our definition has two advantages. First it can more naturally be extended to problems with non-linear constraints with the coherent subspace \mathfrak{s} replaced with a manifold. Second, it facilitates a definition of probabilistic forecast reconciliation to which we now turn our attention.

2.3 Probabilistic forecast reconciliation

Let $(\mathbb{R}^n, \mathcal{F}_{\mathbb{R}^n}, \hat{\nu})$ be a probability triple characterising a probabilistic forecast for all n series. The hat is used for $\hat{\nu}$ analogously with $\hat{\mathbf{y}}$ in the point forecasting case. The objective is to derive a reconciled measure $\tilde{\nu}$, assigning probability to each element of the σ -algebra $\mathcal{F}_{\mathfrak{s}}$.

Definition 2.2. The reconciled probability measure of $\hat{\nu}$ with respect to the mapping $\psi(\cdot)$ is a probability measure $\tilde{\nu}$ on \mathfrak{s} with σ -algebra $\mathcal{F}_{\mathfrak{s}}$ such that

$$\tilde{\nu}(\mathcal{A}) = \hat{\nu}(\psi^{-1}(\mathcal{A})) \quad \forall \mathcal{A} \in \mathcal{F}_{\mathfrak{s}},$$

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

This definition naturally extends forecast reconciliation to the probabilistic setting. In the point forecasting case, the reconciled forecast is obtained by passing an incoherent

forecast through a transformation. Similarly, for probabilistic forecasts, sets of points is mapped to sets of points by a transformation. The same probabilities are assigned to these sets under the base and reconciled measures respectively. Recall that the mapping ψ can also be expressed as a composition of two transformations $s \circ g$. In this case, an m -dimensional reconciled probabilistic distribution ν can be obtained such that $\nu(\mathcal{B}) = \hat{\nu}(g^{-1}(\mathcal{B}))$ for all $\mathcal{B} \in \mathcal{F}_{\mathbb{R}^m}$ and a probabilistic forecast for the full hierarchy can then be obtained via Definition 2.1. This construction will be used in Section 3.

Definition 2.2 can use any continuous mapping ψ , where continuity is required to ensure that open sets in \mathbb{R}^n used to construct $\mathcal{F}_{\mathbb{R}^n}$ are mapped to open sets in \mathfrak{s} . However, hereafter, we restrict our attention to ψ as a linear mapping. This is depicted in Figure 2 when ψ is a projection. This figure is only a schematic, since even the most trivial hierarchy is 3-dimensional. The arrow labelled \mathbf{S} spans an m -dimensional coherent subspace \mathfrak{s} , while the arrow labelled \mathbf{R} spans an $(n - m)$ -dimensional direction of projection. The mapping g collapses all points in the blue shaded region $g^{-1}(\mathcal{B})$, to the black interval \mathcal{B} . Under s , \mathcal{B} is mapped to $s(\mathcal{B})$ shown in red. Under our definition of reconciliation, the same probability is assigned to the red region under the reconciled measure as is assigned to the blue region under the incoherent measure.

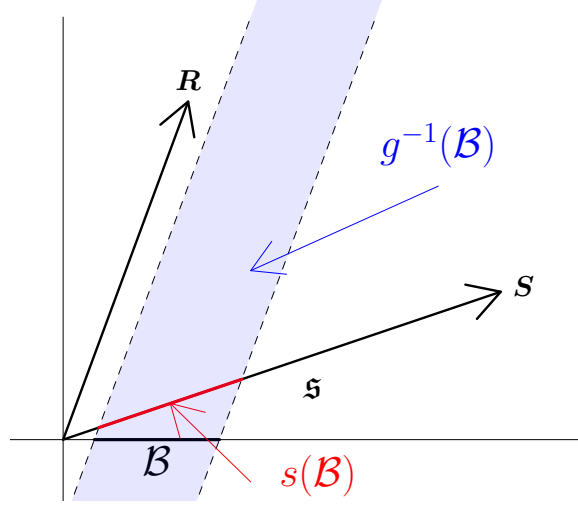


Figure 2: *Summary of probabilistic forecast reconciliation. The probability that \mathbf{y}_{t+h} lies in the red line segment under the reconciled probabilistic forecast is defined to be equal to the probability that \mathbf{y}_{t+h} lies in the shaded blue area under the un-reconciled probabilistic forecast. Note that since the smallest possible hierarchy involves three dimensions, this figure is only a schematic.*

3 Construction of Reconciled Distribution

In this section we derive theoretical results on how distributions on \mathbb{R}^n can be reconciled to a distribution on \mathfrak{s} . In Section 3.1 we show how this can be achieved analytically by a change of coordinates and marginalisation when the density is available. In Section 3.2 we explore this result further in the specific case of elliptical distributions. In Section 3.3 we consider reconciliation in the case where the density may be unavailable but it is possible to draw a sample from the base probabilistic forecast distribution. Throughout we restrict our attention to linear reconciliation.

3.1 Analytical derivation of reconciled densities

The following theorem shows how a reconciled density can be derived from any base probabilistic forecast on \mathbb{R}^n .

Theorem 3.1 (Reconciled density of bottom-level). *Consider the case where reconciliation is carried out using a composition of linear mappings $s \circ g$ where g combines information*

from all levels of the base forecast into a new density for the bottom-level. The density of the bottom-level series under the reconciled distribution is

$$\tilde{f}_{\mathbf{b}}(\mathbf{b}) = |\mathbf{G}^*| \int \hat{f}(\mathbf{G}^-\mathbf{b} + \mathbf{G}_{\perp}\mathbf{a})d\mathbf{a},$$

where \hat{f} is the density of the incoherent base probabilistic forecast, \mathbf{G}^- is an $n \times m$ generalised inverse of \mathbf{G} such that $\mathbf{G}\mathbf{G}^- = \mathbf{I}$, \mathbf{G}_{\perp} is an $n \times (n - m)$ orthogonal complement to \mathbf{G} such that $\mathbf{G}\mathbf{G}_{\perp} = \mathbf{0}$, $\mathbf{G}^* = \begin{pmatrix} \mathbf{G}^- \\ \mathbf{G}_{\perp} \end{pmatrix}$, and \mathbf{b} and \mathbf{a} are obtained via the change of variables

$$\mathbf{y} = \mathbf{G}^* \begin{pmatrix} \mathbf{b} \\ \mathbf{a} \end{pmatrix}.$$

Proof. See Appendix A. □

Theorem 3.2 (Reconciled density of full hierarchy). *Consider the case where a reconciled density for the bottom-level series has been obtained using Theorem 3.1. The density of the full hierarchy under the reconciled distribution is*

$$\tilde{f}_{\mathbf{y}}(\mathbf{y}) = |\mathbf{S}^*| \tilde{f}_{\mathbf{b}}(\mathbf{S}^-\mathbf{y}) \mathbb{1}\{\mathbf{y} \in \mathfrak{s}\},$$

where $\mathbb{1}\{\cdot\}$ equals 1 when the statement in braces is true and 0 otherwise,

$$\mathbf{S}^* = \begin{pmatrix} \mathbf{S}^- \\ \mathbf{S}'_{\perp} \end{pmatrix},$$

\mathbf{S}^- is an $m \times n$ generalised inverse of \mathbf{S} such that $\mathbf{S}^-\mathbf{S} = \mathbf{I}$, and \mathbf{S}_{\perp} is an $n \times (n - m)$ orthogonal complement to \mathbf{S} such that $\mathbf{S}'_{\perp}\mathbf{S} = \mathbf{0}$.

Proof. See Appendix A. □

Example: Gaussian Distribution

Suppose the incoherent base forecasts are Gaussian with mean $\hat{\boldsymbol{\mu}}$, covariance matrix $\hat{\boldsymbol{\Sigma}}$ and density,

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

Then, using Theorem 3.1, the reconciled density for the bottom-level series is given by

$$\tilde{f}_{\mathbf{b}}(\mathbf{b}) = \int (2\pi)^{-\frac{n}{2}} |\hat{\boldsymbol{\Sigma}}|^{-\frac{1}{2}} |\mathbf{G}^*| e^{q/2} d\mathbf{a},$$

where

$$\begin{aligned} q &= \left[\mathbf{G}^* \begin{pmatrix} \mathbf{b} \\ \mathbf{a} \end{pmatrix} - \hat{\boldsymbol{\mu}} \right]' \hat{\boldsymbol{\Sigma}}^{-1} \left[\mathbf{G}^* \begin{pmatrix} \mathbf{b} \\ \mathbf{a} \end{pmatrix} - \hat{\boldsymbol{\mu}} \right] \\ &= \left[\begin{pmatrix} \mathbf{b} \\ \mathbf{a} \end{pmatrix} - \mathbf{G}^{*-1} \hat{\boldsymbol{\mu}} \right]' \left[\mathbf{G}^{*-1} \hat{\boldsymbol{\Sigma}} (\mathbf{G}^{*-1})' \right]^{-1} \left[\begin{pmatrix} \mathbf{b} \\ \mathbf{a} \end{pmatrix} - \mathbf{G}^{*-1} \hat{\boldsymbol{\mu}} \right]. \end{aligned}$$

Noting that

$$\mathbf{G}^{*-1} = \begin{pmatrix} \mathbf{G} \\ \mathbf{G}_{\perp}^{-} \end{pmatrix},$$

where \mathbf{G}_{\perp}^{-} is an $(n - m) \times n$ matrix such that $\mathbf{G}_{\perp}^{-} \mathbf{G}_{\perp} = \mathbf{I}$, q can be rearranged as

$$\left[\begin{pmatrix} \mathbf{b} \\ \mathbf{a} \end{pmatrix} - \begin{pmatrix} \mathbf{G} \\ \mathbf{G}_{\perp}^{-} \end{pmatrix} \hat{\boldsymbol{\mu}} \right]' \left[\begin{pmatrix} \mathbf{G} \\ \mathbf{G}_{\perp}^{-} \end{pmatrix} \hat{\boldsymbol{\Sigma}} \begin{pmatrix} \mathbf{G} \\ \mathbf{G}_{\perp}^{-} \end{pmatrix}' \right]^{-1} \left[\begin{pmatrix} \mathbf{b} \\ \mathbf{a} \end{pmatrix} - \begin{pmatrix} \mathbf{G} \\ \mathbf{G}_{\perp}^{-} \end{pmatrix} \hat{\boldsymbol{\mu}} \right].$$

After the change of variables, the density can be recognised as a multivariate Gaussian in \mathbf{b} and \mathbf{a} . The mean and covariance matrix for the margins of the first m elements are $\mathbf{G}\hat{\boldsymbol{\mu}}$ and $\mathbf{G}\hat{\boldsymbol{\Sigma}}\mathbf{G}'$ respectively. Marginalising out \mathbf{a} , the reconciled forecast for the bottom-level is $\tilde{\mathbf{b}} \sim \mathcal{N}(\mathbf{G}\hat{\boldsymbol{\mu}}, \mathbf{G}\hat{\boldsymbol{\Sigma}}\mathbf{G}')$. Using standard results from matrix algebra of normals, $\tilde{\mathbf{y}} \sim \mathcal{N}(\mathbf{S}\mathbf{G}\hat{\boldsymbol{\mu}}, \mathbf{S}\mathbf{G}\hat{\boldsymbol{\Sigma}}\mathbf{G}'\mathbf{S}')$.

3.2 Elliptical distributions

More generally, consider linear reconciliation of the form $\psi(\hat{\mathbf{y}}) = \mathbf{S}(\mathbf{d} + \mathbf{G}\hat{\mathbf{y}})$. For an elliptical base probabilistic forecast, with location $\hat{\boldsymbol{\mu}}$ and scale $\hat{\boldsymbol{\Sigma}}$, the reconciled probabilistic forecast will also be elliptical with location $\tilde{\boldsymbol{\mu}} = \mathbf{S}(\mathbf{d} + \mathbf{G}\hat{\boldsymbol{\mu}})$ and scale $\tilde{\boldsymbol{\Sigma}} = \mathbf{S}\mathbf{G}\hat{\boldsymbol{\Sigma}}\mathbf{G}'\mathbf{S}'$. This is a consequence of the fact that elliptical distributions are closed under linear transformations and marginalisation. While the base and reconciled distribution may be of a different form, they will both belong to the elliptical family. This leads to the following result.

Theorem 3.3 (Recovering the true density through reconciliation). *Assume the true predictive distribution is elliptical with location $\boldsymbol{\mu}$ and scale $\boldsymbol{\Sigma}$. Then for an elliptical base probabilistic forecast with arbitrary location $\hat{\boldsymbol{\mu}}$ and scale $\hat{\boldsymbol{\Sigma}}$, there exists \mathbf{d}_{opt} and \mathbf{G}_{opt} such that the true predictive distribution is recovered by reconciliation.*

Proof. First consider finding a \mathbf{G}_{opt} for which the following holds,

$$\boldsymbol{\Sigma} = \mathbf{S}\mathbf{G}_{\text{opt}}\hat{\boldsymbol{\Sigma}}\mathbf{G}_{\text{opt}}'\mathbf{S}'.$$

This can be solved as $\mathbf{G}_{\text{opt}} = \mathbf{\Omega}_0^{1/2} \hat{\mathbf{\Sigma}}^{-1/2}$, where $\hat{\mathbf{\Sigma}}^{1/2}$ is any matrix such that $\hat{\mathbf{\Sigma}} = \hat{\mathbf{\Sigma}}^{1/2}(\hat{\mathbf{\Sigma}}^{1/2})'$ (for example a Cholesky factor), $\mathbf{\Omega}_0^{1/2}(\mathbf{\Omega}_0^{1/2})' = \mathbf{\Omega}$ and $\mathbf{\Omega}$ is the true scale matrix for the bottom-level series. To ensure conformability of matrix multiplication, $\mathbf{\Omega}_0^{1/2}$ must be an $m \times n$ matrix; so it can be set to the Cholesky factor of $\mathbf{\Omega}$ augmented with an additional $n - m$ columns of zeros. To reconcile the location, solve the following for \mathbf{d}_{opt}

$$\boldsymbol{\mu} = \mathbf{S}(\mathbf{d}_{\text{opt}} + \mathbf{G}_{\text{opt}}\hat{\boldsymbol{\mu}})$$

which is given by $\mathbf{d}_{\text{opt}} = \boldsymbol{\beta} - \mathbf{G}_{\text{opt}}\hat{\boldsymbol{\mu}}$, where $\boldsymbol{\beta}$ is defined so that $\boldsymbol{\mu} = \mathbf{S}\boldsymbol{\beta}$. \square

While the above theorem is not feasible in practice (exploiting the result requires knowledge of $\boldsymbol{\mu}$ and $\mathbf{\Sigma}$), it does nonetheless have important consequences for the algorithm that we introduce in Section 5. In particular, note that $\mathbf{S}\mathbf{G}_{\text{opt}}$ is not a projection matrix in general. This implies that in the probabilistic forecasting setting, it is advised to include a translation \mathbf{d} in the reconciliation procedure. This holds even if the base forecasts are unbiased (i.e. $\hat{\boldsymbol{\mu}} = \boldsymbol{\mu}$) since in general $\mathbf{S}\mathbf{G}_{\text{opt}}\hat{\boldsymbol{\mu}} \neq \boldsymbol{\mu}$.

Although $\mathbf{S}\mathbf{G}_{\text{opt}}$ is not a projection matrix in general, there are some conditions under which it will be. These are described by the following theorem.

Theorem 3.4 (Optimal Projection for Reconciliation). *Let $\hat{\mathbf{\Sigma}}$ be the scale matrix from an elliptical but incoherent base forecast and assume base forecasts are also unbiased. When the true predictive distribution is also elliptical, then this can be recovered via reconciliation using a projection if $\text{rank}(\hat{\mathbf{\Sigma}} - \mathbf{\Sigma}) \leq n - m$.*

Proof. See Appendix B. \square

3.3 Simulation from a Reconciled Distribution

In practice it is often the case that samples are drawn from a probabilistic forecast since an analytical expression is either unavailable, or relies on unrealistic parametric assumptions. A useful result is the following.

Theorem 3.5 (Reconciled samples). *Suppose that $(\hat{\mathbf{y}}^{[1]}, \dots, \hat{\mathbf{y}}^{[L]})$ is a sample drawn from an incoherent probability measure $\hat{\nu}$. Then $(\tilde{\mathbf{y}}^{[1]}, \dots, \tilde{\mathbf{y}}^{[L]})$ where $\tilde{\mathbf{y}}^{[\ell]} := \psi(\hat{\mathbf{y}}^{[\ell]})$ for $\ell = 1, \dots, L$, is a sample drawn from the reconciled probability measure $\tilde{\nu}$ as defined in Definition 2.2.*

Proof. For any $\mathcal{A} \in \mathcal{F}_g$

$$\begin{aligned} \Pr(\hat{\mathbf{y}} \in \psi^{-1}(\mathcal{A})) &= \lim_{L \rightarrow \infty} \sum_{\ell=1}^L \mathbb{1}\{\hat{\mathbf{y}}^{[\ell]} \in \psi^{-1}(\mathcal{A})\} \\ &= \lim_{L \rightarrow \infty} \sum_{\ell=1}^L \mathbb{1}\{\psi(\hat{\mathbf{y}}^{[\ell]}) \in (\mathcal{A})\} \\ &= \Pr(\tilde{\mathbf{y}} \in (\mathcal{A})) \end{aligned}$$

□

This result implies that reconciling each member of a sample drawn from an incoherent distribution provides a sample from the reconciled distribution. Such a strategy has already been used by Jeon et al. (2019), without formal justification. This result allows coherent forecasts to be built in a general and modular fashion, the mechanism for simulating base forecasts is separated from the question of reconciliation. This will become clear in the simulation study covered in Section 6.

4 Evaluation of Hierarchical Probabilistic Forecasts

An important issue in all forecasting problems is evaluating forecast accuracy. In the probabilistic setting, it is common to evaluate forecasts using proper scoring rules (see Gneiting & Raftery 2007, Gneiting & Katzfuss 2014, and references therein). Throughout, we follow the convention of negatively oriented scoring rules such that smaller values of the score indicate more accurate forecasts. In general, a scoring rule $K(., .)$, is a function taking a probability measure as the first argument and a realisation as the second argument (although for ease of notation we will at times replace the probability measure with its associated density in the first argument). A scoring rule is *proper* if $\mathbb{E}_Q[K(Q, \omega)] \leq \mathbb{E}_Q[K(P, \omega)]$ for all P , where P is any member of some class of probability measures (densities), Q is the true predictive and ω is a realisation. When this inequality is strict for all $P \neq Q$, the scoring rule is said to be *strictly proper*.

Since hierarchical forecasting is inherently a multivariate problem (the linear constraints affect all variables), our focus is on multivariate scoring rules. Arguably the simplest and most common multivariate scoring rule is the log score. The log score simply involves

evaluating the negative log density at the value of the realisation, $LS(P, \omega) = -\log f(\omega)$, where f is the density associated with a distribution P . The log score is more commonly used when a parametric form for the density is available, however this density can also be approximated from a sample of values drawn from the probabilistic forecast (see Jordan et al. 2017).

Alternatively there are a number of other multivariate scoring rules that are difficult to compute using the probabilistic forecast density alone, but can be approximated using a sample drawn from that density. An example is the energy score (ES) (see Székely 2003, Gneiting & Raftery 2007, for details) which is a multivariate generalisation of the popular Cumulative Rank Probability Score (CRPS). The energy score is given by

$$ES(P, \omega) = E_P \|\mathbf{y} - \omega\|^\alpha - \frac{1}{2} E_P \|\mathbf{y} - \mathbf{y}^*\|^\alpha, \quad \alpha \in (0, 2], \quad (1)$$

where \mathbf{y} and \mathbf{y}^* are independent copies drawn from the distribution P . In the empirical results described later, we follow common convention by setting $\alpha = 0.5$. While the expectations in Equation (1) may have no closed form, they can be easily approximated via simulations using a sample drawn from the probabilistic forecast. Other scores with similar behaviour are kernel-based scores (Dawid 2007, Gneiting & Raftery 2007) and the variogram score (Scheuerer & Hamill 2015).

4.1 The Log Score for Hierarchical Time Series

When an expression for the density of an incoherent base forecast is available, Section 3 describes how the density of a reconciled forecast can be recovered. With both densities available, the log score is a natural and straightforward scoring rule to use. However, the following theorem shows that the log score is improper in the setting of comparing incoherent to coherent forecasts.

Theorem 4.1 (Impropriety of log score). *When the true data generating process is coherent, then the log score is improper with respect to the class of incoherent measures.*

Proof. See Appendix C. □

As a result of Theorem 4.1 we recommend avoiding the log score when comparing reconciled and unreconciled probabilistic forecasts.

Table 2: *Properties of scoring rules for reconciled probabilistic forecasts.*

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

If a probabilistic forecast is available for any m series, then a probabilistic forecast for the full hierarchy can be derived. Definition 2.1 provides an example using the bottom-level series. This suggests that it may be adequate to merely compare two coherent forecasts to one another using the bottom-level series only. This is true for the log score.

Consider a coherent probabilistic forecast with density $\tilde{f}_{\mathbf{y}}$ for the full hierarchy and $\tilde{f}_{\mathbf{b}}$ for the bottom-level series. By Theorem 3.2, $\tilde{f}_{\mathbf{y}}(\mathbf{y}) = |\mathbf{S}^*| \tilde{f}_{\mathbf{b}}(\mathbf{S}^-\mathbf{y}) \mathbb{1}\{\mathbf{y} \in \mathfrak{s}\}$. Any realisation \mathbf{y}^* will lie on the coherent subspace and can be written as $\mathbf{S}\mathbf{b}^*$. The expression for the log score is therefore

$$\begin{aligned} \text{LS}(\tilde{f}_{\mathbf{y}}, \mathbf{y}^*) &= -\log \left(|\mathbf{S}^*| \tilde{f}_{\mathbf{b}}(\mathbf{S}^-\mathbf{S}\mathbf{b}^*) \right) \\ &= -\log |\mathbf{S}^*| - \log \tilde{f}_{\mathbf{b}}(\mathbf{b}^*). \end{aligned}$$

For coherent densities, the log score for the full hierarchy differs from the log score for the bottom-level series only by $-\log |\mathbf{S}^*|$. This term is independent of the choice of \mathbf{G} . Consequently, rankings of different reconciliation methods using the log score for the full hierarchy will not change if only the bottom-level series is used.

The same property does not hold for all scores. For example, the energy score is invariant under orthogonal transformations (Székely & Rizzo 2013) but not under linear transformations in general. Therefore it is possible for one method to outperform another when energy score is calculated using the full hierarchy, but for these rankings to change if only bottom-level series are considered. We therefore recommend computing the energy score using the full hierarchy. The properties of multivariate scoring rules in the context of evaluating reconciled probabilistic forecasts are summarised in Table 2.

5 Score Optimal Reconciliation

We now propose an algorithm for finding reconciliation weights by optimising an objective function based on scores. For clarity of exposition, we consider the special case of the energy score. However, the algorithm can be generalised to any score that is computed by sampling from the probabilistic forecast. For example, in the simulations and the empirical application of Sections 6 and 7 we consider optimising with respect to both the energy and variogram scores. We consider linear reconciliation of the form $\tilde{\mathbf{y}} = \psi_{\gamma}(\hat{\mathbf{y}}) = \mathbf{S}(\mathbf{d} + \mathbf{G}\hat{\mathbf{y}})$, where $\gamma := (\mathbf{d}, \text{vec}(\mathbf{G}))$. This allows for more flexibility than a projection, which would imply the constraints $\mathbf{d} = \mathbf{0}$ and $\mathbf{G}\mathbf{S} = \mathbf{I}$. This added flexibility is motivated by Theorem 3.3 which shows that projections in general are not guaranteed to recover the true predictive distribution even in the elliptical case. When making an h -step-ahead forecast at time T , the objective used to determine an optimal value of γ is the total energy score based on in-sample information, given by

$$\mathcal{E}(\gamma) = \sum_{t=T}^{T+R} ES(\tilde{f}_{t+h|t}^{\gamma}, \mathbf{y}_{t+h}), \quad (2)$$

where $\tilde{f}_{t+h|t}^{\gamma}$ is a probabilistic forecast for \mathbf{y}_{t+h} made at time t and reconciled with respect to $\psi_{\gamma}(\cdot)$, and R is the number of score evaluations used in forming the objective function.

One of the challenges in optimising this objective function is that there is, in general, no closed form expression for the energy score. However, it can be easily approximated by simulation as

$$\hat{\mathcal{E}}(\gamma) = \sum_{t=T}^{T+R} \left[\frac{1}{Q} \left(\sum_{q=1}^Q \|\tilde{\mathbf{y}}_{t+h|t}^{[q]} - \mathbf{y}_{t+h}\| - \frac{1}{2} \|\tilde{\mathbf{y}}_{t+h|t}^{[q]} - \tilde{\mathbf{y}}_{t+h|t}^{*[q]}\| \right) \right], \quad (3)$$

where $\tilde{\mathbf{y}}_{t+h|t}^{[q]} = \mathbf{S}(\mathbf{d} + \mathbf{G}\hat{\mathbf{y}}_{t+h|t}^{[q]})$, $\tilde{\mathbf{y}}_{t+h|t}^{*[q]} = \mathbf{S}(\mathbf{d} + \mathbf{G}\hat{\mathbf{y}}_{t+h|t}^{*[q]})$ and $\hat{\mathbf{y}}_{t+h|t}^{[q]}, \hat{\mathbf{y}}_{t+h|t}^{*[q]} \stackrel{iid}{\sim} \hat{f}_{t+h|t}$ for $q = 1, \dots, Q$, and $t = 1, \dots, T - h$.

The objective function is optimised by Stochastic Gradient Descent (SGD). The SGD technique has become increasingly popular in machine learning and statistics over the past decade having been applied to training neural networks (Bottou 2010) and Variational Bayes (Kingma & Welling 2013). The method requires an estimate of the gradient $\partial \hat{\mathcal{E}} / \partial \gamma$ which is computed by automatically differentiating Equation (3) using the header only C++ library of the Stan project (Carpenter et al. 2015). The learning rates used for SGD

are those of the Adam method (see Kingma & Ba 2014, for details). Pseudo-code for the full procedure in the case where $h = 1$ is provided in Algorithm 1 and is implemented in the R package *ProbReco* (Panagiotelis 2020).

Algorithm 1 SGD with Adam for score optimal reconciliation (one-step-ahead forecasts). The initial value of γ is given by OLS reconciliation. Steps 9–14 are the standard steps for SGD with Adam. Squaring \mathbf{g}_j in Step 11 and division and addition in Step 14 are element-wise operations.

```

1: procedure SCOREOPT( $\hat{f}_{2|1}, \dots, \hat{f}_{T|T-1|1}, \mathbf{y}_1, \dots, \mathbf{y}_T, \beta_1, \beta_2, \epsilon, \eta$ ).
2:   for  $t = T : T + R$  do
3:     Find base forecasts  $\hat{f}_{t+1|t}$  using  $t - T + 1, t - T + 2, \dots, t$  as training data.
4:   end for
5:   Initialise  $\mathbf{m}_0 = \mathbf{0}, \mathbf{v}_0 = \mathbf{0}$  and  $\gamma_0 = (\mathbf{0}, \text{vec}((\mathbf{S}'\mathbf{S})^{-1}\mathbf{S}'))$ 
6:   for  $j = 1, 2, 3, \dots$  up to convergence do
7:     Draw  $\hat{\mathbf{y}}_{t+1|t}^{[q]}, \hat{\mathbf{y}}_{t+1|t}^{*[q]} \sim \hat{f}_{t+1|t}$  for  $q = 1, \dots, Q, t = 1, \dots, T - 1$ .
8:     Compute  $\hat{\mathbf{y}}_{t+1|t}^{[q]}$  and  $\hat{\mathbf{y}}_{t+1|t}^{*[q]}$  for  $q = 1, \dots, Q, t = 1, \dots, T - 1$  using  $\gamma_{j-1}$ .
9:      $\mathbf{g}_j \leftarrow \partial \hat{\mathcal{E}} / \partial \gamma \Big|_{\gamma = \gamma_{j-1}}$  ▷ Compute gradient
10:     $\mathbf{m}_j \leftarrow \beta_1 \mathbf{m}_{j-1} + (1 - \beta_1) \mathbf{g}_j$  ▷ Moving average of gradient
11:     $\mathbf{v}_j \leftarrow \beta_2 \mathbf{v}_{j-1} + (1 - \beta_2) \mathbf{g}_j^2$  ▷ Moving average of squared gradient
12:     $\hat{\mathbf{m}}_j \leftarrow \mathbf{m}_j / (1 - \beta_1^j)$  ▷ Bias correct
13:     $\hat{\mathbf{v}}_j \leftarrow \mathbf{v}_j / (1 - \beta_2^j)$  ▷ Bias correct
14:     $\gamma_j \leftarrow \gamma_{j-1} + \eta \frac{\hat{\mathbf{m}}_j}{(\hat{\mathbf{v}}_j + \epsilon)}$  ▷ Update weights
15:   end for
16:   Set the reconciled forecast as  $\tilde{f}_{T+1|T}^{\gamma_{\text{opt}}}$  where  $\gamma_{\text{opt}}$  is the converged value of  $\gamma$ .
17: end procedure

```

While Algorithm 1 is not the first instance of calibrating parameters by optimising scoring rules (see Gneiting et al. 2005, for an earlier example), to the best of our knowledge it is the first instance of doing so using SGD and the first application to forecast reconciliation. It is amenable to parallel computing architectures: the loop beginning at line 2 of the pseudo-code of Algorithm 1 can be done in parallel as can the computation of the gradient. Finally, the total score in Equation (2) can be replaced with a weighted sum where

appropriate; for instance weights that decay for scores computed further in the past will favour choices of γ that produced better forecasting performance for more recent forecast windows.

6 Simulations

The aim of the simulations that follow is to demonstrate probabilistic forecast reconciliation including the algorithm discussed in Section 5. For all simulations, the tuning parameters for the SGD are set as $\eta = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 1 \times 10^{-8}$, which are the values recommended by Kingma & Ba (2014) and used in popular software packages such as TensorFlow, Keras and Torch amongst others. Convergence is achieved when the change in all gradients is less than 10% of the step size η . The number of sample periods used to construct the objective function is $R = 250$, while the number of draws used to estimate each score is $Q = 250$. All estimation of base models uses a sample size of $T = 500$. All forecast evaluations are carried out using a rolling window, also of size $W = 500$.

6.1 Data Generating Processes

The data generating process we consider corresponds to the 3-level hierarchical structure presented in Figure 1. Bottom-level series are first generated from $\text{ARIMA}(p, d, q)$ processes, which are in-turn aggregated to form the middle and top-level series. The orders p and q are randomly selected from $\{1, 2\}$ for each bottom-level series. The AR and MA parameters are randomly and uniformly generated from $[0.3, 0.5]$ and $[0.3, 0.7]$ respectively, and only accepted if they belong to the stationary and invertible region. In addition a non-stationary case where d is randomly chosen for each bottom-level from $\{0, 1\}$ was considered, these results are omitted for brevity. A complete set of results are available at the github repository <https://git.io/JJwQB>.

We consider a multivariate Gaussian and a non-Gaussian setting for the errors driving the ARIMA processes. Specifically, the non-Gaussian errors are drawn from a meta-distribution of a Gumbel copula with Beta(1, 3) margins. After simulating from the ARIMA models, additional noise is added to ensure bottom-level series have a lower signal-to-noise ratio than top level series with details provided in Appendix D. For each series the first

500 observations are ignored to avoid the impact of initial values.

6.2 Modelling and Base Forecasts

We fit univariate ARIMA models to each series using the `ARIMA()` function in the `fable` package (O’Hara-Wild et al. 2020) in R (R Core Team 2018). Note that the order of the ARIMA models is not set to the true order but is chosen using the algorithm of Hyndman & Khandakar (2008), allowing for the possibility of misspecification. Indeed, an advantage of forecast reconciliation is the ability to down-weight the forecasts of series within the hierarchy that come from misspecified models. We also considered exponential smoothing (ETS) models using the `ETS()` function in the `fable` package. These are omitted for brevity; please refer to <https://git.io/JJwQB> for a full set of results.

Let $\hat{\mathbf{y}}_{t+h|t} = (\hat{y}_{1,t+h|t}, \dots, \hat{y}_{n,t+h|t})'$, where $\hat{y}_{i,t+h|t}$ is the h -step-ahead point forecast for series i , and $\mathbf{E} := \{e_{i,t}\}_{i=1,\dots,n;t=1,\dots,T}$ where $e_{i,t}$ is the residual stacked in the $(n \times T)$ matrix. For each series and model, base probabilistic forecasts for $h = 1$ are constructed in the following four ways:

- **Independent Gaussian:** The base probabilistic forecast is made up of independent Gaussian distributions with the forecast mean and variance of variable i given by $\hat{y}_{i,t+h|t}$ and $\hat{\sigma}_{i,t+h|t}^2$, where $\hat{\sigma}_{i,t+h|t}^2$ is the sample variance of the residuals in the i th row of \mathbf{E} .
- **Joint Gaussian:** The base probabilistic forecast is a multivariate Gaussian distribution with the forecast mean $\hat{\mathbf{y}}_{t+h|t}$ and variance covariance matrix $\hat{\Sigma}$, where $\hat{\Sigma}$ is the variance covariance matrix of the residuals.
- **Independent Bootstrap:** A draw from the base probabilistic forecast is made independently for each variable as $\hat{y}_{i,t+h|t} + e_{i,\tau}$ where τ is drawn randomly (with replacement) from $1, 2, \dots, T$.
- **Joint Bootstrap:** A draw from the joint probabilistic forecast is made as $\hat{\mathbf{y}}_{t+h|t} + \mathbf{e}_\tau$ where \mathbf{e}_τ is the τ th column of \mathbf{E} , where τ is drawn randomly (with replacement) from $1, 2, \dots, T$.

We restrict our attention to the case of $h = 1$ although these methods can be generalised to larger h using the recursive method (Hyndman & Athanasopoulos 2018). For multi-step-ahead forecasts, methods (c) and (d) should be sampled in blocks to preserve serial

dependence in the residuals.

6.3 Reconciliation

For each DGP, model and method for obtaining base forecasts, reconciled probabilistic forecasts are obtained using each of the following techniques.

- **Base:** The base forecasts with no reconciliation.
- **JPP:** The best method of Jeon et al. (2019). This is equivalent to reconciling quantiles. A sample is drawn from the base forecast, these are ranked, one variable at a time (so that the smallest value drawn from each variable are put together, etc.). These are then pre-multiplied by $\mathbf{S}(\mathbf{S}'\mathbf{W}\mathbf{S})^{-1}\mathbf{S}'\mathbf{W}$ where \mathbf{W} is a diagonal matrix with elements $(1/4^2, 1/2^2, 1/2^2, 1, 1, 1, 1)$. These are the squared reciprocals of the number of bottom-level series used to form an aggregate.
- **BTTH:** The method of Ben Taieb et al. (2020). This is a method whereby draws from the probabilistic forecasts of the bottom-level series are permuted so that they have the same empirical copula as the residuals. These are then aggregated to form a sample from the distribution of all series. The mean is adjusted to be equivalent to the mean that would be obtained using the MinT method of Wickramasuriya et al. (2019) described in Table 1.
- **BottomUp:** Reconciliation via premultiplication by $\mathbf{S}\mathbf{G}$ where $\mathbf{G} = (\mathbf{0}_{m \times (n-m)}, \mathbf{I}_{m \times m})$.
- **OLS:** Reconciliation via pre-multiplication by $\mathbf{S}(\mathbf{S}'\mathbf{S})^{-1}\mathbf{S}'$.
- **MinTShr:** Reconciliation via pre-multiplication using the shrinkage estimator of the covariance matrix used by Wickramasuriya et al. (2019) but applied to probabilistic rather than point forecasting.
- **ScoreOptE:** The algorithm described in Section 5 used to optimise energy score.
- **ScoreOptV:** The algorithm described in Section 5 used to optimise variogram score.

Note that JPP and BTTH are two methods previously existing in the literature. The methods BottomUp, OLS, and MinTShr have been used extensively in the point forecasting literature but their application to probabilistic forecasting for general base forecasts is, to the best of our knowledge, novel in this paper.

In addition to these methods, two further reconciliation methods were considered; WLS, which reconciles via pre-multiplication by the same matrix used in Jeon et al. (2019) but

without any reordering of the draws, and MinTSam which uses a sample estimate of the covariance matrix rather than a shrinkage estimator. These methods were mostly dominated by OLS and MinTShr respectively and are therefore omitted for brevity; please refer to <https://git.io/JJwQB> for a full set of results.

6.4 Results for Gaussian probabilistic forecasts

Table 3 shows the mean energy score for different reconciliation methods and different methods of generating base forecasts. When base probabilistic forecasts are generated independently, score optimisation with the energy score (ScoreOptE) performs best, while when base forecasts are generated jointly, the MinT method for reconciliation using the shrinkage estimator (MinTShr) yields the most accurate forecasts. The bottom-up method as well as BTTH and JPP fail to even improve upon base forecasts in all cases. As expected score optimisation using the variogram score does not perform as well as score optimisation using energy score, when evaluation is carried out with respect to the latter. However, the results are quite close suggesting that score optimisation is fairly robust to using an alternative proper score.

Suggestions for below. (1) We agreed with Tas to swap the columns so that they are Ind. Gaussian — Ind. Bootstrap — Joint Gaussian — Joint Bootstrap. (2) Was thinking about leaving Base as is and presenting everything as percentage difference to Base Ind. Gaussian (could be messy and hard to read) (3) Maybe present these as the sample plot below 3. Rob any other ideas?

I think we should just show the graph sketched in Fig 3, but showing energy scores not ratios.

To assess significant differences between the reported results, we use post-hoc Nemenyi tests (Hollander et al. 2013). The Nemenyi test is a non-parametric test that identifies groups of forecasts which cannot be significantly distinguished from one another. We use the implementation of the tests available in the `tsutils` R package (Kourentzes 2019). Figure 4 reports the results which should be looked at column-wise. A blue square indicates that the method in the corresponding row, is statistically indistinguishable from the method in that column. For all four methods of generating base forecasts, MinTShr, ScoreOptE and OLS significantly outperform base forecasts, bottom-up forecasts, BTTH and JPP.

Table 3: Mean Energy score for Gaussian DGP. Values in bold indicate the best reconciliation method for a given base forecasting method.

Method	Ind. Bootstrap	Ind. Gaussian	Joint Bootstrap	Joint Gaussian
Base	11.3256	11.3247	11.0662	11.0556
BottomUp	11.9759	11.9705	11.6429	11.6203
BTTH	21.7508	21.8051	21.8542	21.8826
JPP	22.8651	22.8986	22.8573	22.8599
MinTShr	10.8865	10.8905	10.7475	10.7428
OLS	11.1284	11.1264	10.8396	10.8292
ScoreOptE	10.8053	10.8297	10.8433	10.8406
ScoreOptV	11.2066	11.2074	11.0413	11.0262

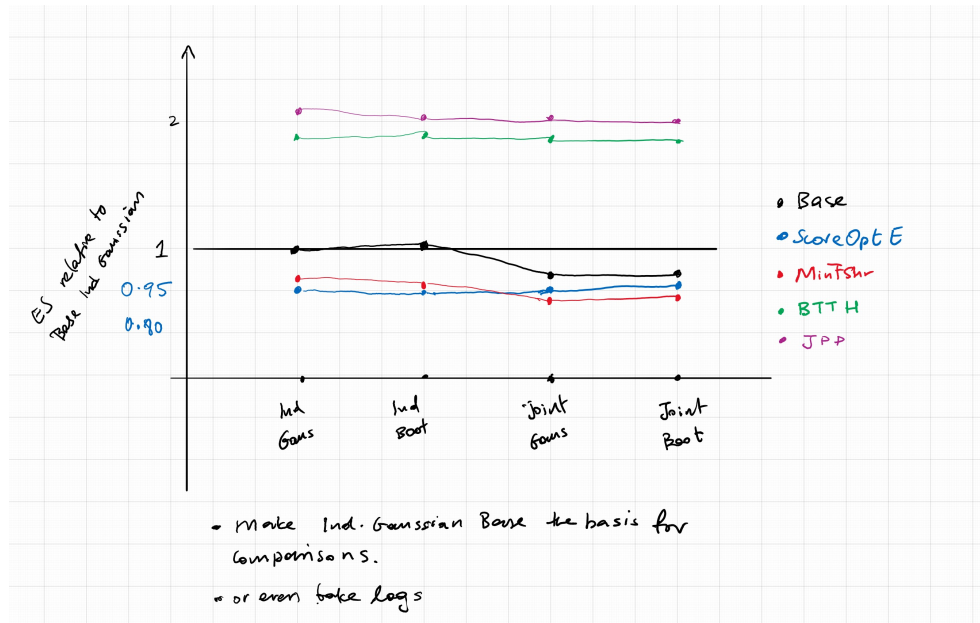


Figure 3

I had the idea of joining these like shown on the bottom of page 24 (in the annotated by George pdf) with the average rank shown on the rows of the left hand side - as we really don't care that much about these (we do not even comment on them). From discussions with Tas he suggested an alternative may be a 32×32 where everything is compared to everything. It may be that the best way to present these is the current Figure 4. Rob what

do you think?

I think what is here is probably ok without complicating things. But show Nemenyi matrices side-by-side to save space. No need to repeat row labels in that case – as per George’s diagram on p24.

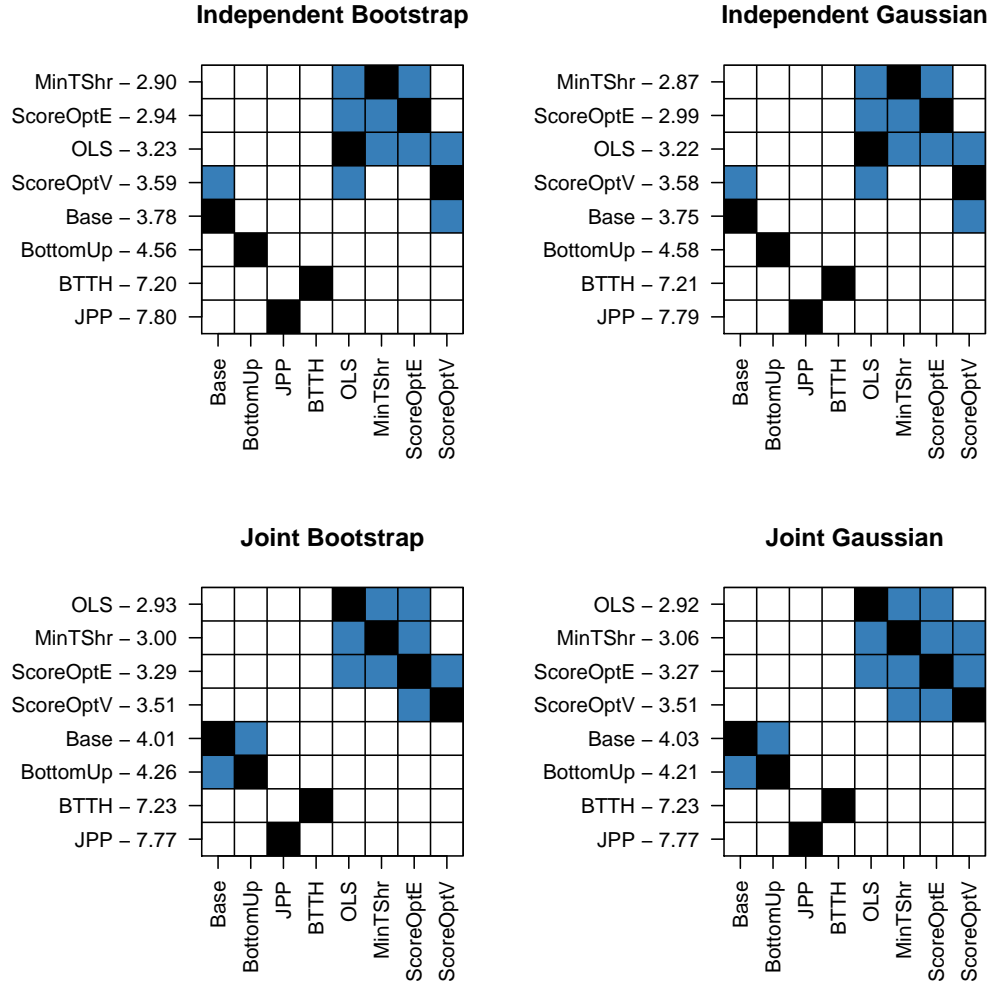


Figure 4: Nemenyi matrix for Energy Score for Gaussian DGP.

Table 4 and Figure 5 report the same output but using the variogram score for evaluation. For this specific DGP, base model and score, BTTH significantly outperforms all other methods. However, this result was not observed when using BTTH for any other simulation scenario, including those reported only in the online supplement. Excluding this result, score optimisation with respect to the variogram score is the best performing method with MinTShr and OLS also performing well. Score optimisation, OLS, MinTShr

and BTTH all lead to significantly improvements relative to base, bottom-up and JPP.

Suggest Table 4 be replaced with a graph similar to Fig 3.

Table 4: *Mean Variogram score for Gaussian DGP. Values in bold indicate the best reconciliation method for a given base forecasting method.*

Method	Ind. Bootstrap	Ind Gaussian	Joint Bootstrap	Joint Gaussian
Base	694.5012	694.6362	694.7861	694.5443
BottomUp	888.7264	891.1163	697.9468	697.7384
BTTH	630.8518	632.3063	632.7646	631.8727
JPP	1032.9742	1039.0791	1033.3112	1038.1928
MinTShr	683.3065	683.0490	677.8987	677.4923
OLS	713.4531	713.5246	680.3692	680.1407
ScoreOptE	686.0447	686.8067	683.5021	684.2645
ScoreOptV	675.1864	675.2327	673.5918	672.3905

6.5 Results for non-Gaussian probabilistic forecast

Table 5 reports the mean energy score for the non-Gaussian DGP. Overall, the results are quite similar to the Gaussian DGP. The best performing reconciliation method is ScoreOptE when base probabilistic forecasts are independent, and MinTShr when base forecasts are dependent. The Nemenyi matrix is omitted for brevity; please refer to <https://git.io/JJwQB> for a full set of results. However, these lead to similar conclusion to Figure 4. The methods ScoreOptE, MinTShr and OLS are statistically indistinguishable from one another but are significantly better than base forecasts and the bottom-up method. The methods BTTH and JPP lead to a statistically significant deterioration in forecast quality relative to base forecasts.

Suggest Tables 5 and 6 be replaced with a graph similar to Fig 3.

Finally, Table 6 and Figure 6 report results for the non-Gaussian DGP using the variogram score to evaluate forecasts. In this case, score optimisation with respect to the variogram score yields the best performance when base forecasts are dependent, while MinTShr yields the best performance when base forecasts are independent. In contrast to the Gaussian DGP, the JPP method leads to significant improvements over base forecasts,

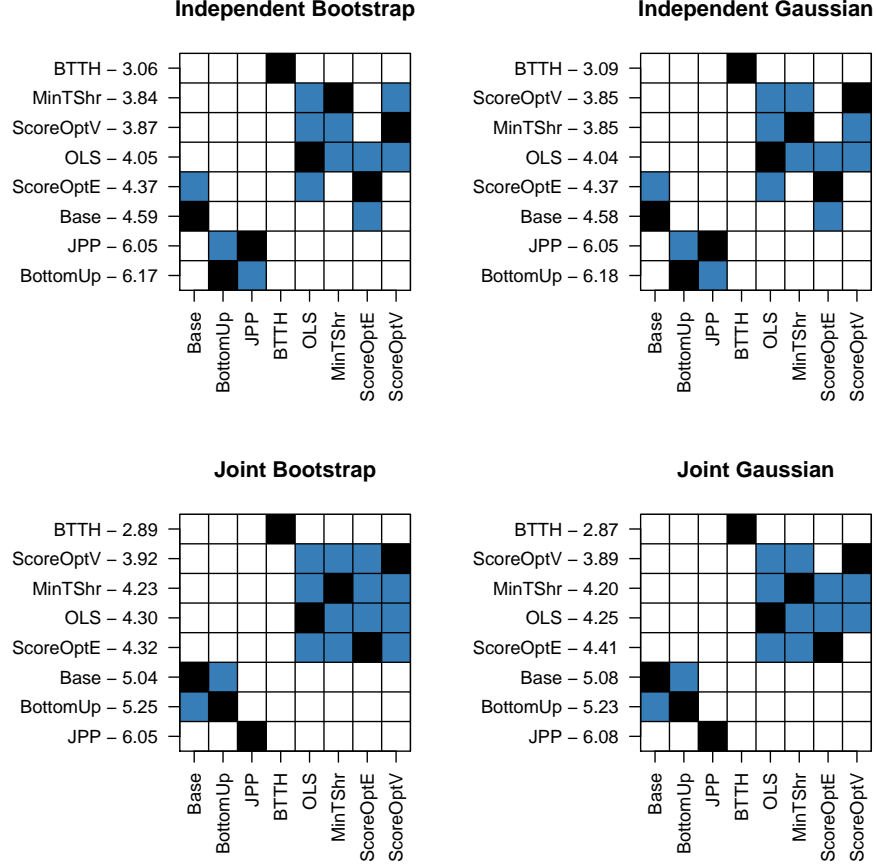


Figure 5: Nemenyi matrix for Variogram score for Gaussian Stationary DGP.

Table 5: Mean Energy score for non-Gaussian DGP. Values in bold indicate the best reconciliation method for a given base forecasting method.

Method	Ind. Bootstrap	Ind. Gaussian	Joint Bootstrap	Joint Gaussian
Base	1.4169	1.4246	1.3850	1.3912
BottomUp	1.5071	1.5305	1.4658	1.4731
BTTH	2.7704	2.9271	2.7817	2.9335
JPP	2.8805	2.9737	2.8788	2.9672
MinTShr	1.3514	1.3504	1.3297	1.3367
OLS	1.3643	1.3638	1.3409	1.3483
ScoreOptE	1.3391	1.3379	1.3384	1.3391
ScoreOptV	1.3772	1.3772	1.3746	1.3738

while the BTTH method leads to a significantly worse performance than base forecasts.

Table 6: *Mean Variogram score for non-Gaussian DGP. Values in bold indicate the best reconciliation method for a given base forecasting method.*

Method	Ind. Bootstrap	Ind. Gaussian	Joint Bootstrap	Joint Gaussian
Base	26.2480	26.2638	26.2339	26.2643
BottomUp	28.5977	28.7956	27.6925	27.7625
BTTH	29.2544	29.1507	29.0767	29.0854
JPP	26.3968	26.3993	26.3833	26.4030
MinTShr	25.4412	25.4215	25.4799	25.5248
OLS	25.6498	25.6362	25.6860	25.7348
ScoreOptE	25.6992	25.6954	25.5763	25.6889
ScoreOptV	25.5433	25.5325	25.4541	25.4798

Overall, the main conclusion from the simulation study is that score optimisation leads to significant improvements in forecast performance over base forecasts irrespective of whether the DGP is Gaussian or non-Gaussian and irrespective of whether the energy or variogram score is used for evaluation. For the DGPs considered in the simulation study, MinTShr and to a lesser extent OLS also provided significant improvements over base and bottom-up forecasts. While the existing probabilistic forecast reconciliation methods considered in the literature (BTTH and JPP) performed well in some scenarios (particularly BTTH for the Gaussian DGP evaluated by variogram score), overall results for these methods was mixed and even led to a statistically significant deterioration in forecast quality relative to base forecasts in some settings.

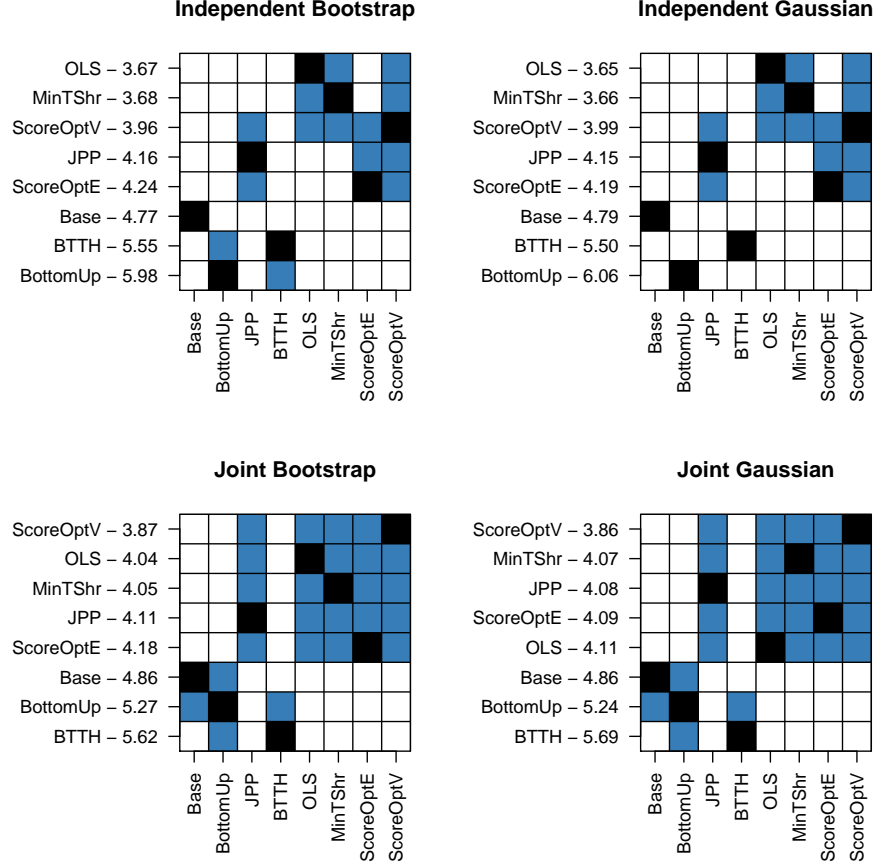


Figure 6: Nemenyi matrix for Variogram score with a non-Gaussian DGP.

7 Forecasting Australian Electricity Generation

7.1 Data Description

To demonstrate the potential of the proposed methods, we consider an application to forecasting Australian electricity generation from different sources of energy. Daily time series were obtained from opennem.org.au, a website that compiles publicly available data from the Australian Energy Market Operator (AEMO). Probabilistic day-ahead forecasts are crucial inputs into operational and planning decisions that ensure the efficiency and stability of the power network. This has become a more challenging problem with the growth in intermittent sources of generation such as wind and solar.

The hierarchy comprises three levels of aggregation.

1. *Total* generation is the sum of generation from *Renewable* and *non-Renewable* sources.

2. *Renewable* generation is the sum of *Batteries*, *Hydro (inc. Pumps)*, *Solar*, *Wind* and *Biomass*. *Non-Renewable* is the sum of *Coal*, *Gas* and *Distillate*
3. *Battery* generation is given by *Battery (Discharging)* minus *Battery (Charging)*, *Hydro (inc. Pumps)* is *Hydro* generation minus *Pumps* (energy used to pump water upstream), while *Solar* generation is the sum of *Solar (Rooftop)* and *Solar (Utility)*. *Coal* generation is the sum of *Black Coal* and *Brown Coal*, while *Gas* is the sum of *Gas (OCGT)*, *Gas (CCGT)*, *Gas (Steam)*, *Gas (Reciprocating)*.

In total, there are $n = 23$ series of which $m = 15$ are bottom-level series.

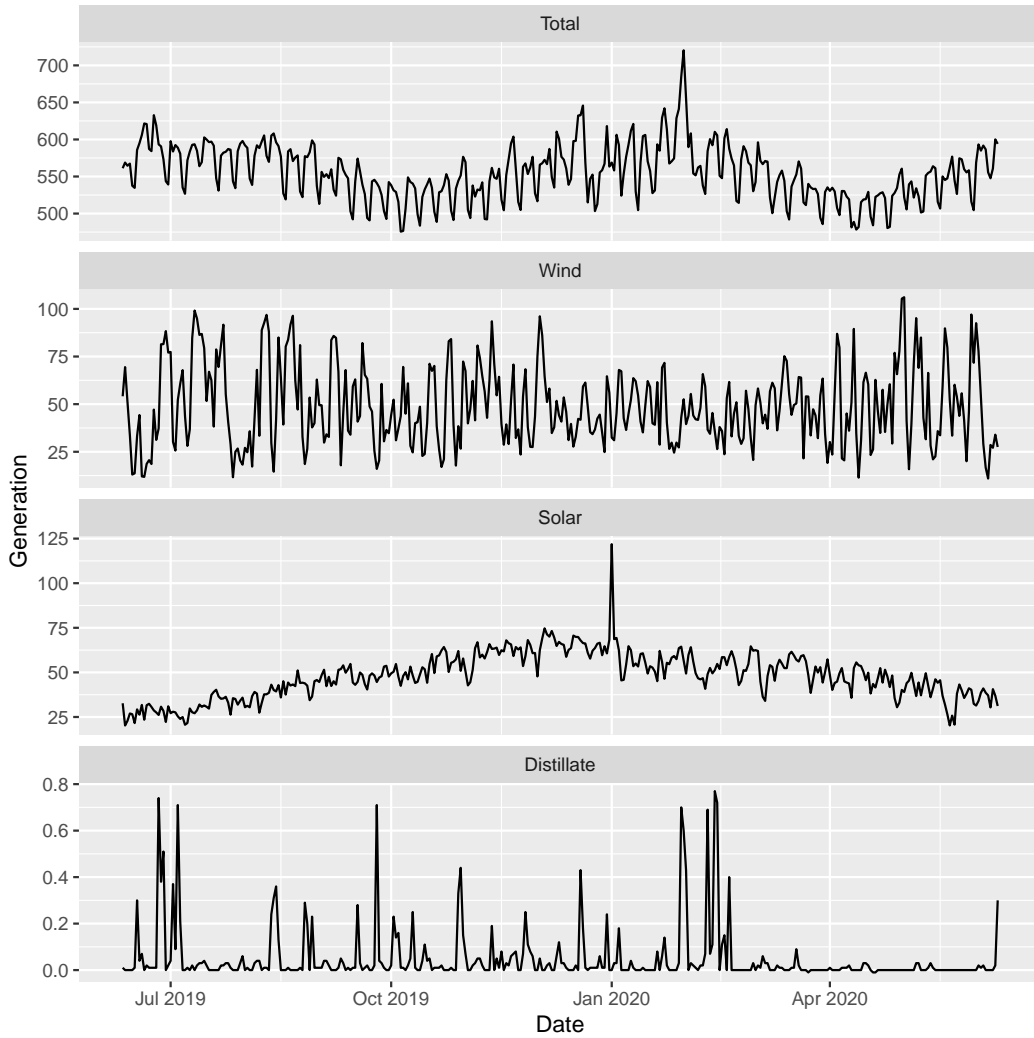


Figure 7: Time series plots for selected series from 11 June 2019 to 10 June 2020.

Figure 7 shows time plots for some selected series². The series exhibit some interesting and unique features. At the aggregate level, *Total* generation shows strong weekly seasonality, with troughs corresponding to weekends. An annual seasonal pattern is also displayed with peaks occurring during the months of June–August as well as December–February. These periods correspond to the winter and summer months in Australia for which electricity demand peaks for heating and cooling purposes respectively. As expected, generation from *Solar* peaks during the summer months of December–February. There are also some unusually large spikes observed in both the *Total* and *Solar* series during February and January 2020 respectively. *Wind* displays higher volatility (especially outside the summer months), while generation from *Distillate* exhibits aperiodic spikes.

The diversity and prominence of the features in each series and each level of aggregation highlights the importance of modelling and forecasting each series on its own merits and then applying a reconciliation approach.

7.2 Base Forecasts

The forecast evaluation is based on a rolling window. Each training window consists of 140 days (20 weeks) of data. One-step-ahead forecasts were generated leading to 170 daily forecasts for evaluation. Each series was independently modelled using a one-layer feed-forward neural network with up to 28 lags of the target variable as inputs. This was implemented using the `NNETAR` function in the `fable` package. Neural networks are used to highlight the versatility of reconciliation to different forecasting approaches. While including more layers or meteorological variables as predictors will probably lead to improved base forecasts, the primary objective is to assess the effectiveness of different forecast reconciliation methods.

Four situations were considered where base forecasts are assumed to be either Gaussian or bootstrapped from residuals, and either dependent or independent. The histograms in Figure 8 demonstrate departures from normality while the correlation heatmap in Figure 9 demonstrates departures from independence. Therefore, independent Gaussian probabilistic forecasts are likely to represent a case of severe misspecification.

²Time plots for the remaining series are available from <https://git.io/JJwd0>.

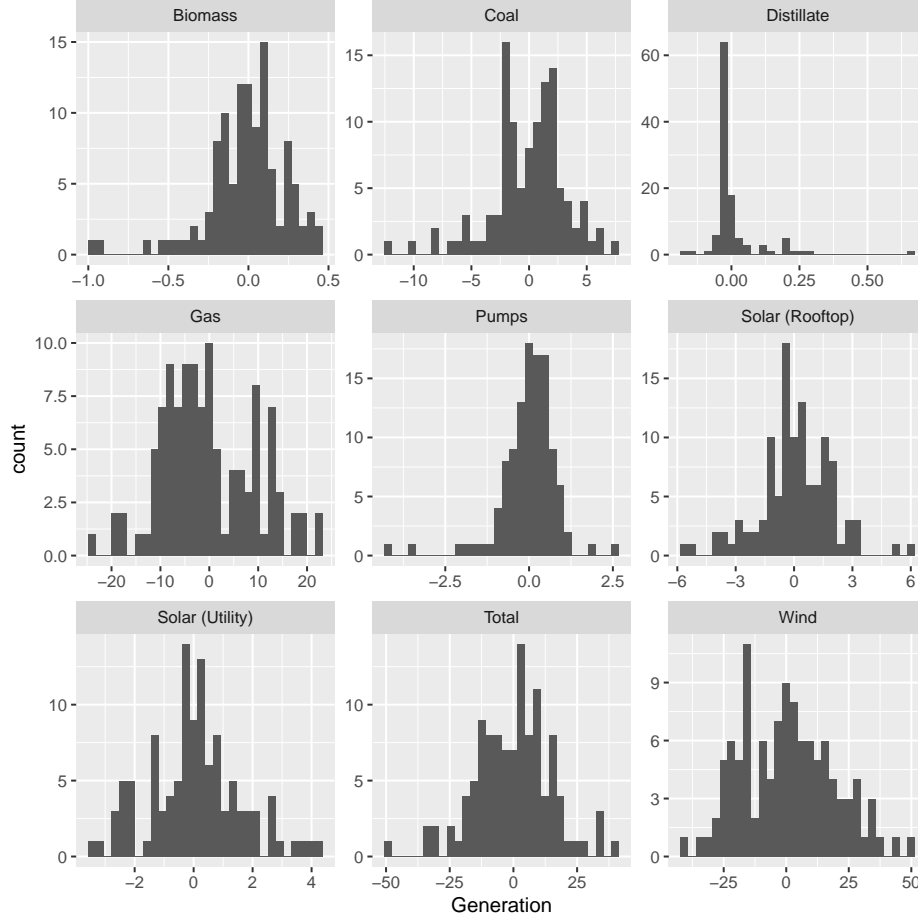


Figure 8: *Densities of residuals for selected series from a typical training window of 2 October 2019 to 21 January 2020.*

7.3 Reconciliation

The same methods were used for reconciliation as in the simulation study with score optimisation based on an objective function containing 56 days (8 weeks) of score evaluations. For brevity, only the results for energy score are presented here; please refer to <https://git.io/JJwQB> for a full set of results.

The mean energy score for all four base forecasting methods is summarised in Table 7. When base forecasts are generated assuming both independence and a Gaussian distribution, score optimisation achieves a mean energy score that is considerably smaller than all other competing methods, with MinT providing the second smallest value. Figure 10 (left) provides the Nemenyi matrix for this set of base forecasts, and shows that the superior

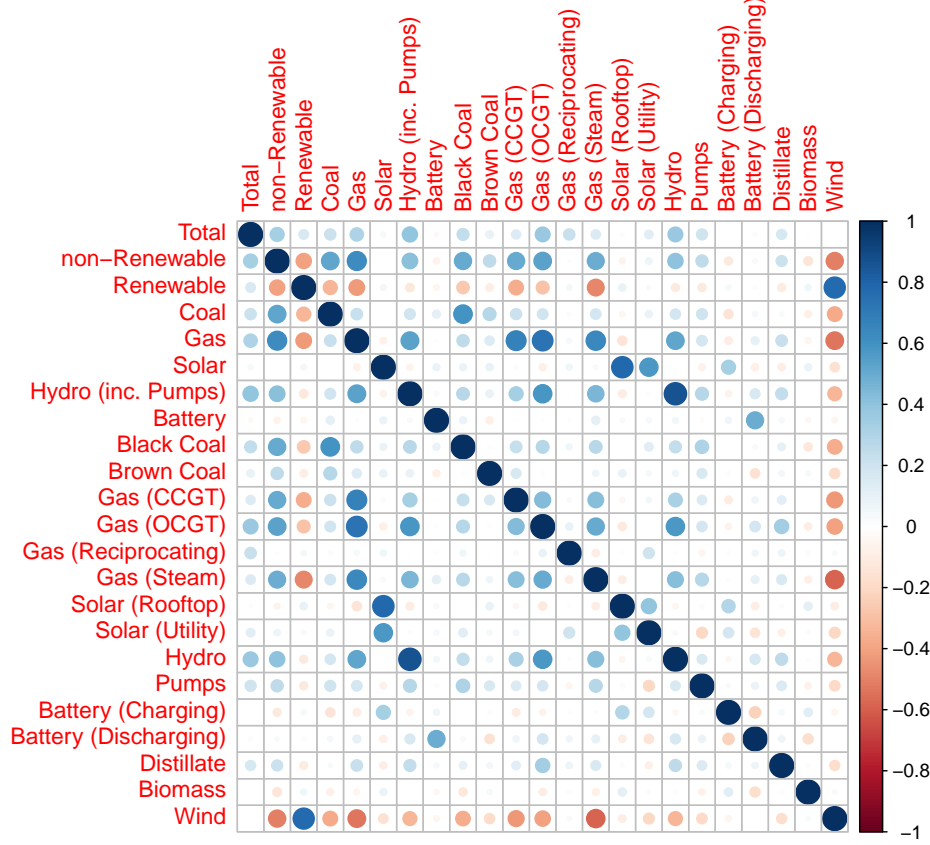


Figure 9: *Correlation heatmap of residuals from a typical training window of 2 October 2019 to 21 January 2020. Blue circles indicate positive correlation, while red circles indicate negative correlation with larger circles indicating stronger correlations.*

forecasting performance of score optimisation is statistically significant. This suggests that score optimisation is best for guarding against severe model misspecification.

For all other methods the best performing method is OLS. This difference is statistically significant as seen in Figure 10 (right), which shows the Nemenyi matrix for jointly bootstrapped base forecasts. The corresponding figures for joint Gaussian and independent bootstrap look mostly similar to Figure 10; please refer to <https://git.io/JJwQB> for a full set of results. Although score optimisation does not improve upon base, the differences are not significant. For all base forecasts, both JPP and BTTH are significantly worse than base forecasts.

Table 7: *Mean Energy Score for electricity application. Values in bold indicate the best reconciliation method for a given base forecasting method.*

RecoMethod	Ind. Bootstrap	Ind. Gaussian	Joint Bootstrap	Joint Gaussian
Base	36.41	95.02	36.54	36.32
BottomUp	35.89	90.64	35.72	35.51
BTTH	58.57	269.40	58.23	59.07
JPP	55.15	292.27	55.16	55.89
MinTShr	37.89	54.80	38.52	38.33
OLS	35.42	73.72	34.71	34.48
ScoreOptE	38.23	39.96	39.39	39.07

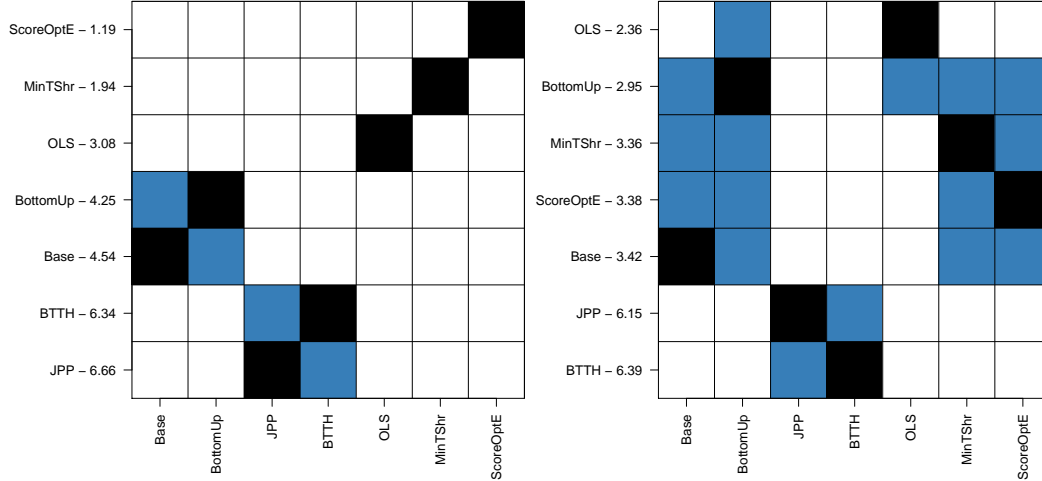


Figure 10: *Nemenyi matrices for Energy score for the electricity application. Left: base forecasts are independent and Gaussian. Right: base forecasts are obtained by jointly bootstrapping residuals.*

8 Conclusions

This paper introduces a rigorous formulation of forecast reconciliation in the probabilistic setting. It can be applied when the probabilistic forecast is either available as a density or when a sample has been drawn from the probabilistic forecast. In the elliptical case, we prove that reconciliation can recover the correct probability distribution as long as the base forecast is of the correct form, irrespective of the scale and location of the base forecast. Probably due to this reason, score optimisation works well in applications even when the base forecasts are assumed to be independent.

We also prove that the log score is not proper when comparing incoherent and coherent forecasts. Consequently, we introduce a new algorithm that trains reconciliation weights by minimising the energy score or variogram score. Since the scores are approximated by Monte Carlo simulation, stochastic gradient descent is used for optimisation. This method is shown to lead to significant improvements over base forecasts, bottom-up methods and existing probabilistic reconciliation approaches across a wide variety of simulated and empirical examples.

An interesting result is that projection methods with certain optimality properties in the point forecasting setting, also work well when extended to the probabilistic case. In particular, a simple least squares projection is the best performing method in the high-dimensional empirical example, provided the base forecasts are not too badly misspecified. This may arise since projections implicitly provide constrained versions of the reconciliation weights. A promising future research avenue may involve regularised versions of score optimisation that add an L_1 or L_2 penalty to the objective function. Alternatively, early stopping (Bühlmann & Yu 2003, Yao et al. 2007) of the gradient descent may lead to a better bias-variance tradeoff in learning reconciliation weights.

A final important avenue of future research is the development of probabilistic forecast reconciliation for domains other than the real line. These may include domains constrained above zero, discrete domains, or domains that are a mixture of continuous distributions and discrete point masses. While such problems are challenging, the geometric interpretation of probabilistic forecast introduced in this paper, lays the foundation for this research agenda.

A Proof of Theorem 3.1 and Theorem 3.2

Consider the region \mathcal{I} given by the Cartesian product of intervals $(l_1, u_1), (l_2, u_2), \dots (l_m, u_m)$. We derive the probability, under the reconciled measure, that the bottom-level series lie in \mathcal{I} , i.e. $\Pr(\boldsymbol{\ell} \succ \mathbf{b} \succ \mathbf{u})$, where $\boldsymbol{\ell} = (\ell_1, \ell_2, \dots, \ell_m)$, $\mathbf{u} = (u_1, u_2, \dots, u_m)$ and \succ denotes element-wise inequality between vectors. The pre-image of \mathcal{I} under g can similarly be denoted as all points \mathbf{y} satisfying $\boldsymbol{\ell} \succ \mathbf{G}\mathbf{y} \succ \mathbf{u}$. Using Definition 2.2,

$$\Pr(\boldsymbol{\ell} \succ \mathbf{b} \succ \mathbf{u}) = \int_{\boldsymbol{\ell} \succ \mathbf{G}\mathbf{y} \succ \mathbf{u}} \hat{f}(\mathbf{y}) d\mathbf{y},$$

where \hat{f} is the density of the base probabilistic forecast. Now consider a change of variables to an n -dimensional vector \mathbf{z} where $\mathbf{y} = \mathbf{G}^* \mathbf{z}$. Recall, $\mathbf{G}^* = (\mathbf{G}^- \vdash \mathbf{G}_\perp)$, \mathbf{G}^- is a generalised inverse of \mathbf{G} , and \mathbf{G}_\perp is an orthogonal complement of \mathbf{G} . By the change of variables

$$\begin{aligned} \Pr(\boldsymbol{\ell} \succ \mathbf{b} \succ \mathbf{u}) &= \int_{\boldsymbol{\ell} \succ \mathbf{G}\mathbf{y} \succ \mathbf{u}} \hat{f}(\mathbf{y}) d\mathbf{y} \\ &= \int_{\boldsymbol{\ell} \succ \mathbf{G}\mathbf{G}^* \mathbf{z} \succ \mathbf{u}} \hat{f}(\mathbf{G}^* \mathbf{z}) |\mathbf{G}^*| d\mathbf{z} \\ &= \int_{\boldsymbol{\ell} \succ \mathbf{z}_1 \succ \mathbf{u}} \hat{f}(\mathbf{G}^* \mathbf{z}) |\mathbf{G}^*| d\mathbf{z}, \end{aligned}$$

where \mathbf{z}_1 denotes the first m elements of \mathbf{z} . Letting \mathbf{a} denote the last $n - m$ elements of \mathbf{z} the integral above can be written as

$$\Pr(\mathbf{b} \in \mathcal{I}) = \int_{\boldsymbol{\ell} \succ \mathbf{z}_1 \succ \mathbf{u}} \int \hat{f}(\mathbf{G}^- \mathbf{z}_1 + \mathbf{G}_\perp \mathbf{a}) |\mathbf{G}^*| d\mathbf{a} d\mathbf{z}_1$$

Replacing \mathbf{z}_1 with \mathbf{b} , it can be seen that the term inside the outer integral is a density for the bottom-level series. Therefore

$$\tilde{f}_b(\mathbf{b}) = \int \hat{f}(\mathbf{G}^- \mathbf{b} + \mathbf{G}_\perp \mathbf{a}) |\mathbf{G}^*| d\mathbf{a}, \quad (4)$$

is the density of \mathbf{b} . To obtain the density of the full hierarchy we first augment the density in Equation (4) by $n - m$ variables denoted \mathbf{u}

$$f(\mathbf{b}, \mathbf{u}) = \tilde{f}_b(\mathbf{b}) \mathbb{1}\{\mathbf{u} = \mathbf{0}\}, \quad (5)$$

such that the density $f(\mathbf{b}, \mathbf{u})$ is a density for n -dimensional vector that is degenerate across the dimensions corresponding to \mathbf{u} . Using the change of variables,

$$\mathbf{y} = (\mathbf{S} \vdash \mathbf{S}_{\perp}^{-}) \begin{pmatrix} \mathbf{b} \\ \mathbf{u} \end{pmatrix},$$

where \mathbf{S}_{\perp}^{-} is a generalised inverse such that $\mathbf{S}_{\perp}' \mathbf{S}_{\perp}^{-} = \mathbf{I}$ and noting the inverse of $(\mathbf{S} \vdash \mathbf{S}_{\perp})$ is given by

$$\mathbf{S}^* := \begin{pmatrix} \mathbf{S}^{-} \\ \mathbf{S}_{\perp}' \end{pmatrix},$$

it can be seen that $\mathbf{b} = \mathbf{S}^{-} \mathbf{y}$ and $\mathbf{u} = \mathbf{S}_{\perp}' \mathbf{y}$. Applying this change of variables yields the density

$$\tilde{f}_{\mathbf{y}}(\mathbf{y}) = |\mathbf{S}^*| \tilde{f}_{\mathbf{b}}(\mathbf{S}^{-} \mathbf{y}) \mathbb{1}\{\mathbf{S}_{\perp}' \mathbf{y} = \mathbf{0}\}.$$

Since \mathbf{S}_{\perp}' is the orthogonal complement of \mathbf{S} and since the columns of \mathbf{S} span the coherent subspace, the statement $\mathbf{S}_{\perp}' \mathbf{y} = \mathbf{0}$ is equivalent to the statement $\mathbf{y} \in \mathfrak{s}$. As such, the reconciled density is given by

$$\tilde{f}_{\mathbf{y}}(\mathbf{y}) = |\mathbf{S}^*| \tilde{f}_{\mathbf{b}}(\mathbf{S}^{-} \mathbf{y}) \mathbb{1}\{\mathbf{y} \in \mathfrak{s}\}.$$

B Proof of Theorem 3.4

Let

$$\hat{\Sigma} = \Sigma + D = S\Omega S' + D.$$

If reconciliation is carried out via a projection onto \mathfrak{s} , then $SGS = S$ and

$$\begin{aligned}\tilde{\Sigma} &= SG\hat{\Sigma}G'S' \\ &= SGS\Omega S'G'S' + SGDG'S' \\ &= S\Omega S' + SGDG'S' \\ &= \Sigma + SGDG'S' .\end{aligned}$$

Therefore to recover the true predictive using a projection, some G_{opt} must be found such that $G_{\text{opt}}D = 0$. Let the eigenvalue decomposition of D be given by $R\Lambda R'$, where R is an $n \times q$ matrix with $q = \text{rank}(D)$ and Λ is an $q \times q$ diagonal matrix containing non-zero eigenvalues of D . By the rank nullity theorem, R will have an orthogonal complement R_{\perp} of dimension $n \times (n - q)$. If $q = n - m$ then the number of columns of R_{\perp} is m and G_{opt} can be formed as the $m \times n$ matrix $(R'_{\perp}S)^{-1}R'_{\perp}$. If $q < n - m$ the number of columns of R_{\perp} is greater than m , and any m columns of R_{\perp} can be used to form G_{opt} in a similar fashion. However when $q > n - m$, the number of columns of R_{\perp} is less than m and no such $m \times n$ matrix G_{opt} can be formed. Therefore the true predictive can only be recovered via a projection when $\text{rank}(D) \leq n - m$.

With respect to the location, if SG is a projection then reconciled forecasts will be unbiased as long as the base forecasts are also unbiased. When base forecasts are biased they can be bias corrected before reconciliation as described by Panagiotelis et al. (2020) in the point forecasting setting.

C Proof of Theorem 4.1

The proof relies on the following change of variables,

$$\mathbf{y} = (\mathbf{S} \vdash \mathbf{S}_\perp) \begin{pmatrix} \mathbf{b} \\ \mathbf{u} \end{pmatrix}.$$

Also recall from the proof of Theorem 3.2 that $\mathbf{S}^* = (\mathbf{S} \vdash \mathbf{S}_\perp)^{-1}$

Let the density of the true predictive $f(\mathbf{y})$ after a change of variables, be given by $|\mathbf{S}^*|^{-1} f_{\mathbf{b}}(\mathbf{b}) \mathbb{1}\{\mathbf{u} = \mathbf{0}\}$. To prove that the log score is improper we construct an incoherent base density \hat{f} such that $E_f[LS(\hat{f}, \mathbf{y})] < E_f[LS(f, \mathbf{y})]$. This incoherent density is constructed, so that after the same change of variables it can be written as $|\mathbf{S}^*|^{-1} \hat{f}_{\mathbf{b}}(\mathbf{b}) \hat{f}_{\mathbf{u}}(\mathbf{u})$. We require $\hat{f}_{\mathbf{u}}(\mathbf{0}) > 1$, i.e., \mathbf{u} is highly concentrated around $\mathbf{0}$ but still non-degenerate. An example is an independent normal with mean 0 and variances less than $(2\pi)^{-1}$. Now, let \mathbf{y}^* be a realisation from f . Let the first m elements of $\mathbf{S}^* \mathbf{y}^*$ be \mathbf{b}^* , and the remaining elements be \mathbf{u}^* . The log score for f is thus,

$$\begin{aligned} LS(f, \mathbf{y}^*) &= -\log f(\mathbf{y}^*) \\ &= -\log |\mathbf{S}^*| - \log f_{\mathbf{b}}(\mathbf{b}^*) - \log(\mathbb{1}\{\mathbf{u}^* = \mathbf{0}\}) \\ &= -\log |\mathbf{S}^*| - \log f_{\mathbf{b}}(\mathbf{b}^*), \end{aligned} \tag{6}$$

where the third term in Equation (6) is equal to zero since the fact that $\mathbf{y}^* \in \mathfrak{s}$ implies that $\mathbf{u}^* = \mathbf{0}$. The log score for \hat{f} is

$$LS(\hat{f}, \mathbf{y}^*) = -\log |\mathbf{S}^*| - \log f_{\mathbf{b}}(\mathbf{b}^*) - \log f_{\mathbf{u}}(\mathbf{0}).$$

Since $f_{\mathbf{u}}(\mathbf{0}) > 1$ by construction, $-\log f_{\mathbf{u}}(\mathbf{0}) < 0$, therefore

$$LS(\hat{f}, \mathbf{y}^*) < -\log |\mathbf{S}^*| - \log f_{\mathbf{b}}(\mathbf{b}^*) = LS(f, \mathbf{y}^*)$$

Since this holds for any possible realisation, it will also hold after taking expectations (by the monotonicity of expectations). Thus \hat{f} violates the condition for a proper scoring rule.

D Data generating process

To ensure that bottom-level series are noisier than top level series (a feature often observed empirically), noise is added to the bottom-level series in the following manner

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

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

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

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

where $w_{AA,t}, w_{AB,t}, w_{BA,t}, w_{BB,t}$ are generated from ARIMA processes as described in Section 6.1 with innovations $\varepsilon_{AA,t}, \varepsilon_{AB,t}, \varepsilon_{BA,t}, \varepsilon_{BB,t}$.

For the Gaussian DGP, $u_t \sim \mathcal{N}(0, \sigma_u^2)$ and $v_t \sim \mathcal{N}(0, \sigma_v^2)$ and $\{\varepsilon_{AA,t}, \varepsilon_{AB,t}, \varepsilon_{BA,t}, \varepsilon_{BB,t}\} \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \Sigma) \forall t$. We follow Wickramasuriya et al. (2019) and set

$$\Sigma = \begin{pmatrix} 5.0 & 3.1 & 0.6 & 0.4 \\ 3.1 & 4.0 & 0.9 & 1.4 \\ 0.6 & 0.9 & 2.0 & 1.8 \\ 0.4 & 1.4 & 1.8 & 3.0 \end{pmatrix}$$

and $\sigma_u^2 = 28$ and $\sigma_v^2 = 22$. This ensures that the following inequalities are satisfied,

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

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

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

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

For the non-Gaussian case, errors are generated from a Gumbel copula with Beta margins as described in Section 6.1. Rather than add Gaussian noise, we simulate u_t and v_t from skew t distributions using the `sn` package (Azzalini 2020). The scale, skew and degrees of freedom parameters are chosen as 0.5, 1.5 and 4 and 0.9, 2 and 8 for u_t and v_t respectively. Monte Carlo simulations show that these values satisfy the inequalities described above.

References

- Abramson, B. & Clemen, R. (1995), ‘Probability forecasting’, *International Journal of Forecasting* **11**(1), 1–4.
- 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**(1), 60–74.
- Azzalini, A. (2020), *sn: The Skew-Normal and Related Distributions such as the Skew-t*, Università di Padova, Italia. R package version 1.6-1.
URL: <https://CRAN.R-project.org/package=sn>
- Ben Taieb, S., Huser, R., Hyndman, R. J. & Genton, M. G. (2017), ‘Forecasting uncertainty in electricity smart meter data by boosting additive quantile regression’, *IEEE Transactions on Smart Grid* **7**(5), 2448–2455.
- Ben Taieb, S. & Koo, B. (2019), Regularized regression for hierarchical forecasting without unbiasedness conditions, in ‘KDD ’19: Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining’, pp. 1337–1347.
- Ben Taieb, S., Taylor, J. W. & Hyndman, R. J. (2020), ‘Hierarchical probabilistic forecasting of electricity demand with smart meter data’, *Journal of the American Statistical Association* . in press.
- Böse, J.-H., Flunkert, V., Gasthaus, J., Januschowski, T., Lange, D., Salinas, D., Schelter, S., Seeger, M. & Wang, Y. (2017), ‘Probabilistic demand forecasting at scale’, *Proceedings of the VLDB Endowment* **10**(12), 1694–1705.
- Bottou, L. (2010), Large-scale machine learning with stochastic gradient descent, in Y. Lechevallier & G. Saporta, eds, ‘Proceedings of COMPSTAT’2010’, Physica-Verlag HD, pp. 177–186.
- Bühlmann, P. & Yu, B. (2003), ‘Boosting with the L_2 loss: regression and classification’, *Journal of the American Statistical Association* **98**(462), 324–339.

- Carpenter, B., Hoffman, M. D., Brubaker, M., Lee, D., Li, P. & Betancourt, M. (2015), ‘The Stan math library: Reverse-mode automatic differentiation in C++’.
URL: <http://arxiv.org/abs/1509.07164>
- Dawid, A. P. (2007), ‘The geometry of proper scoring rules’, *Annals of the Institute of Statistical Mathematics* **59**(1), 77–93.
- Dunn, D. M., Williams, W. H. & Dechaine, T. L. (1976), ‘Aggregate versus subaggregate models in local area forecasting’, *Journal of the American Statistical Association* **71**(353), 68–71.
- Gneiting, T. & Katzfuss, M. (2014), ‘Probabilistic forecasting’, *Annual Review of Statistics and Its Application* **1**, 125–151.
- Gneiting, T. & Raftery, A. E. (2007), ‘Strictly proper scoring rules, prediction, and estimation’, *Journal of the American Statistical Association* **102**(477), 359–378.
- Gneiting, T., Raftery, A. E., Westveld III, A. H. & Goldman, T. (2005), ‘Calibrated probabilistic forecasting using ensemble model output statistics and minimum CRPS estimation’, *Monthly Weather Review* **133**(5), 1098–1118.
- Gross, C. W. & Sohl, J. E. (1990), ‘Disaggregation methods to expedite product line forecasting’, *Journal of Forecasting* **9**(3), 233–254.
- Hollander, M., Wolfe, D. A. & Chicken, E. (2013), *Nonparametric statistical methods*, John Wiley & Sons.
- 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*, 2nd edn, OTexts, Melbourne, Australia.
- 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., 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.
- Jeon, J., Panagiotelis, A. & Petropoulos, F. (2019), ‘Probabilistic forecast reconciliation with applications to wind power and electric load’, *European Journal of Operational Research* **279**(2), 364–379.
- Jordan, A., Krüger, F. & Lerch, S. (2017), ‘Evaluating probabilistic forecasts with the R package scoringRules’.
URL: <http://arxiv.org/abs/1709.04743>
- Kingma, D. P. & Ba, J. (2014), ‘Adam: A method for stochastic optimization’.
URL: <http://arxiv.org/abs/1412.6980>
- Kingma, D. P. & Welling, M. (2013), ‘Auto-encoding variational Bayes’.
URL: <http://arxiv.org/abs/1312.6114>
- Kourentzes, N. (2019), *tsutils: Time Series Exploration, Modelling and Forecasting*. R package version 0.9.0.
URL: <https://CRAN.R-project.org/package=tsutils>
- McLean Sloughter, J., Gneiting, T. & Raftery, A. E. (2013), ‘Probabilistic wind vector forecasting using ensembles and Bayesian model averaging’, *Monthly Weather Review* **141**(6), 2107–2119.
- Nystrup, P., Lindström, E., Pinson, P. & Madsen, H. (2020), ‘Temporal hierarchies with autocorrelation for load forecasting’, *European Journal of Operational Research* **280**(3), 876 – 888.
- O’Hara-Wild, M., Hyndman, R. & Wang, E. (2020), *fable: Forecasting Models for Tidy Time Series*. R package version 0.2.0.
URL: <https://CRAN.R-project.org/package=fable>
- Panagiotelis, A. (2020), *ProbReco: Score Optimal Probabilistic Forecast Reconciliation*. R package version 0.1.0.
URL: <https://github.com/anastasiospanagiotelis/ProbReco>

- Panagiotelis, A., Athanasopoulos, G., Gamakumara, P. & Hyndman, R. J. (2020), ‘Forecast reconciliation: A geometric view with new insights on bias correction’, *International Journal of Forecasting* . forthcoming.
- Pinson, P., Madsen, H., Papaefthymiou, G. & 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.
URL: <https://www.R-project.org/>
- Rossi, B. (2014), Density forecasts in economics, forecasting and policymaking, Technical report, Els Opuscles del CREI.
URL: <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.722.6744>
- 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).
- Scheuerer, M. & Hamill, T. M. (2015), ‘Variogram-based proper scoring rules for probabilistic forecasts of multivariate quantities’, *Monthly Weather Review* **143**(4), 1321–1334.
- 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.
- Székely, G. J. (2003), E-statistics: The energy of statistical samples, Technical Report 03-05, Department of Mathematics and Statistics, Bowling Green State University.
- Székely, G. J. & Rizzo, M. L. (2013), ‘Energy statistics: A class of statistics based on distances’, *Journal of Statistical Planning and Inference* **143**(8), 1249–1272.
- Van Erven, T. & Cugliari, J. (2015), Game-theoretically optimal reconciliation of contemporaneous hierarchical time series forecasts, in ‘Modeling and Stochastic Learning for Forecasting in High Dimensions’, Springer, pp. 297–317.

- Wickramasuriya, S. L., Athanasopoulos, G. & Hyndman, R. J. (2019), ‘Optimal forecast reconciliation for hierarchical and grouped time series through trace minimization’, *Journal of the American Statistical Association* **114**(526), 804–819.
- Wytock, M. & Kolter, J. Z. (2013), Large-scale probabilistic forecasting in energy systems using sparse Gaussian conditional random fields, *in* ‘2013 IEEE 52nd Annual Conference on Decision and Control (CDC)’, pp. 1019–1024.
- Yao, Y., Rosasco, L. & Caponnetto, A. (2007), ‘On early stopping in gradient descent learning’, *Constructive Approximation* **26**(2), 289–315.
- Zarnowitz, V. & Lambros, L. A. (1987), ‘Consensus and uncertainty in economic prediction’, *Journal of Political Economy* **95**(3), 591–621.