

NHR-TFNet: Forecasting Hierarchical Time Series using Non-Linear Mappings

Shanika L Wickramasuriya^{a*}, Kasun Bandara^b, Hansika Hewamalage^c, Maneesha Perera^d

^a Department of Econometrics and Business Statistics, Monash University, Melbourne, Australia

^b School of Computing and Information Systems, Melbourne Centre for Data Science, The University of Melbourne, Melbourne, Australia

^c School of Computer Science & Engineering, University of New South Wales, Sydney, Australia

^d Department of Mechanical Engineering, School of Electrical, Mechanical and Infrastructure Engineering, The University of Melbourne, Melbourne, Australia

Conceptualization: Shanika L Wickramasuriya, Kasun Bandara and Hansika Hewamalage

Methodology: Shanika L Wickramasuriya, Kasun Bandara and Hansika Hewamalage

Software: Shanika L Wickramasuriya, Kasun Bandara, Hansika Hewamalage and Maneesha Perera

Formal analysis: Shanika L Wickramasuriya, Kasun Bandara, Hansika Hewamalage and Maneesha Perera

Drafting the paper: Shanika L Wickramasuriya, Kasun Bandara, Hansika Hewamalage and Maneesha Perera

*Postal address: Department of Econometrics and Business Statistics, 900 Dandenong Road, Caulfield East, VIC 3145, Australia. E-mail address: shanika.wickramasuriya@monash.edu

NHR-TFNet: Forecasting Hierarchical Time Series using Non-linear Mappings

Abstract

Forecast reconciliation is a procedure of forecasting multivariate time series with linear aggregation constraints to ensure that forecasts satisfy the same set of constraints. Most state-of-the-art reconciliation methods reconcile the forecasts using linear mappings to project the base forecasts onto a coherent subspace. They also use assumptions such as the forecasts are unbiased and residuals are jointly covariance stationary. In this study, we propose a non-linear forecast reconciliation approach, NHR-TFNet, using machine learning methods, relaxing the assumptions of traditional linear reconciliation approaches. We introduce a novel loss function incorporating non-linear mappings to obtain a set of coherent forecasts from the individual base forecasts. To obtain the weights of the non-linear mapping between the base forecasts, we train a feed-forward neural network with the proposed loss function using transformed fitted values of the base models as inputs to the network. We evaluate the proposed methodology against state-of-the-art linear reconciliation approaches on multiple benchmark datasets. The results indicate that the proposed method works well when there is a reasonable number of observations in comparison to the number of bottom-level and non-bottom level series.

Keywords: Hierarchical time series, Coherent forecasts, Reconciliation, Non-linear mappings, Feed forward neural networks

1. Introduction

Nowadays, many operational forecasting tasks are associated with huge collections of time series spanning across several operational dimensions such as product categories, geographical areas, etc., which form a hierarchical structure. Often, the forecasts are collated at multiple levels for better decision-making and strategic planning from a business perspective. For instance, the product assortment structure in a retail business often follows a hierarchical structure, where it is required to produce sales forecasts at the product, store, state and country levels [4]. Also, in the electricity domain, generating energy demand forecasts at the smart-meter, grid, and regional levels is important for better demand planning and efficient resource management [27]. The prevalence of hierarchical forecasting problems is further demonstrated by recently held M5 forecasting competition [20] and Wikipedia Web Traffic forecasting competition [26], which were based on hierarchically structured time series data.

Hierarchical time series are often organised with several layers of granularity, where each level of the hierarchy can exhibit diverse patterns such as linear or nonlinear trends,

1
2
3 different seasonal shapes, level shifts, etc. For example, time series at the bottom level of
4 the hierarchy are generally noisy and sparse, whereas time series at the higher aggregation
5 levels of the hierarchy can show strong time series patterns (high signal-to-noise ratio) that
6 are relatively easier to forecast. In a hierarchical forecasting setting, it is also required to
7 generate forecasts that are coherent across the entire aggregation structure, which means the
8 sum of the forecasts must be consistent with the aggregation structure of the hierarchy [30].
9 Therefore, generating accurate, but also coherent forecasts across the whole time series
10 hierarchy is a challenging task.

11
12 The traditional approaches to forecasting hierarchical time series are mostly involved in
13 generating forecasts for one level of aggregation and then either aggregated for higher lev-
14 els (bottom-up approach), or disaggregated for lower levels (top-down approach), to obtain
15 coherent forecasts across the hierarchy [30, 12]. As these methods produce forecasts con-
16 sidering only a single aggregation level, they are unable to exploit the information available
17 across the rest of the levels in the hierarchy. Alternatively, methods have been introduced
18 to combine forecasts across the entire hierarchy by assigning weights to each individual
19 forecast in the hierarchy, where the weights are obtained through an optimisation process
20 that produces coherent forecasts. Typically, such statistically solid hierarchical forecasting
21 methods pursue a two-step process: 1) *the forecasting step* to produce base forecasts for
22 each time series in the hierarchy; and 2) *the reconciliation step* to ensure the coherency
23 of the forecasts across the hierarchy [11, 15, 30]. During the reconciliation process, these
24 methods are benefited from utilising information across the hierarchy, thus often outper-
25 forming bottom-up and top-down approaches in the literature [30]. However, these methods
26 use strong assumptions, e.g., model forecasts being unbiased, model residuals being jointly
27 covariance stationary, and focus only on linear mappings to project the base forecasts onto
28 the coherent subspace.

29
30 In parallel to these developments, with the advent of Big Data, many recent advance-
31 ments in time series forecasting have been shifting to data-driven workhorses such as machine
32 learning. In hierarchical settings, the non-parametric modelling capabilities of these models
33 provide an opportunity to apply them as non-linear mapping functions to base forecasts,
34 obviating the assumptions of traditional linear reconciliation techniques. Nevertheless, the
35 advantages of using machine learning techniques as non-linear mapping functions for fore-
36 cast reconciliation are yet to be translated into hierarchical forecasting research. Moreover,
37 the time series forecasting techniques used to produce base forecasts have also evolved from
38 traditional univariate models to global models that are trained across sets of many related
39 time series [18]. Thus, the concept of sharing cross-series information can be applied to
40 build competitive global models by clustering all the series in the hierarchy and applying a
41 global model to each cluster allowing better control of the overall model complexity [3, 10].

42
43 In the recent literature, the incorporation of structural information for hierarchical fore-
44 casting has been explored [24, 21, 8] and has shown that exploiting hierarchical information
45 in a non-linear manner leads to better accuracy than following a conventional linear recon-
46 ciliation process. For example, [24] uses XGBoost and random forest as non-linear models
47 to non-linearly map the forecasts in the hierarchy to generate better forecasts at the bottom
48 level of the hierarchy. Whereas [8] uses a regularized loss function to capture the non-linear

1
2
3 relationships among the time series of the hierarchy, which also maintains the forecast co-
4 herency of the hierarchy. Also, [21] proposes a framework that simultaneously learns from
5 all the time series in the hierarchy to produce coherent, probabilistic forecasts, removing
6 the additional post-processing reconciliation step used in the traditional reconciliation al-
7 gorithms. However, the use of non-linear model along with a regularized loss function to
8 minimise the error of the entire hierarchy and a generic hierarchical forecasting framework
9 that can be used with both traditional univariate models and global models as base models
10 has not yet been thoroughly studied. To this end, in this paper, we propose NHR-TFNet, a
11 forecasting framework that attempts to account for the non-linear relationships that exist in
12 hierarchical time series. The proposed NHR-TFNet framework follows a two-step procedure
13 to produce coherent forecasts across the structure. Firstly, the base forecasts for each time
14 series in the hierarchy are generated. NHR-TFNet framework allows both univariate and
15 global models to be employed as base forecasting models. Secondly, to obtain the weights
16 of the non-linear mapping between the base forecasts, we train a feed-forward multi-layer
17 perceptron neural network (MLP) using the transformed fitted values of the base models as
18 inputs to the network. In order to train the MLP, we introduce a novel loss function that
19 accounts for the forecast loss and the loss associated with the coherency error, which acts
20 as a penalty if the fit of the model deviates from the inherent hierarchical structure. As the
21 base forecast models, we use univariate statistical forecasting techniques such as exponential
22 smoothing methods (ETS) [14, 13] and autoregressive integrated moving average (ARIMA)
23 models [6] from the `forecast` package [13] and a global model DeepAR [22] from the `GluonTS`
24 package [1]. The NHR-TFNet framework is evaluated using multiple hierarchically struc-
25 tured time series databases, which contain various seasonal and trend patterns, exhibiting
26 different time series characteristics. The source code relevant to NHR-TFNet framework is
27 available at <https://github.com/ManeeshaPerera/hierarchical-reconciliation-ML/>

28
29 The rest of the paper is organised as follows. In Section 2, we formally define hierar-
30 chical time series and discuss the state-of-the-art reconciliation methods. Here, we will also
31 introduce the proposed loss function that is used to train the MLP neural network. Next
32 in Section 3, we explore the proposed NHR-TFNet framework in detail and highlight the
33 key model development strategies used in our architecture. In Section 4, we summarise
34 the benchmark datasets and explain the experimental setup used in this study to evaluate
35 NHR-TFNet against other hierarchical reconciliation methods. Section 5 summarises the
36 results obtained by NHR-TFNet on multiple hierarchical time series datasets and analyse
37 the main observations of this study. Finally, Section 6 concludes the paper.

38 2. Problem Statement and Related work

39
40 This section mathematically defines hierarchical time series and briefly outlines widely
41 used traditional forecast reconciliation approaches. We then introduce the loss function of
42 the non-linear forecast reconciliation approach proposed in this paper which is then followed
43 by a discussion of closely related relevant work in the literature.

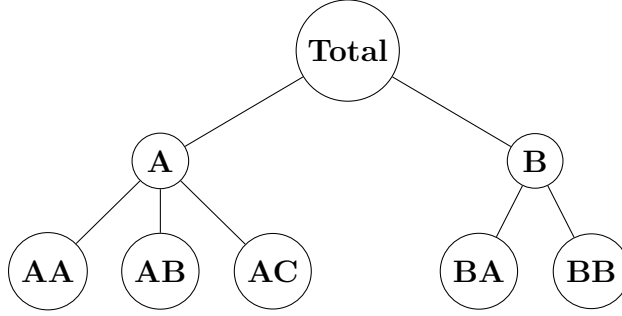


Figure 1: A two-level hierarchical structure.

2.1. Forecasting hierarchical time series

Let $\mathbf{b}_t \in \mathbb{R}^n$ be a vector of observations collected from the n series at the most disaggregated level of the structure at time t , $\mathbf{y}_t \in \mathbb{R}^m$ be a vector of all observations obtained by aggregating the series at the most disaggregated level at time t , for $t = 1, 2, \dots, T$, where T is the length of the series. These vectors hold the following relationship:

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

where \mathbf{S} is often referred to as the *summing matrix* of order $m \times n$, which consists of a set of linear constraints present in the structure. To elaborate on these notations, let's consider the structure given in Figure 1. For this structure, $m = 8, n = 5, \mathbf{b}_t = [y_{AA,t}, y_{AB,t}, y_{AC,t}, y_{BA,t}, y_{BB,t}]^\top, \mathbf{y}_t = [y_t, y_{A,t}, y_{B,t}, \mathbf{b}_t^\top]^\top$, and

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

where \mathbf{I}_k denotes an identity matrix of order $k \times k$. In general, a structure with linear aggregation constraints can be decomposed as $\mathbf{S} = \begin{bmatrix} \mathbf{C} \\ \mathbf{I}_n \end{bmatrix}$, where $\mathbf{C} \in \mathbb{R}^{m^* \times n}$, and $m^* = m - n$.

The point forecasts produced for all m series in the structure (also known as *base forecasts*) may not satisfy the aggregation constraints present in the data unless extremely simple forecasting methods are used. We denote $\hat{\mathbf{y}}_{T+h|T} \in \mathbb{R}^m$ as the vector consisting of h -steps-ahead *incoherent base forecasts* made using information available up to and including time T . To ensure that these forecasts satisfy the constraints, the base forecasts need to be adjusted and we called this process *forecast reconciliation*. It involves defining a mapping function ψ such that $\tilde{\mathbf{y}}_{T+h|T} = \psi(\hat{\mathbf{y}}_{T+h|T})$, where $\tilde{\mathbf{y}}_{T+h|T}$ satisfies the constraints and we refer to these as *reconciled forecasts*. The mapping function $\psi(\cdot)$ can be assumed to compose of two mappings $\psi(\hat{\mathbf{y}}_{T+h|T}) = s[g(\hat{\mathbf{y}}_{T+h|T})]$, where g combines the base forecasts to produce a new set of bottom level forecasts which are then aggregated by s .

When g is assumed to be linear, we can rewrite the mapping function to obtain the reconciled forecasts as $\tilde{\mathbf{y}}_{T+h|T} = \mathbf{S}\mathbf{G}\hat{\mathbf{y}}_{T+h|T}$, where \mathbf{G} depends on the forecast reconciliation

approach. Many existing point forecast reconciliation approaches in the literature belong to this category and they are summarized in Table 1. Among them, the first five choices satisfy the constraints $\mathbf{G}\mathbf{S} = \mathbf{I}$ which ensures that the reconciled forecasts are unbiased provided that the base forecasts are unbiased. In addition, the estimation of covariance matrices for MinT(Sample) and MinT(Shrink) assumes that the in-sample base forecast errors are jointly covariance stationary.

We could provide an alternative representation for the first five choices of the \mathbf{G} matrix where we could interpret the reconciliation process as an adjustment done to the bottom-level base forecasts as

$$\mathbf{G}\hat{\mathbf{y}}_{T+h|T} = \hat{\mathbf{b}}_{T+h|T} - \mathbf{J}\mathbf{X}\mathbf{U}^\top \hat{\mathbf{y}}_{T+h|T},$$

where \mathbf{X} depends on the forecast reconciliation approach and are given in Table 1, $\hat{\mathbf{b}}_{T+h|T}$ is the h -steps ahead base forecasts at the bottom level. and $\mathbf{U}^\top = \begin{bmatrix} \mathbf{I}_n & | & -\mathbf{C} \end{bmatrix}$. $\mathbf{U}^\top \hat{\mathbf{y}}_{T+h|T}$ quantifies the amount of aggregation inconsistency. For the hierarchy in Figure 1,

$$\mathbf{U}^\top \hat{\mathbf{y}}_{T+h|T} = \begin{bmatrix} \hat{y}_{T+h|T} - \hat{y}_{AA,T+h|T} - \hat{y}_{AB,T+h|T} - \hat{y}_{AC,T+h|T} - \hat{y}_{BA,T+h|T} - \hat{y}_{BB,T+h|T} \\ \hat{y}_{A,T+h|T} - \hat{y}_{AA,T+h|T} - \hat{y}_{AB,T+h|T} - \hat{y}_{AC,T+h|T} \\ \hat{y}_{B,T+h|T} - \hat{y}_{BA,T+h|T} - \hat{y}_{BB,T+h|T} \end{bmatrix}.$$

In theory, the matrix \mathbf{U} is the null space of \mathbf{S}^\top and the above choice is only one possibility. We have chosen this representation here because $\mathbf{U}^\top \hat{\mathbf{y}}_{T+h|T}$ is interpretable as the amount of aggregation inconsistency of the bottom-level forecasts.

Table 1: Traditional choices of the \mathbf{G} matrix.

Reconciliation method	\mathbf{G}	\mathbf{X}
BU (bottom-up)	$\begin{bmatrix} \mathbf{0}_{n \times m^*} & & \mathbf{I}_n \end{bmatrix}$	$\mathbf{0}$
OLS (ordinary least squares) [11]	$(\mathbf{S}^\top \mathbf{S})^{-1} \mathbf{S}^\top$	$\mathbf{U}(\mathbf{U}^\top \mathbf{U})^{-1}$
WLS (weighted least squares) [15]	$(\mathbf{S}^\top \hat{\mathbf{\Lambda}}_1^{-1} \mathbf{S})^{-1} \mathbf{S}^\top \hat{\mathbf{\Lambda}}_1^{-1}$	$\hat{\mathbf{\Lambda}}_1 \mathbf{U}(\mathbf{U}^\top \hat{\mathbf{\Lambda}}_1 \mathbf{U})^{-1}$
MinT(Sample) [30]	$(\mathbf{S}^\top \hat{\mathbf{W}}_{1,\text{sam}}^{-1} \mathbf{S})^{-1} \mathbf{S}^\top \hat{\mathbf{W}}_{1,\text{sam}}^{-1}$	$\hat{\mathbf{W}}_{1,\text{sam}} \mathbf{U}(\mathbf{U}^\top \hat{\mathbf{W}}_{1,\text{sam}} \mathbf{U})^{-1}$
MinT(Shrink) [30]	$(\mathbf{S}^\top \hat{\mathbf{W}}_{1,\text{shr}}^{-1} \mathbf{S})^{-1} \mathbf{S}^\top \hat{\mathbf{W}}_{1,\text{shr}}^{-1}$	$\hat{\mathbf{W}}_{1,\text{shr}} \mathbf{U}(\mathbf{U}^\top \hat{\mathbf{W}}_{1,\text{shr}} \mathbf{U})^{-1}$
EMinT-U [28]	$\mathbf{B}_h^\top \hat{\mathbf{Y}}_h (\hat{\mathbf{Y}}_h^\top \hat{\mathbf{Y}}_h)^{-1}$	

$\hat{\mathbf{W}}_{1,\text{sam}}$ and $\hat{\mathbf{W}}_{1,\text{shr}}$ are the sample, and shrinkage covariance matrix, respectively of 1-step-ahead in-sample base forecast errors. $\hat{\mathbf{\Lambda}}_1 = \text{diag}(\hat{\mathbf{W}}_{1,\text{sam}})$, $\mathbf{Y}_h = [\mathbf{y}_h, \mathbf{y}_{h+1}, \dots, \mathbf{y}_T]^\top \in \mathbb{R}^{(T-h+1) \times m}$, $\mathbf{B}_h = [\mathbf{b}_h, \mathbf{b}_{h+1}, \dots, \mathbf{b}_T]^\top \in \mathbb{R}^{(T-h+1) \times n}$, $\hat{\mathbf{Y}}_h = [\hat{\mathbf{y}}_{h|0}, \hat{\mathbf{y}}_{h+1|1}, \dots, \hat{\mathbf{y}}_{T|T-h}]^\top \in \mathbb{R}^{(T-h+1) \times m}$. $\text{diag}(\mathbf{A})$ constructs a diagonal matrix using the diagonal elements of the square matrix \mathbf{A} .

2.2. Proposed loss function

Motivated by this second interpretation of forecast reconciliation, we consider a non-linear mapping for g and are interested in minimise the following loss function:

$$\min_{\boldsymbol{\theta}} \frac{1}{T} \left[\sum_{t=1}^T \left\| \mathbf{b}_t - \tilde{\mathbf{b}}_{t|t-1} \right\|_2^2 + \lambda \left\| \mathbf{a}_t - \mathbf{C} \tilde{\mathbf{b}}_{t|t-1} \right\|_2^2 \right], \quad (1)$$

where $\tilde{\mathbf{b}}_{t|t-1} = \hat{\mathbf{b}}_{t|t-1} + \mathbf{g}(\mathbf{U}^\top \hat{\mathbf{y}}_{t|t-1}, \boldsymbol{\theta})$, $\lambda > 0$ is the penalty parameter, $\mathbf{g}(\cdot, \boldsymbol{\theta}) = [g_1(\cdot, \boldsymbol{\theta}_1), g_2(\cdot, \boldsymbol{\theta}_2), \dots, g_n(\cdot, \boldsymbol{\theta}_n)]^\top$, $g_j(\cdot, \boldsymbol{\theta}_j)$ is a non-linear mapping function with parameter vector $\boldsymbol{\theta}_j$ for $j = 1, 2, \dots, n$, $\boldsymbol{\theta} = [\boldsymbol{\theta}_1^\top, \boldsymbol{\theta}_2^\top, \dots, \boldsymbol{\theta}_n^\top]^\top$, $\mathbf{y}_t = [\mathbf{a}_t^\top, \mathbf{b}_t^\top]^\top$, and $\|\cdot\|_2$ is the l_2 -norm. The intuition behind this objective function is to non-linearly map a given set of base forecasts so that the mean squared forecast error of the bottom level is minimized while incorporating a penalty for not being able to minimize the mean squared forecast error of the aggregated series. We could rearrange the terms of the loss function and rewrite it as

$$\min_{\boldsymbol{\theta}} \frac{1}{T} \left[\sum_{t=1}^T \left\| \boldsymbol{\Lambda}^{1/2} (\mathbf{y}_t - \mathbf{S} \tilde{\mathbf{b}}_{t|t-1}) \right\|_2^2 \right], \quad (2)$$

where $\boldsymbol{\Lambda} = \begin{bmatrix} \lambda \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$ and it can be viewed as a nonlinear weighted least squares approach where the weights are given by $\boldsymbol{\Lambda}$. It is interesting to note that MinT(Sample) can be viewed as a special case of this approach. If restricted to linear mappings with no hidden layers in the MLP network and include conditions for unbiasedness, then the solution reduces to MinT(Sample) for any choice of λ .

After estimating the best non-linear mapping by minimizing the in-sample 1-step-ahead forecast errors as shown in Eq. (1), we obtain the coherent forecasts for the structure as

$$\tilde{\mathbf{y}}_{T+h|T} = \mathbf{S} \left[\hat{\mathbf{b}}_{T+h|T} + \mathbf{g}(\mathbf{U}^\top \hat{\mathbf{y}}_{T+h|T}, \hat{\boldsymbol{\theta}}) \right],$$

where $\hat{\boldsymbol{\theta}}$ is the solution of Eq. (1). For simplicity, we could use the same non-linear mapping for each forecast horizon. Alternatively, we could modify the proposed objective function to accommodate h -step-ahead in-sample forecast errors which leads to different non-linear mappings for each h .

2.3. Related literature

Burba and Chen [7] introduced a non-linear approach based on an encoder-decoder neural network, where the encoder is a trainable feed-forward neural network that takes base forecasts as the input and outputs the bottom-level reconciled forecasts, and the decoder is the summing matrix which obtains forecasts for all the levels using encoded bottom level forecasts. They have evaluated the performance of this approach using two loss functions: mean absolute scaled error and mean logarithm of absolute error and showed superior performances over traditional approaches when the hierarchies are complex. Shiratori et al. [23]

1
2
3 considered a loss function similar to Eq. (1) but it was used to estimate a prediction model
4 for the bottom-level series while incorporating upper-level forecasts through the penalty
5 function. They have used a different penalty parameter for each level in the hierarchy which
6 can be computationally costly for large structures. Using 2-level hierarchies constructed
7 from synthetic and real data, they have demonstrated that their method yields comparable
8 performances to BU and MinT(Shrink). Spiliotis et al. [24] proposed a non-linear reconcilia-
9 tion approach based on decision tree-based machine learning models: XGBoost and random
10 forest. Their approach models each of the bottom-level series separately on 1-step ahead
11 base forecasts obtained for each series in the structure for a holdout set. The loss function
12 used in this study is interested in minimizing the mean squared combination forecast error
13 for each bottom-level series and has not taken reconciliation error into account.
14
15
16
17
18

19 **3. NHR-TFNet Framework**

20
21 In this section, we describe the main components of the NHR-TFNet framework. NHR-
22 TFNet is composed of three components, namely: 1) the base-forecasting layer, 2) the
23 MLP training layer, and 3) the post-processing layer. In the following, we first discuss the
24 base-forecasting layer that generates base forecasts for each time series in the hierarchy and
25 provide fitted values of the base models. Next, we explain the MLP training layer, which
26 is the main learning component of NHR-TFNet. Then, we discuss the post-processing layer
27 of NHR-TFNet and describe the hyper-parameter selection method used in NHR-TFNet.
28 Finally, we demonstrate the forecast reconciliation stage of this trained network.
29
30
31

32 *3.1. Base-forecasting layer*

33
34 The base-forecasting layer produces base forecasts and model fits for each time series in
35 the hierarchy. In our experiments, we use both univariate and global models to generate base
36 forecasts. As the univariate models, we apply ETS and ARIMA implementations from the
37 `forecast` package [13], whereas for the global model, we use DeepAR [22] implementation
38 from the `GluonTS` package [1]. As global models are trained across a set of many time series,
39 prior to applying the DeepAR model, we first perform k -means clustering on the set of time
40 series features extracted from all the time series in the hierarchy. We use the `tsfeatures`
41 package developed by Hyndman et al. [16] to extract features from each time series. Based
42 on the time series features, we next perform the k -means clustering with a cluster size of
43 20 (i.e., the number of clusters). Although we have fixed the number of clusters here to
44 reduce the complexity of finding the optimal number of clusters and having a large number
45 of possible clusters, we can apply techniques such as the Elbow method to determine the
46 optimal number of clusters.
47
48
49
50

51 *3.2. MLP training layer*

52
53 For adjusting the base forecasts according to the aggregation constraints in the hierarchy,
54 we use an MLP based model. Neural Network based models have the support for producing
55 multiple outputs as a function of the same set of inputs. At time step t , the inputs $\bar{\mathbf{x}}_{t|t-1}$
56 to our network are the transformed fitted values from the base forecasting method. Given
57
58
59
60
61
62
63
64
65

a simple hierarchical structure with two levels as shown in Figure 1, the transformation of the fitted values can be done as per Eq. (3). In Eq. (3), $\bar{y}_{t|t-1}$, $\bar{y}_{A,t|t-1}$ and $\bar{y}_{B,t|t-1}$ for $t = 1, 2, \dots, T$ refer to the transformed fitted values. If the residuals of the base forecasting method are jointly covariance stationary, considering the transformed fitted values in this way allows us to feed reasonably stationary data as inputs to the MLP model. The corresponding outputs of the network $\tilde{\mathbf{b}}_{t|t-1}$ are the adjusted fitted values at the bottom-level as shown in Eq. (4).

$$\begin{aligned}\bar{y}_{t|t-1} &= \hat{y}_{t|t-1} - [\hat{y}_{AA,t|t-1} + \hat{y}_{AB,t|t-1} + \hat{y}_{AC,t|t-1} + \hat{y}_{BA,t|t-1} + \hat{y}_{BB,t|t-1}] \\ \bar{y}_{A,t|t-1} &= \hat{y}_{A,t|t-1} - [\hat{y}_{AA,t|t-1} + \hat{y}_{AB,t|t-1} + \hat{y}_{AC,t|t-1}] \\ \bar{y}_{B,t|t-1} &= \hat{y}_{B,t|t-1} - [\hat{y}_{BA,t|t-1} + \hat{y}_{BB,t|t-1}] \\ \bar{\mathbf{x}}_{t|t-1} &= [\bar{y}_{t|t-1}, \bar{y}_{A,t|t-1}, \bar{y}_{B,t|t-1}, \hat{y}_{AA,t|t-1}, \hat{y}_{AB,t|t-1}, \hat{y}_{AC,t|t-1}, \hat{y}_{BA,t|t-1}, \hat{y}_{BB,t|t-1}]^\top \quad (3)\end{aligned}$$

$$\tilde{\mathbf{b}}_{t|t-1} = [\tilde{y}_{AA,t|t-1}, \tilde{y}_{AB,t|t-1}, \tilde{y}_{AC,t|t-1}, \tilde{y}_{BA,t|t-1}, \tilde{y}_{BB,t|t-1}]^\top \quad (4)$$

Figure 2 shows the overall architecture of the MLP model used for the reconciliation in this work. The model contains a number of stacks each of which consists of a dense layer, a batch normalisation layer, rectified linear unit (ReLU) activation followed by a dropout layer. The batch normalisation step helps to speed up the neural network model training by standardizing the inputs to layers of the network while also handling internal covariate shift, which is the change of data distribution from one layer of the network to the next [17]. Also, both batch normalisation and the dropout layer help to reduce overfitting in the network. The dropout layer randomly drops nodes from the dense layers during the model training to ensure model generalisation and to reduce node inter-dependence [25]. A final dense layer, placed after the final stack of the network, having the same size as the number of time series in the bottom-most level of the hierarchy, reshapes the model outputs to the required final output shape.

The MLP model also uses a skip connection to the outputs which directly adds up the original fitted values of the bottom-level series to the outputs of the MLP model. Hence, the skip connection helps the model to learn only the adjustment that needs to be performed on top of the fitted values of the bottom-level series provided by the base forecasting method. Theoretically, this makes the learning process of the MLP model easier [9]. On top of the aforementioned input transformation, inputs are further transformed using min-max normalisation as well.

3.3. Post-processing layer

The outputs of the model, which are the adjusted fitted values of the bottom-most level in the hierarchy, are in the original scale. Once these values are retrieved, they are hierarchically aggregated in a bottom-up fashion to construct the values for all the levels of aggregation.

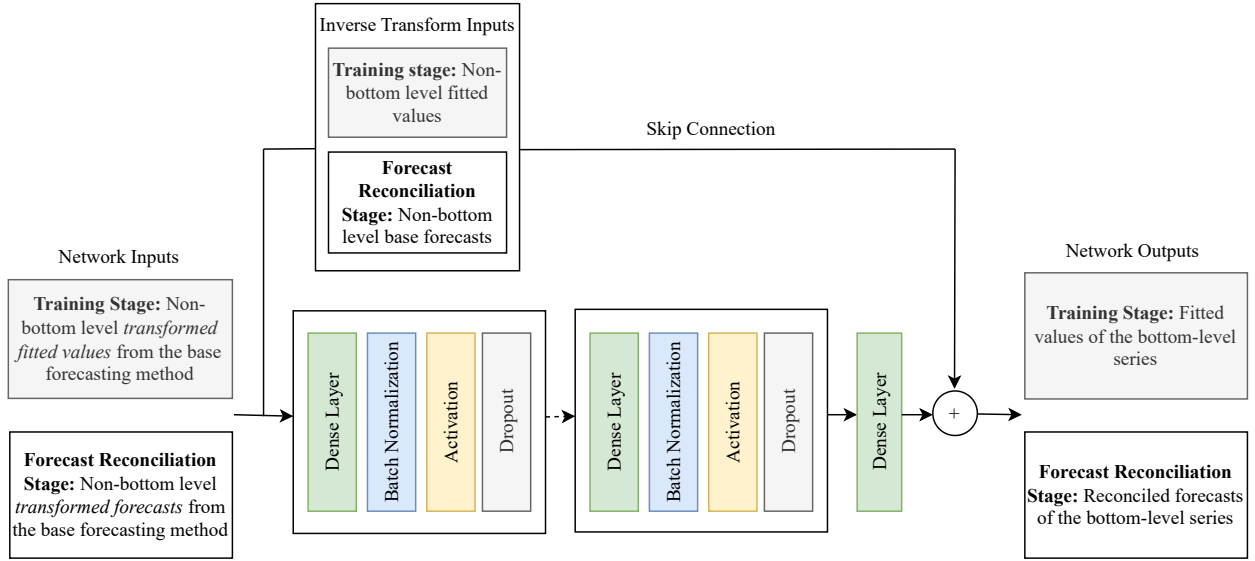


Figure 2: An overview of the proposed architecture for the NHR-TFNet framework, which consists of a base-forecasting layer, MLP training layer, and a post-processing layer.

Table 2: Minimum and maximum values for each hyper-parameter.

Hyper-parameter	Minimum	Maximum
Number of layers	1	5
Number of epochs	10	200
Dropout rate	0	0.5
Max normalization constraint	0	10
Learning rate	0.0001	0.1
Batch size	1	size of input data
λ (proposed loss function)	0.01	5

3.4. Hyper-parameter automation and optimisation

Several hyper-parameters need to be tuned for the MLP model. Table 2 shows the hyper-parameters, minimum and maximum values we use for each parameter. To find the best values for the hyper-parameters we use the Bayesian optimization algorithm [5] from `hyperopt` Python package. Once the optimal parameter values are found, we train the MLP using the optimal parameters five times with different seeds to account for the random weight initialization of the network. The average across the adjusted bottom-level fitted values from the five iterations is taken as the final adjusted bottom-level fitted values.

3.5. Forecast reconciliation stage

Once the network is trained, the base forecasts and transformed base forecasts (by following similar equations as transformed fitted values) are fed into the network to compute the reconciled forecasts at the bottom level. These forecasts are then summed appropriately to obtain the forecasts for other series in the structure.

Table 3: Description of the datasets.

Dataset	Frequency	Total number of			
		Levels	Series	Bottom-level series	Non-bottom level series
Prison	4 (quarterly)	5	121	64	57
Labour	4 (quarterly)	4	57	32	25
Tourism	12 (monthly)	3	85	77	8
Wikipedia	7 (weekly)	6	1095	913	182

4. Experiments

4.1. Datasets

We use four state-of-the-art hierarchical time series datasets for evaluation. We first pre-process the datasets using the `tsclean` function from the `forecast` package [13] to remove any outliers that maybe present in the data. Figure 3 shows the time series datasets after pre-processing. Table 3 shows the frequency of the time series, the number of levels, the total number of time series in the hierarchy, bottom-level and non-bottom level for the four datasets. We can observe that these datasets show varying levels of non-linear trends and seasonality. For the Wikipedia dataset, we can observe that some series at the bottom level are intermittent.

Prison Quarterly prison population in Australia from Q1 2005 to Q4 2016 [12]. The levels of the hierarchy are Australia (total), state, gender, legal status and indigenous status.

Labour Quarterly employed individuals in Australia from Q1 1987 to Q4 2018 [2]. The levels of the hierarchy are Australia (total), occupation category, employment status and gender.

Tourism Monthly domestic visitor nights in Australia from January 1998 to December 2019 [29]. The levels of the hierarchy are Australia (total), state and region.

Wikipedia Daily pageviews for the most popular social network articles on Wikipedia from 01-01-2016 to 29-06-2017 [19]. The levels of the hierarchy are access, agent, language, purpose and article.

4.2. Evaluation

4.2.1. Benchmark Methods

We evaluate the proposed NHR-TFNet with six state-of-the-art hierarchical forecasting benchmarks: BU, OLS, WLS, MinT(Sample), MinT(Shrink) and EMinT-U. For the prison and Wikipedia datasets, MinT(Sample) cannot be calculated as the number of time series in the hierarchy is greater than the number of observations. Therefore the comparisons for these two datasets are done using the other five benchmarks.

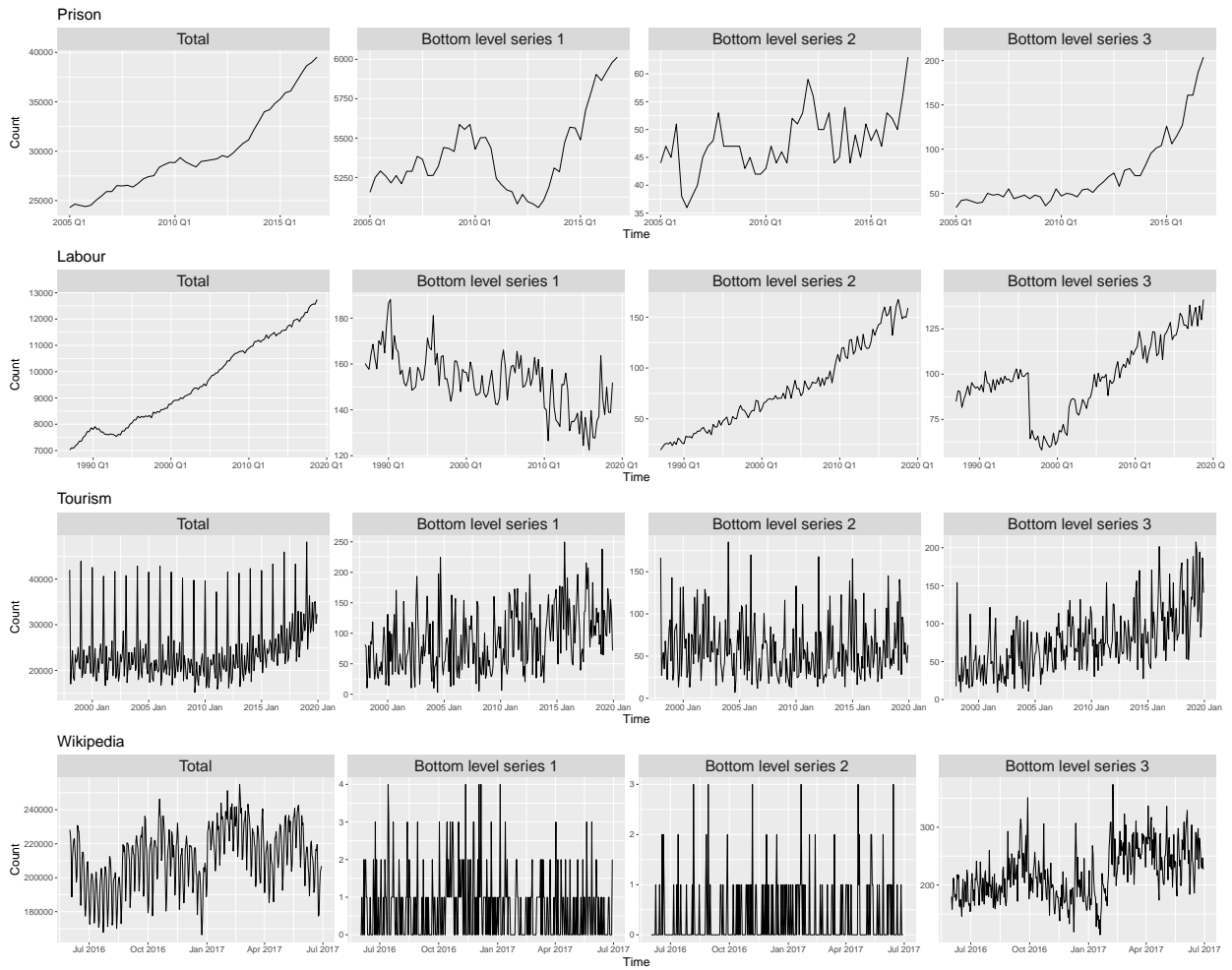


Figure 3: A few selected series at the top and bottom level of the hierarchy for each dataset used in this study.

4.2.2. Rolling Origin Evaluation

In this work to evaluate the forecasting performance, we conduct a rolling origin evaluation where the forecasting origin is updated at each rolling window and the forecasts are produced from each origin. An example of this evaluation is illustrated in Figure 4. To start the rolling origin evaluation, for each dataset we pick a reasonable training window size considering the trade-off between the number of rolling windows and the number of observations available to train the base forecasting methods and derive the fitted values to train the proposed MLP model. Table 4 shows the number of rolling windows and the number of observations available for the first rolling window (i.e., training window size of the first rolling window).

For each rolling window, we fit the base forecasting methods (ARIMA, ETS, or DeepAR) and compute the fitted values and 1-step-ahead forecasts. We next apply the hierarchical forecasting benchmarks and the proposed MLP using the fitted values and forecasts for a

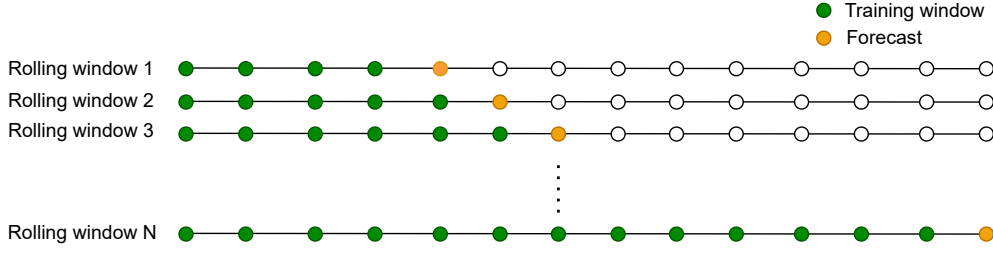


Figure 4: An illustration of the rolling origin evaluation.

Table 4: Number of rolling windows for each dataset and observations present for training in the first rolling window.

Dataset	Number of rolling windows	Min. window size
Prison	24	24
Labour	60	68
Tourism	120	144
Wikipedia	70	324

given rolling window. The same steps are conducted for all rolling windows of a dataset. For each rolling window, the mean squared error (MSE) is calculated and the average across series in different levels and average across all the series in the structure for a given dataset is obtained. Finally, the percentage relative improvement in average loss (PRIAL) is calculated as shown in Eq. (5) for each level in the structure and for the entire hierarchy (i.e., overall) using the average MSEs of the forecasts from the base forecasting methods (i.e., base forecasts) and the average MSEs after applying the reconciliation methods (i.e., reconciled forecasts). A positive value of PRIAL indicates the accuracy of the reconciled forecasts has increased while a negative value indicates the accuracy has decreased.

$$\text{PRIAL} = \frac{\text{MSE}(\text{base forecasts}) - \text{MSE}(\text{reconciled forecasts})}{\text{MSE}(\text{base forecasts})} \times 100\%. \quad (5)$$

5. Discussion

Table 5 reports the forecasting performance of the various reconciliation methods for the four datasets. The row labeled *Base MSE* reports the MSEs of the base forecasting methods and *Rank* reports the ranking of the reconciliation methods where rank one is given to the best-performing method whereas rank 6 or 7 is given to the least performing method. For the Prison, Labour and Tourism datasets, we set the λ value to take $[0.1, 0.9]$ whereas for the Wikipedia dataset, we set a slightly higher range i.e., $[0.01, 5]$ because some series at the bottom-level are intermittent so it may be better to give a high importance to the non-bottom level series.

According to Table 5, we observe that the proposed NHR-TFNet framework obtains the best PRIAL rank on the Prison, Labour, and Tourism datasets when DeepAR, ETS, and

Table 5: Average MSE of the base forecasts and PRIAL from reconciliation methods for the four datasets.

Prison		BU	OLS	WLS	MinT(Sample)	MinT(Shrink)	EMinT-U	NHR-TFNet
ARIMA								
	Base MSE ($\times 10^3$)							2.8
	PRIAL	-48.5	7.5	-2.3		-0.2	-83.5	-3.1
	Rank	5	1	3		2	6	4
ETS								
	Base MSE ($\times 10^3$)							3.0
	PRIAL	-46.6	6.8	3.8		10.3	-139.2	-45.5
	Rank	5	2	3		1	6	4
DeepAR								
	Base MSE ($\times 10^3$)							4.3
	PRIAL	23.6	4.7	21.9		23.0	-32.4	24.1
	Rank	2	5	4		3	6	1
Labour		BU	OLS	WLS	MinT(Sample)	MinT(Shrink)	EMinT-U	NHR-TFNet
ARIMA								
	Base MSE ($\times 10^2$)							3.9
	PRIAL	-6.8	3.8	2.2	-33.1	5.8	-141.6	-2.7
	Rank	5	2	3	6	1	7	4
ETS								
	Base MSE ($\times 10^2$)							3.8
	PRIAL	-2.8	2.2	2.0	-30.3	3.2	-198.1	3.5
	Rank	5	3	4	6	2	7	1
DeepAR								
	Base MSE ($\times 10^2$)							5.3
	PRIAL	15.0	5.0	12.4	-1.3	17.0	-116.9	13.4
	Rank	2	5	4	6	1	7	3
Tourism		BU	OLS	WLS	MinT(Sample)	MinT(Shrink)	EMinT-U	NHR-TFNet
ARIMA								
	Base MSE ($\times 10^4$)							5.6
	PRIAL	-71.7	2.5	-27.7	0.3	-11.9	-35.9	-3.1
	Rank	7	1	5	2	4	6	3
ETS								
	Base MSE ($\times 10^4$)							5.3
	PRIAL	-38.7	0.8	-17.7	2.1	-13.1	-116.8	1.0
	Rank	6	3	5	1	4	7	2
DeepAR								
	Base MSE ($\times 10^4$)							11.8
	PRIAL	-23.2	1.8	-9.8	16.1	12.6	-6.5	23.0
	Rank	7	4	6	2	3	5	1
Wikipedia		BU	OLS	WLS	MinT(Sample)	MinT(Shrink)	EMinT-U	NHR-TFNet
ARIMA								
	Base MSE ($\times 10^5$)							1.6
	PRIAL	4.4	2.1	8.8		13.4	-294.3	6.5
	Rank	4	5	2		1	6	3
ETS								
	Base MSE ($\times 10^5$)							1.2
	PRIAL	1.6	0.6	2.4	13	2.7	-604.4	-5.2
	Rank	3	4	2		1	6	5
DeepAR								
	Base MSE ($\times 10^5$)							5.0
	PRIAL	19.9	3.3	19.9		47.4	-49.0	15.5
	Rank	3	5	2		1	6	4

1
2
3
4 DeepAR models are used as the respective base forecasting models. For the Prison, Labour
5 and Wikipedia datasets, in general, the best-performing methods use a linear mapping
6 function such as MinT(Shrink). Notably, the proposed NHR-TFNet framework performs
7 among the top four reconciliation methods. This observation can be attributed to these
8 three datasets having a limited number of observations in comparison to the number of
9 bottom-level and non-bottom level series so learning a less complex linear mapping function
10 is beneficial. For the Tourism dataset, MinT(Sample) and NHR-TFNet show competitive
11 performances as this dataset has a reasonable number of observations in comparison to
12 the number of bottom-level and non-bottom level series. For all the datasets, the worst-
13 performing method is EMinT-U. This result is expected as this method works well when the
14 series are jointly covariance stationary which is not valid for these datasets. With respect to
15 base forecast models, we see that MSE is highest in DeepAR. This result is not surprising
16 given that all the datasets have a limited number of time series, thus learning a global model
17 such as DeepAR can be challenging in practice as they require a sufficient amount of time
18 series to estimate their numerous model parameters.

19
20 Table 6 summarizes the performance across various levels in the hierarchy for the Prison
21 dataset. The bold entries highlight the best-performing method and the underline numbers
22 highlight the second-best method. On average, for ARIMA base forecasts, OLS performs
23 the best, and MinT(Shrink) is the second best in most of the levels. For ETS and DeepAR
24 base forecasts, MinT(Shrink) is the best in most of the levels except the top-level. This
25 observed pattern is aligned with the rankings noted in Table 5. To save space we report
26 PRIAL for each level in [Appendix A](#) for the remaining datasets.

32 33 6. Conclusions and Future Work

34
35 Widely used traditional methods for forecast reconciliation use a linear function to map
36 all the base forecasts into the bottom-level series which are then added up by a summing
37 matrix to produce coherent forecasts for the whole hierarchy. In this paper, we explored the
38 forecasting performance when a non-linear mapping function is used instead. We trained a
39 feed-forward neural network with a skip connection where an adjustment estimated using
40 transformed fitted values is added to the original base forecasts at the bottom-level series
41 which are then summed to produce forecasts for the upper-level series. We evaluated the pro-
42 posed methodology using four benchmark datasets. The results indicated that NHR-TFNet
43 showed competitive results for the datasets having a reasonable number of observations
44 compared to the number of bottom-level and non-bottom level series.

45
46 An important perspective of this study is to explore a slightly modified version of the
47 proposed loss function. It attempts to find a non-linear mapping by minimizing the in-sample
48 1-step-ahead mean squared error of the whole structure while penalizing for incoherence in
49 the forecasts. However, this approach will produce near coherent forecasts, and the deviation
50 from coherence depends on the tuning parameter.

Table 6: Average MSEs of base forecasts and PRIAL from reconciliation methods for the Prison dataset.

Method	Australia	State	Gender	Legal	Indigenous
ARIMA					
BU	-125.3	-46.6	-37.4	-4.4	0.0
OLS	5.7	3.3	4.0	13.2	11.2
WLS	-26.0	-3.5	<u>2.8</u>	11.0	<u>13.3</u>
MinT(Shrink)	-19.7	<u>-2.6</u>	2.3	<u>12.1</u>	14.5
EMinT-U	-101.4	-79.1	-79.4	-66.8	-85.9
NHR-TFNet	<u>-9.1</u>	-9.3	-4.5	3.9	5.4
Base MSE ($\times 10^3$)	<i>88.1</i>	<i>7.9</i>	<i>3.5</i>	<i>2.3</i>	<i>0.9</i>
ETS					
BU	-105.8	-40.6	-30.4	-6.4	0.0
OLS	<u>4.7</u>	<u>4.5</u>	8.1	9.0	<u>9.8</u>
WLS	-4.7	3.0	<u>9.3</u>	<u>9.3</u>	9.6
MinT(Shrink)	6.5	8.2	13.3	13.4	13.7
EMinT-U	-154.6	-122.4	-133.1	-140.2	-133.5
NHR-TFNet	-83.5	-43.0	-30.8	-25.1	-11.6
Base MSE ($\times 10^3$)	<i>113.7</i>	<i>8.6</i>	<i>3.8</i>	<i>2.2</i>	<i>0.8</i>
DeepAR					
BU	<u>43.1</u>	<u>12.5</u>	9.0	0.9	0.0
OLS	11.4	-1.3	0.6	-2.4	-2.2
WLS	39.1	11.6	<u>9.4</u>	2.2	<u>1.1</u>
MinT(Shrink)	39.1	11.8	9.8	6.4	5.9
EMinT-U	-9.4	-34.8	-45.2	-63.9	-77.9
NHR-TFNet	46.9	8.7	5.7	0.0	-0.4
Base MSE ($\times 10^3$)	<i>241.9</i>	<i>10.7</i>	<i>4.4</i>	<i>2.1</i>	<i>0.8</i>

7. Acknowledgment

This research was funded by an IIF-SAS research award and by an Infrastructure Data Science research award at Meta. The research was undertaken using the LIEF HPC-GPGPU Facility hosted at the University of Melbourne. This Facility was established with the assistance of LIEF Grant LE170100200

References

- [1] Alexandrov, A., Benidis, K., Bohlke-Schneider, M., Flunkert, V., Gasthaus, J., Januschowski, T., Maddix, D.C., Rangapuram, S., Salinas, D., Schulz, J., Stella, L., Türkmen, A.C., Wang, Y., 2020. GluonTS: Probabilistic and neural time series modeling in Python. *Journal of Machine Learning Research* 21, 1–6. URL: <http://jmlr.org/papers/v21/19-820.html>.
- [2] Australian Bureau of Statistics, 2022. Employed persons by age, occupation sub-major group of main job (ANZSCO) and sex. Accessed on 12-06-2022.
- [3] Bandara, K., Bergmeir, C., Smyl, S., 2020. Forecasting across time series databases using recurrent neural networks on groups of similar series: A clustering approach. *Expert Systems with Applications* 140, 112896.
- [4] Bandara, K., Hewamalage, H., Godahewa, R., Gamakumara, P., 2022. A fast and scalable ensemble of global models with long memory and data partitioning for the M5 forecasting competition. *International Journal of Forecasting* 38, 1400–1404.

- 1
2
3
4 [5] Bergstra, J., Yamins, D., Cox, D., 2013. Making a science of model search: hyperparameter optimization
5 in hundreds of dimensions for vision architectures, in: International Conference on Machine Learning,
6 PMLR. pp. 115–123. URL: <http://hyperopt.github.io/hyperopt>.
- 7 [6] Box, G.E.P., Jenkins, G.M., Reinsel, G.C., Ljung, G.M., 2015. Time Series Analysis: Forecasting and
8 Control. John Wiley & Sons.
- 9 [7] Burba, D., Chen, T., 2021. A trainable reconciliation method for hierarchical time-series.
10 [arXiv:2101.01329](https://arxiv.org/abs/2101.01329).
- 11 [8] Han, X., Dasgupta, S., Ghosh, J., 2021. Simultaneously reconciled quantile forecasting of hierarchically
12 related time series, in: Proceedings of the 24th International Conference on Artificial Intelligence and
13 Statistics (AISTATS), San Diego, CA, USA.
- 14 [9] He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: 2016 IEEE
15 Conference on Computer Vision and Pattern Recognition (CVPR), pp. 770–778. doi:[10.1109/CVPR.
16 2016.90](https://doi.org/10.1109/CVPR.2016.90).
- 17 [10] Hewamalage, H., Bergmeir, C., Bandara, K., 2021. Recurrent neural networks for time series forecasting:
18 current status and future directions. *International Journal of Forecasting* 37, 388–427.
- 19 [11] Hyndman, R.J., Ahmed, R.A., Athanasopoulos, G., Shang, H.L., 2011. Optimal combination forecasts
20 for hierarchical time series. *Computational Statistics & Data Analysis* 55, 2579–2589.
- 21 [12] Hyndman, R.J., Athanasopoulos, G., 2021. Forecasting: Principles and Practice. OTexts: Melbourne,
22 Australia. OTexts.com/fpp3. Accessed on 12-06-2022.
- 23 [13] Hyndman, R.J., Khandakar, Y., 2008. Automatic time series forecasting: the forecast package for R.
24 *Journal of Statistical Software* 26, 1–22.
- 25 [14] Hyndman, R.J., Koehler, A.B., Ord, K.J., Snyder, R.D., 2008. Forecasting with Exponential Smooth-
26 ing: The State Space Approach. Springer Science & Business Media.
- 27 [15] Hyndman, R.J., Lee, A.J., Wang, E., 2016. Fast computation of reconciled forecasts for hierarchical
28 and grouped time series. *Computational Statistics & Data Analysis* 97, 16–32.
- 29 [16] Hyndman, R.J., Wang, E., Kang, Y., Talagala, T., 2021. tsfeatures: Time series feature extraction.
30 URL: <https://github.com/robjhyndman/tsfeatures/>. r package version 0.1.
- 31 [17] Ioffe, S., Szegedy, C., 2015. Batch normalization: accelerating deep network training by reducing internal
32 covariate shift, in: Proceedings of the 32nd International Conference on International Conference
33 on Machine Learning, JMLR.org. pp. 448–456.
- 34 [18] Januschowski, T., Gasthaus, J., Wang, Y., Salinas, D., Flunkert, V., Bohlke-Schneider, M., Callot, L.,
35 2020. Criteria for classifying forecasting methods. *International Journal of Forecasting* 36, 167–177.
- 36 [19] Mahsa Ashouri, R.J.H., Shmueli, G., 2022. Fast forecast reconciliation using linear models. *Journal of
37 Computational and Graphical Statistics* 31, 263–282. doi:[10.1080/10618600.2021.1939038](https://doi.org/10.1080/10618600.2021.1939038).
- 38 [20] Makridakis, S., Spiliotis, E., Assimakopoulos, V., 2022. M5 accuracy competition: results, findings,
39 and conclusions. *International Journal of Forecasting* 38, 1346–1364.
- 40 [21] Rangapuram, S.S., Werner, L.D., Benidis, K., Mercado, P., Gasthaus, J., Januschowski, T., 2021. End-
41 to-End learning of coherent probabilistic forecasts for hierarchical time series, in: Meila, M., Zhang, T.
42 (Eds.), Proceedings of the 38th International Conference on Machine Learning, PMLR. pp. 8832–8843.
- 43 [22] Salinas, D., Flunkert, V., Gasthaus, J., Januschowski, T., 2020. DeepAR: probabilistic forecasting with
44 autoregressive recurrent networks. *International Journal of Forecasting* 36, 1181–1191.
- 45 [23] Shiratori, T., Kobayashi, K., Takano, Y., 2020. Prediction of hierarchical time series using structured
46 regularization and its application to artificial neural networks. *PLoS ONE* 15, 1–23.
- 47 [24] Spiliotis, E., Abolghasemi, M., Hyndman, R.J., Petropoulos, F., Assimakopoulos, V., 2021. Hierarchical
48 forecast reconciliation with machine learning. *Applied Soft Computing* 112, 107756.
- 49 [25] Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R., 2014. Dropout: a simple
50 way to prevent neural networks from overfitting. *Journal of Machine Learning Research* 15, 1929–1958.
- 51 [26] Suilin, A., 2018. Kaggle web traffic. <https://github.com/Arturus/kaggle-web-traffic>. Accessed:
52 2020-02-10.
- 53 [27] Taieb, S.B., Taylor, J.W., Hyndman, R.J., 2021. Hierarchical probabilistic forecasting of electricity
54 demand with smart meter data. *Journal of the American Statistical Association* 116, 27–43.
- 55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

[28] Wickramasuriya, S.L., 2021. Properties of point forecast reconciliation approaches. [arXiv:2103.11129](https://arxiv.org/abs/2103.11129).
[29] Wickramasuriya, S.L., 2023. Probabilistic forecast reconciliation under the gaussian framework. *Journal of Business & Economic Statistics* , 1–14.
[30] Wickramasuriya, S.L., Athanasopoulos, G., Hyndman, R.J., 2019. Optimal forecast reconciliation for hierarchical and grouped time series through trace minimization. *Journal of the American Statistical Association* 114, 804–819.

Appendix A. Level-wise forecasting performances of the datasets

Table A.1: Average MSEs of the base forecasts and PRIAL from reconciliation methods for the Labour dataset.

Method	Total Employment	Occupation	Employment Status	Gender
ARIMA				
BU	-36.1	-3.3	0.7	0.0
OLS	<u>3.5</u>	<u>4.6</u>	4.6	2.3
WLS	-7.0	3.6	<u>5.3</u>	<u>2.9</u>
MinT(Sample)	-46.8	-30.7	-28.7	-32.1
MinT(Shrink)	4.6	5.3	7.8	4.7
EMinT-U	-145.4	-143.2	-140.4	-138.7
NHR-TFNet	-18.3	-2.3	2.5	1.2
Base MSE ($\times 10^2$)	<i>37.6</i>	<i>7.8</i>	<i>4.2</i>	<i>1.8</i>
ETS				
BU	-14.5	-3.5	2.5	0.0
OLS	3.9	1.6	3.2	0.4
WLS	3.1	<u>1.0</u>	3.7	0.1
MinT(Sample)	-11.3	-39.4	-28.7	-35.1
MinT(Shrink)	<u>7.7</u>	<u>1.0</u>	4.6	<u>0.9</u>
EMinT-U	-194.0	-201.2	-198.6	-196.7
NHR-TFNet	12.1	-0.8	<u>4.5</u>	1.3
Base MSE ($\times 10^2$)	<i>37.6</i>	<i>7.6</i>	<i>4.0</i>	<i>1.7</i>
DeepAR				
BU	27.9	21.5	<u>2.8</u>	0.0
OLS	7.4	14.0	-2.4	-3.1
WLS	20.9	20.5	1.8	-0.4
MinT(Sample)	9.7	7.2	-12.4	-17.5
MinT(Shrink)	<u>26.0</u>	24.4	7.1	4.1
EMinT-U	-144.2	-87.4	-117.7	-118.3
NHR-TFNet	14.9	<u>23.8</u>	7.1	<u>3.1</u>
Base MSE ($\times 10^2$)	<i>88.6</i>	<i>10.8</i>	<i>4.3</i>	<i>1.8</i>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Table A.2: Average MSEs of the base forecasts and PRIAL from reconciliation methods for the Tourism dataset.

Method	Australia	States	Regions
ARIMA			
BU	-122.2	-31.0	0.0
OLS	-2.5	<u>9.1</u>	5.1
WLS	-54.8	-3.1	<u>6.4</u>
MinT(Sample)	<u>-9.5</u>	11.8	8.1
MinT(Shrink)	-29.7	5.7	8.1
EMinT-U	-44.4	-24.5	-31.4
NHR-TFNet	-12.6	7.6	5.1
Base MSE ($\times 10^4$)	<i>242.2</i>	<i>21.2</i>	<i>1.1</i>
ETS			
BU	-59.2	-20.4	0.0
OLS	-1.2	3.4	3.3
WLS	-29.4	-6.5	<u>3.4</u>
MinT(Sample)	2.0	<u>2.2</u>	2.4
MinT(Shrink)	-22.6	-3.8	3.8
EMinT-U	-127.7	-103.5	-102.4
NHR-TFNet	<u>0.3</u>	1.3	2.9
Base MSE ($\times 10^4$)	<i>249.4</i>	<i>18.1</i>	<i>0.9</i>
DeepAR			
BU	-36.6	0.3	0.0
OLS	1.0	3.7	1.2
WLS	-19.0	7.6	2.5
MinT(Sample)	<u>13.2</u>	<u>25.8</u>	<u>9.1</u>
MinT(Shrink)	8.6	22.9	10.2
EMinT-U	-4.4	-2.5	-30.6
NHR-TFNet	22.5	30.5	6.2
Base MSE ($\times 10^4$)	<i>635.7</i>	<i>37.7</i>	<i>1.3</i>

Table A.3: Average MSEs of the base forecasts and PRIAL from reconciliation methods for the Wikipedia dataset.

Method	Total	Access	Agent	Language	Purpose	Article
ARIMA						
BU	7.6	8.6	-2.7	3.9	6.3	0.0
OLS	4.4	6.3	-2.5	1.6	0.1	0.0
WLS	11.9	<u>14.2</u>	<u>3.7</u>	<u>8.9</u>	<u>7.1</u>	<u>1.5</u>
MinT(Shrink)	15.2	19.2	9.2	14.8	11.2	4.7
EMinT-U	-316.0	-258.0	-285.0	-293.5	-313.0	-323.0
NHR-TFNet	<u>13.0</u>	10.9	-1.0	6.7	3.5	0.0
Base MSE ($\times 10^4$)	<i>3811.9</i>	<i>1239.1</i>	<i>670.2</i>	<i>137.1</i>	<i>14.1</i>	<i>1.7</i>
ETS						
BU	-0.4	4.9	4.6	0.8	-3.3	0.0
OLS	-1.9	2.8	2.3	-0.2	-1.1	0.5
WLS	<u>0.4</u>	<u>5.2</u>	<u>5.5</u>	<u>1.5</u>	-2.0	<u>1.3</u>
MinT(Shrink)	0.8	5.5	5.8	1.9	-1.7	1.7
EMinT-U	-694.0	-576.5	-558.6	-621.1	-631.7	-512.0
NHR-TFNet	-8.1	-4.4	-4.0	-4.6	-8.7	-0.2
Base MSE ($\times 10^4$)	<i>2799.7</i>	<i>897.6</i>	<i>552.0</i>	<i>100.8</i>	<i>10.6</i>	<i>1.5</i>
DeepAR						
BU	33.6	15.8	18.2	21.0	8.6	0.0
OLS	17.4	-2.0	0.5	1.5	-2.2	-2.2
WLS	<u>40.0</u>	14.9	16.8	17.9	6.1	-3.0
MinT(Shrink)	62.3	46.3	47.9	45.3	30.8	8.1
EMinT-U	-52.9	-30.3	-27.0	-48.6	-107.9	-183.3
NHR-TFNet	-8.2	<u>22.4</u>	<u>25.7</u>	<u>25.0</u>	<u>13.3</u>	<u>1.3</u>
Base MSE ($\times 10^4$)	<i>11921.9</i>	<i>4766.0</i>	<i>2841.5</i>	<i>404.1</i>	<i>30.9</i>	<i>2.1</i>