

Regularized Regression for Hierarchical Forecasting Without Unbiasedness Conditions

Souhaib Ben Taieb*

University of Mons
Mons, Belgium
souhaib.bentaieb@umons.ac.be
Monash University
Melbourne, Australia
souhaib.bentaieb@monash.edu

Bonsoo Koo

Monash University
Melbourne, Australia
bonsoo.koo@monash.edu

ABSTRACT

Forecasting a large set of time series with hierarchical aggregation constraints is a central problem for many organizations. However, it is particularly challenging to forecast these hierarchical structures. In fact, it requires not only good forecast accuracy at each level of the hierarchy, but also the coherency between different levels, i.e. the forecasts should satisfy the hierarchical aggregation constraints. Given some incoherent base forecasts, the state-of-the-art methods compute revised forecasts based on forecast combination which ensures that the aggregation constraints are satisfied. However, these methods assume the base forecasts are unbiased and constrain the revised forecasts to be also unbiased. We propose a new forecasting method which relaxes these unbiasedness conditions, and seeks the revised forecasts with the best tradeoff between bias and forecast variance. We also present a regularization method which allows us to deal with high-dimensional hierarchies, and provide its theoretical justification. Finally, we compare the proposed method with the state-of-the-art methods both theoretically and empirically. The results on both simulated and real-world data indicate that our methods provide competitive results compared to the state-of-the-art methods.

CCS CONCEPTS

• **Mathematics of computing** → **Time series analysis**; • **Computing methodologies** → *Regularization*.

KEYWORDS

time series, hierarchical forecasting, regularization, sparsity

ACM Reference Format:

Souhaib Ben Taieb and Bonsoo Koo. 2019. Regularized Regression for Hierarchical Forecasting Without Unbiasedness Conditions. In *The 25th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '19)*.

*Souhaib is an Adjunct Senior Research Fellow in the Department of Econometrics and Business Statistics at Monash University.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

KDD '19, August 4–8, 2019, Anchorage, AK, USA

© 2019 Copyright held by the owner/author(s). Publication rights licensed to ACM.

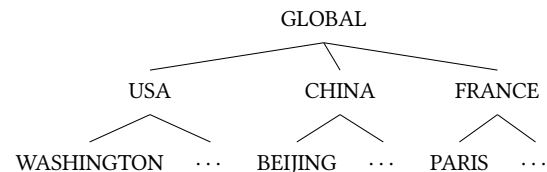
ACM ISBN 978-1-4503-6201-6/19/08...\$15.00

<https://doi.org/10.1145/3292500.3330976>

August 4–8, 2019, Anchorage, AK, USA. ACM, New York, NY, USA, 11 pages.
<https://doi.org/10.1145/3292500.3330976>

1 INTRODUCTION

Forecasting a large set of time series with hierarchical aggregation constraints is a central problem for many organizations. For example, the store-level sales of a multinational company can be grouped in a hierarchical structure composed of countries, cities and regions [13]. Smart grid is another context where a hierarchy with different levels of aggregation naturally arises. The bottom level can include demand for single households or buildings and supply from single solar panels or wind turbines, while the intermediate levels could be the total demand and supply of entire regions, states and countries as a whole [1]. The next Figure gives an example of hierarchical aggregation.



It is particularly challenging to forecast these hierarchical structures (see Section 2.1). In fact, in addition to providing good forecast accuracy, it is often essential for decision-making purposes to produce “coherent” forecasts, i.e. the forecasts of aggregates should be equal to the sum of the corresponding disaggregated forecasts. Unfortunately, independently forecasting all the series in a hierarchy, also called base forecasts, is unlikely to deliver coherent forecasts. Another feature of these hierarchical structures is the highly noisy series in the most disaggregated bottom level. As a result, computing bottom-up forecasts, i.e. summing the individual bottom-level forecasts is unlikely to provide good accuracy at upper levels in the hierarchy [11].

The state-of-the-art methods for hierarchical forecasting combine both base and bottom-up forecasts [10]. Specifically, starting from (probably) incoherent base forecasts, coherent revised forecasts are computed by applying forecast combination and using a bottom-up procedure. In other words, the base forecasts are adjusted so that they become coherent.

[17] proposed a new forecasting method called MinT which computes the optimal combination weights for the revised forecasts. More precisely, MinT assumes the base forecasts are unbiased and seeks the best combination weights under the constraint that the

revised forecasts will also be unbiased. In other words, MinT seeks the revised forecasts with minimum variance among all the unbiased revised forecasts (see Section 2.2). Although the MinT optimal combination weights are available in closed-form, the two conditions are hard to justify in practice. In fact, the base forecasts can be biased and revised forecasts with low mean squared forecast errors is often more important than the unbiasedness property.

We make the following contributions to hierarchical forecasting:

- First, we propose a new forecasting method which does not assume or constrain the base or the revised forecasts to be unbiased. To do so, we formulate the hierarchical forecasting problem as an empirical risk minimization problem which directly minimizes the mean squared forecast errors. We also provide the associated closed-form solution (see Section 3.1).
- Second, we address the problem of forecasting high dimensional hierarchies, i.e. hierarchies for which the total number of series is much larger than the number of historical observations. We adjust our objective function by adding a regularization term and provide its theoretical justification. We also derive its asymptotic properties including its limiting distribution (see Section 3.2).
- Third, we provide a comparison between our methods and the state-of-the-art MinT method, including similarities and differences, as well as the asymptotic properties of the methods. For example, we show that MinT is a particular case of our method, and under certain conditions, it is asymptotically equivalent to our method (see Section 4).
- Finally, we evaluate and compare the different hierarchical forecasting methods using both simulated and real-world data sets (see Section 5).

2 PRELIMINARIES

2.1 Hierarchical Forecasting

We let $\mathbf{b}_t \in \mathbb{R}^m$ represent the observations at time t for the m series in the most disaggregated (bottom) level. Then $\mathbf{a}_t = \mathbf{A}\mathbf{b}_t \in \mathbb{R}^k$ contains the observations at time t for the k aggregated series, where $\mathbf{A} \in \{0, 1\}^{k \times m}$. Each entry A_{ij} is equal to 1 if the i th aggregate series contains the j th bottom-level series, where $i = 1, \dots, k$ and $j = 1, \dots, m$. Finally, $\mathbf{y}_t = \mathbf{S}\mathbf{b}_t = (\mathbf{a}'_t \mathbf{b}'_t) \in \mathbb{R}^n$ contains observations both at the aggregate and bottom levels, where $\mathbf{S}' = [\mathbf{A}' \mathbf{I}_m]$ and \mathbf{I}_m is an $m \times m$ identity matrix. To avoid pathological hierarchies, we will assume that $m \geq 2$, $k \geq 1$ and $\sum_{j=1}^m A_{i,j} > 1$ for $i = 1, \dots, k$.

Figure 1 gives the tree and the matrix representation of a hierarchical time series with $m = 4$ bottom level series with $\mathbf{b}_t = (y_{aa,t}, y_{ab,t}, y_{ba,t}, y_{bb,t})'$, and $k = 3$ aggregate series with $\mathbf{a}_t = (y_t, y_{a,t}, y_{b,t})'$, $y_t = y_{a,t} + y_{b,t}$, $y_{a,t} = y_{aa,t} + y_{ab,t}$ and $y_{b,t} = y_{ba,t} + y_{bb,t}$.

While $\mathcal{I}_T = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_T\}$ denotes data observed up to time T , let $\hat{\mathbf{y}}_T(h)$ denote h -period ahead forecasts based on \mathcal{I}_T and $h = 1, \dots, H$ where h is the forecast horizon. The optimal h -period ahead forecasts which minimizes the conditional expectation:

$$E \left[\|\mathbf{y}_{T+h} - \hat{\mathbf{y}}_T(h)\|_2^2 | \mathcal{I}_T \right], \quad (1)$$

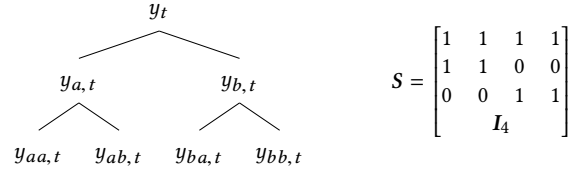


Figure 1: The tree and the matrix representation of a hierarchical time series with $m = 4$ bottom-level series and $k = 3$ aggregate series.

where $\|\cdot\|_p$ denotes an L_p norm over the relevant space, are given by

$$\boldsymbol{\mu}_T(h) = E[\mathbf{y}_{T+h} | \mathcal{I}_T] = \mathbf{S} E[\mathbf{b}_{T+h} | \mathcal{I}_T]. \quad (2)$$

See [6] for more details on optimal forecasts.

In other words, if (1) is the error measure, any hierarchical forecasting method aims to provide the best estimate of $\boldsymbol{\mu}_T(h)$.

A natural plug-in estimator for expression (2), known as *bottom-up* (BU), is given by

$$\hat{\mathbf{y}}_T(h) = \mathbf{S} \hat{\mathbf{b}}_T(h), \quad (3)$$

where $\hat{\mathbf{b}}_T(h)$ is an estimate of $E[\mathbf{b}_{T+h} | \mathcal{I}_T]$. The bottom-up method forecasts each of the bottom-level series, and then uses simple aggregation to obtain forecasts at higher levels of the hierarchy. Although this method is conceptually simple, it can provide poor forecasts on highly disaggregated and noisy data, especially at upper levels of the hierarchy.

Another approach, known as *base*, consists in forecasting all series in the hierarchy without considering the hierarchical constraints. In other words, we compute

$$\hat{\mathbf{y}}_T(h) = (\hat{\mathbf{a}}_T(h)' \hat{\mathbf{b}}_T(h)')', \quad (4)$$

where $\hat{\mathbf{a}}_T(h)$ and $\hat{\mathbf{b}}_T(h)$ are estimates of $E[\mathbf{a}_{T+h} | \mathcal{I}_T]$ and $E[\mathbf{b}_{T+h} | \mathcal{I}_T]$, respectively. This approach has the advantage that the typical poor forecasts for highly noisy series at the bottom level will not affect the forecasts for the series at upper levels of the hierarchy. However, it is unlikely that the resulting set of forecasts will be coherent, i.e. satisfy the aggregation constraints $\hat{\mathbf{a}}_T(h) = \mathbf{A}\hat{\mathbf{b}}_T(h)$. Imposing the aggregation constraints seems sensible since the optimal forecasts in expression (2) are coherent by definition. Furthermore, coherent forecasts will allow coherent decisions over the entire hierarchy.

2.2 The MinT approach to hierarchical forecasting

The state-of-the-art method for hierarchical forecasting combines the best of both *bottom-up* and *base* forecasts, given in (3) and (4), respectively. Starting from a set of (probably incoherent) base forecasts $\hat{\mathbf{y}}_T(h)$, [10] proposed the computation of revised forecasts

$$\tilde{\mathbf{y}}_T(h) = \mathbf{S}\mathbf{P}\hat{\mathbf{y}}_T(h), \quad (5)$$

for some appropriately chosen matrix $\mathbf{P} \in \mathcal{P} \subseteq \mathbb{R}^{m \times n}$ where \mathcal{P} is the domain of \mathbf{P} . This approach has multiple advantages: (1) the forecasts are coherent by construction; (2) a forecast combination from all levels is applied through the weight matrix \mathbf{P} ; and (3) many

hierarchical forecasting methods are represented as particular cases, including the bottom-up forecasts, for which

$$P = [\mathbf{0}_{m \times k} \mid I_m], \quad (6)$$

where $\mathbf{0}_{m \times k}$ is an $m \times k$ zero matrix.

More recently, [17] proposed to optimally combine the base forecasts. Specifically, assuming the base forecasts $\hat{\mathbf{y}}_T(h)$ are unbiased, [17] computes the weight P giving the minimum variance unbiased revised forecasts. In other words, assuming $E[\hat{\mathbf{y}}_T(h)|\mathcal{I}_T] = E[\mathbf{y}_{T+h}|\mathcal{I}_T]$, denoted as **A1**, [17] considered the following optimization problem:

$$\begin{aligned} \min_{P \in \mathcal{P}} E[\|\mathbf{y}_{T+h} - \tilde{\mathbf{y}}_T(h)\|_2^2 | \mathcal{I}_T] \\ \text{subject to } E[\tilde{\mathbf{y}}_T(h)|\mathcal{I}_T] = E[\mathbf{y}_{T+h}|\mathcal{I}_T] \quad (\mathbf{C1}). \end{aligned} \quad (7)$$

Under the unbiasedness conditions, **A1** and **C1**, [17] showed that minimizing (7) reduces to the following problem:

$$\min_{P \in \mathcal{P}} \text{Tr}(\text{Var}[\mathbf{y}_{T+h} - \tilde{\mathbf{y}}_T(h)|\mathcal{I}_T]) \text{ subject to } SPS = S, \quad (8)$$

where $\text{Tr}(\cdot)$ denotes the trace of a matrix.

The previous optimization problem has a closed-form solution given by

$$P^* = (S'W_c^{-1}S)^{-1}S'W_c^{-1}, \quad (9)$$

where W_c^{-1} is the Moore-Penrose generalized inverse of $W_c = E[\hat{\mathbf{e}}_T(h)\hat{\mathbf{e}}_T(h)'|\mathcal{I}_T]$ and $\hat{\mathbf{e}}_T(h) = \mathbf{y}_{T+h} - \hat{\mathbf{y}}_T(h)$. We refer to this approach as the infeasible MinT (minimum trace) reconciliation given that W_c is unavailable in practice.

In [17], the authors proposed to compute

$$\hat{P}_{\text{MinT}} = (S'\hat{W}^{-1}S)^{-1}S'\hat{W}^{-1}, \quad (10)$$

where \hat{W} is an estimate of W_c . A natural estimator could be $\hat{W} = T^{-1} \sum_{t=1}^T \hat{\mathbf{e}}_t(h)\hat{\mathbf{e}}_t(h)'$. However, since W_c is hard to estimate for $h > 1$, [17] assumed $W_c = k_h W_c^{(1)}$ where $k_h > 0$ and $W_c^{(1)}$ is the covariance matrix of the one-step-ahead base forecast errors, i.e. $h = 1$. Among the various covariance estimators for $W_c^{(1)}$ considered in [17], the most effective one is the shrinkage estimator with diagonal target given by

$$\hat{W} = (1 - \alpha)\hat{W}_s + \alpha\hat{W}_d, \quad \hat{W}_s = \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{e}}_t(1)\hat{\mathbf{e}}_t(1)', \quad (11)$$

where $\hat{W}_d = \text{diag}(\hat{W}_s)$ and $\alpha \in (0, 1]$. This covariance matrix is always invertible which is particularly useful for the computation of (9) when $n > T$ or $n = O(T)$ since then \hat{W}_s is singular, and hence non-invertible.

3 A NEW METHOD FOR HIERARCHICAL FORECASTING

3.1 Relaxation of the unbiasedness assumptions

We can decompose the objective function of MinT in (7) as

$$E[\|\mathbf{y}_{T+h} - \tilde{\mathbf{y}}_T(h)\|_2^2 | \mathcal{I}_T] \quad (12)$$

$$= \|SP(E[\hat{\mathbf{y}}_T(h)|\mathcal{I}_T] - E[\mathbf{y}_{T+h}|\mathcal{I}_T]) + (S - SPS)E[\mathbf{b}_{T+h}|\mathcal{I}_T]\|_2^2 \quad (13)$$

$$+ \text{Tr}(\text{Var}[\mathbf{y}_{T+h} - \tilde{\mathbf{y}}_T(h)|\mathcal{I}_T]), \quad (14)$$

where (13) and (14) are the bias and variance terms of the revised forecasts $\tilde{\mathbf{y}}_T(h)$, respectively.

Since MinT assumes the base forecasts are unbiased, the first term in (13) cancels. And since MinT requires the revised forecasts to be unbiased, equivalently $S = SPS$, the last term in (13) also cancels. As a result, minimizing (12) reduces to minimizing (14) with the constraint $S = SPS$. This is exactly the optimization problem given in (8).

We propose a new forecasting algorithm which do not impose the two unbiasedness conditions of MinT, **A1** and **C1**. In other words, we do not assume the base forecasts are unbiased, and we do not seek to produce the minimum variance unbiased revised forecasts. Instead, we aim to find the best tradeoff between bias and estimation variance of the revised forecasts by directly minimizing the mean squared forecast errors in (12).

Let $\hat{\mathbf{y}}_t(h)$ be the h -period ahead base forecasts based on $\mathcal{I}_t = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_t\}$ computed for $t = T_1, \dots, T - h$ where $T_1 < T$ is the number of observations used for model fitting. We consider the following empirical risk minimization (ERM) problem:

$$\min_{P \in \mathcal{P}} L_T(P), \quad (15)$$

where

$$L_T(P) = \frac{1}{(T - T_1 - h + 1)n} \sum_{t=T_1}^{T-h} \|\mathbf{y}_{t+h} - SP\hat{\mathbf{y}}_t(h)\|_2^2, \quad (16)$$

$$= \left\| Y - \hat{Y}P'S' \right\|_F^2 / Nn, \quad (17)$$

where $\|\cdot\|_F$ is the Frobenius norm defined as $\|X\|_F = \sqrt{\text{Tr}(X'X)}$, $N = T - T_1 - h + 1$, and

$$\begin{aligned} Y &= [\mathbf{y}_{T_1+h}, \mathbf{y}_{T_1+h+1}, \dots, \mathbf{y}_{T-h}]' \in \mathbb{R}^{N \times n}, \\ \hat{Y} &= [\hat{\mathbf{y}}_{T_1}(h), \hat{\mathbf{y}}_{T_1+1}(h), \dots, \hat{\mathbf{y}}_{T-h}(h)]' \in \mathbb{R}^{N \times n}. \end{aligned}$$

PROPOSITION 1. *If the matrix $\hat{Y}'\hat{Y}$ is invertible, the solution to (15) is given by*

$$\hat{P}_{\text{ERM}} = (S'S)^{-1}S'\hat{Y}'\hat{Y}(\hat{Y}'\hat{Y})^{-1} \quad (18)$$

$$= B'\hat{Y}(\hat{Y}'\hat{Y})^{-1}, \quad (19)$$

where $S'S$ is invertible by construction.

PROOF. Let $\text{vec}(X)$ denote the vectorization of the matrix X , i.e. all the columns of X are sequentially stacked in a one-column vector. Using the fact that $\|X\|_F^2 = \|\text{vec}(X)\|_2^2$, we can write (17) multiplied by Nn as

$$\left\| Y - \hat{Y}P'S' \right\|_F^2 = \left\| \text{vec}(Y) - \text{vec}(\hat{Y}P'S') \right\|_2^2 \quad (20)$$

$$= \left\| \text{vec}(Y) - (S \otimes \hat{Y})\text{vec}(P') \right\|_2^2, \quad (21)$$

The vector $\text{vec}(P')$ which minimizes (21) is given by

$$\text{vec}(P') = \left[(S \otimes \hat{Y})'(S \otimes \hat{Y}) \right]^{-1} (S \otimes \hat{Y})' \text{vec}(Y) \quad (22)$$

$$= \left[(S'S)^{-1}S' \otimes (\hat{Y}'\hat{Y})^{-1}\hat{Y}' \right] \text{vec}(Y) \quad (23)$$

$$= \text{vec} \left((\hat{Y}'\hat{Y})^{-1}\hat{Y}'YS(S'S)^{-1} \right). \quad (24)$$

Unvectorizing and transposing (24) gives expression (18). Finally, by noting that $Y = BS'$, (18) can be simplified into (19). \square

REMARK 2. When $\hat{Y}'\hat{Y}$ is not invertible, \hat{P}_{ERM} in (19) is not unique. This can notably happen with coherent base forecasts or in high-dimensional setting.

If the base forecasts are coherent, we can write $\hat{Y} = DS'$. In that case, the columns of \hat{Y} are linearly dependent. Consequently, we will have

$$\text{rank}(\hat{Y}'\hat{Y}) = \text{rank}(\hat{Y}) = \text{rank}(SD') \leq \min(\text{rank}(S), \text{rank}(D)) \leq m.$$

Therefore, $\hat{Y}'\hat{Y} \in \mathbb{R}^{n \times n}$ is not invertible. In other words, incoherent base forecasts is a necessary condition for $\hat{Y}'\hat{Y}$ to be invertible.

Obtaining coherent base forecasts is unlikely to happen in practice due to estimation errors in \hat{Y} . However, asymptotically, i.e. when $T_1 \rightarrow \infty$, and if the model is well-specified, the base forecasts $\hat{y}_t(h)$ will converge to the conditional mean given in expression (2). In that case, the rows of matrix D will contain $E[\mathbf{b}_{t+h}|\mathcal{I}_t]$ at times $t = T_1 + h, \dots, T - h$. This implies that the model is unidentified in the limit when the forecasts are coherent eventually as the number of observations T_1 tends to infinity.

Nevertheless, we can compute a unique solution using the Moore-Penrose generalized inverse. If $\hat{Y} = U_{\hat{Y}} \Sigma_{\hat{Y}}^{-1} V_{\hat{Y}}'$ is the thin SVD of \hat{Y} , i.e. without the zeroes in $\Sigma_{\hat{Y}}$, then this is equivalent to computing

$$\hat{P}_{ERM} = V_{\hat{Y}} \Sigma_{\hat{Y}}^{-1} U_{\hat{Y}}' B. \quad (25)$$

Note that the same observation applies to the MinT solution in (9). In fact, if the model is well-specified, W_c is asymptotically singular since $W_c = S V_b S'$ where V_b is the covariance of the random error terms in the bottom level.

3.2 Regularization for high-dimensional hierarchies

The ERM method has the advantage of relaxing the unbiasedness assumptions of MinT. The associated optimization problem in (15) involves the estimation of $m \times n = m^2 + mk$ parameters where m is the number of bottom-level series. However, common hierarchies can have a large number of bottom-level series m compared to the number of historical observations. The larger number of parameters to estimate results in an increase in estimation variance due to the accumulation of estimation errors. This is in turn transmitted into higher forecast variance, leading to poor forecast accuracy.

To overcome this issue, we propose to apply regularization to our objective in (17). More precisely, we propose a new method, called ERMreg, which solves the following problem:

$$\hat{P}_{ERMreg} = \underset{P \in \mathcal{P}}{\text{argmin}} \left\{ \left\| Y - \hat{Y} P' S' \right\|_F^2 / Nn + \lambda \|\text{vec}(P)\|_1 \right\}, \quad (26)$$

where $\lambda \geq 0$ is a hyperparameter which controls the amount of regularization to apply. When $\lambda = 0$, the problem reduces to (15) with a closed-form solution given by (19).

The major benefit of imposing a complexity penalty lies in the reduction in estimation variance by sacrificing some bias. In fact, parsimonious models often provide better forecasts than complex true models. Imposing an L_1 penalty is motivated by the fact that,

asymptotically, the optimal forecasts are given by $P = P_{BU} = [0_{m \times k} \mid I_m]$, which is sparse.

By replacing the first term in (26) with (21), the previous problem can be reduced to a standard LASSO problem with $\text{vec}(Y)$ as dependent variable and $(S \otimes \hat{Y})$ as design matrix. In other words, we need to solve a LASSO problem with $N \times n$ observations and $m^2 + mk$ variables.

We can use various efficient LASSO solvers for high-dimensional problems, for example cyclical coordinate descent methods (see Section 5.4 in [7] and [5]) and since $(S \otimes \hat{Y})$ is a sparse matrix, exploiting sparse matrix algebra can further reduce the computational load.

Adding an L_1 penalty term to the objective function will encourage sparsity in P . This is notably motivated by the fact that the optimal forecasts have a sparse P matrix. However, this will also have the effect of shrinking the revised forecasts towards zero, which is not desirable when dealing with strictly positive observations.

In traditional forecast combination, the combination weights are often penalized towards the average combination, i.e. the same weight is given to all forecasts [16]. In the context of hierarchical forecasting, average combination is not possible since the forecasts are associated to different variables and their aggregates.

We propose to penalize the matrix P towards the $P_{BU} = [0_{m \times k} \mid I_m]$, i.e. no combination is applied.

In other words, we propose to solve the following problem:

$$\hat{P}_{ERMregbu} = \underset{P \in \mathcal{P}}{\text{argmin}} \left\{ \left\| Y - \hat{Y} P' S' \right\|_F^2 / Nn + \lambda \|\text{vec}(P - P_{BU})\|_1 \right\}. \quad (27)$$

Note that we do not use P_{BU} as weight matrix, but we rather penalize our weights towards P_{BU} where the amount of regularization is controlled by λ .

Using P_{BU} as reference matrix is motivated by the fact that, asymptotically, the optimal forecasts given in (2) can be computed as revised forecasts with $P = P_{BU}$. Furthermore, the matrix P_{BU} has the advantage of parsimony. If $\hat{y}_T(h) = \mu_T(h)$, any matrix P that satisfy $PS = I$ will produce bottom-up revised forecasts. But, P_{BU} is the sparsest (in L_0 "norm") matrix P for which $PS = I$, as summarized in the following theorem.

THEOREM 3. $P_{BU} = [0_{m \times k} \mid I_m]$ is the unique solution of the following optimization problem:

$$\min_{P \in \mathcal{P}} \sum_{i=1}^m \sum_{j=1}^n |P_{ij}|^0 \text{ subject to } PS = I. \quad (28)$$

The proof is given in the Appendix B.1.

Using the reparametrization $C = P - P_{BU}$, Problem (27) can be rewritten in the same form as (26). Furthermore, Problem (27) can be written as a standard LASSO problem as stated by the next Proposition.

PROPOSITION 4. Let $\beta = \text{vec}(P' - P'_{BU})$. Problem (27) can be formulated as

$$\min_{\beta \in \mathbb{R}^p} \left\{ \|z - X\beta\|_2^2 / Nm + \lambda \|\beta\|_1 \right\}, \quad (29)$$

where $X = (S \otimes \hat{Y})$ and $z = \text{vec}(Y) - X \text{vec}(P'_{BU})$.

PROOF. The result follows from plugging $\beta = \hat{P}' - \hat{P}'_{\text{BU}}$ in (27) and applying vectorization as in (21). \square

REMARK 5. As stated by Proposition 4, Problem (27) can be reduced to a standard LASSO problem and therefore, one can use readily available standard algorithms embedded in statistical software. Furthermore, thanks to the L_1 regularization, a unique solution to (29) may exist even though $(X'X)$ is singular, where X is defined in (29).

However, standard inference tools based on well-established statistical properties of LASSO estimation are not directly applicable since it is subject to the issue of nearly singular design introduced in [12]. Using the method developed in [12], for the purpose of statistical inference, the following theorem provides the limiting distribution of appropriately standardized estimator of β in (29) even in the nearly singular design.

THEOREM 6. Under the assumptions **D1-D2** in Appendix A.2,

$$r_q(\hat{\beta} - \beta) \xrightarrow{d} \arg\min Q_0(g)$$

with

$$Q_0(g) = -2g'V + g'Dg + \lambda_0 \sum_{j=1}^{m \times n} \{g_j \text{sgn}(\beta_j) + |g_j|I(\beta_j = 0)\},$$

where $M = \text{plim}_{q \rightarrow \infty} 1/q \sum_{t=1}^q \mathbf{x}_t \mathbf{x}_t'$, $q = Nm$, V is defined as a zero-mean multivariate normal random vector such that $\text{Var}(g'V) = \sigma^2 g'Dg$ for g satisfying $Mg = \mathbf{0}$ and $\text{sgn}(\cdot)$ denotes the sign function.

The proof is given in Appendix B.2.

4 COMPARISON WITH MINT

In this section, we give more details about the differences between the MinT method and our new forecasting methods presented in Sections 2.2 and 3, respectively.

4.1 Unbiasedness assumptions

Under some regularity conditions (see Appendix A), the objective function in (16) is the sample version of the following population objective

$$\mathcal{L}(P) = E \left[\|\mathbf{y}_{T+h} - \tilde{\mathbf{y}}_T(h)\|_2^2 \right] \quad (30)$$

$$= E \left[E \left[\|\mathbf{y}_{T+h} - \tilde{\mathbf{y}}_T(h)\|_2^2 \mid \mathcal{I}_T \right] \right]. \quad (31)$$

Let us compare our population objective in (31) with the MinT population objective function in (7). First, we can see that we have relaxed the unbiasedness assumptions of MinT. Secondly, our objective involves *averaged* forecast errors while MinT deals with *path-dependent* forecast errors.

Although *path-dependent* forecast errors are more general than *averaged* forecast errors [14], it is hard to compute a reliable estimate of the *path-dependent* forecast errors. This is due to fact that we only observe one realization of the underlying data generating process (see Section 7.12 in [8] for a related discussion).

Furthermore, the MinT solution depends on the conditional covariance matrix W_c defined in (9). However, by the law of large numbers, \hat{W} in (11) asymptotically converges to the unconditional covariance matrix $W_u = E \left[\hat{\mathbf{e}}_T(h) \hat{\mathbf{e}}_T'(h) \right]$. In other words, MinT asymptotically minimize (31) under the constraints **A1** and **C1**. This is summarized in the following lemma.

LEMMA 7. Suppose that the conditions **A1** and **C1** hold. Then, the weight matrix that minimizes the criterion (31) is given by

$$P_U = \left(S' W_u^{-1} S \right)^{-1} S' W_u^{-1}, \quad (32)$$

where $W_u = E \left[\hat{\mathbf{e}}_T(h) \hat{\mathbf{e}}_T'(h) \right]$ and W_u^{-1} is the Moore-Penrose generalized inverse of W_u .

PROOF. Using expression (14) and Lemma 1 in [17], we can write

$$E \left[E \left[\|\mathbf{y}_{T+h} - \tilde{\mathbf{y}}_T(h)\|_2^2 \mid \mathcal{I}_T \right] \right] \quad (33)$$

$$= E \left[\text{Tr}(\text{Var}[\mathbf{y}_{T+h} - \tilde{\mathbf{y}}_T(h) \mid \mathcal{I}_T]) \right] \quad (34)$$

$$= \text{Tr} \left[E \left(\text{Var}[\mathbf{y}_{T+h} - \tilde{\mathbf{y}}_T(h) \mid \mathcal{I}_T] \right) \right] \quad (35)$$

$$= \text{Tr} \left[SPW_u P' S' \right], \quad (36)$$

where $W_u = E \left[\hat{\mathbf{e}}_T(h) \hat{\mathbf{e}}_T'(h) \right]$ with $\hat{\mathbf{e}}_T(h) = \mathbf{y}_{T+h} - \hat{\mathbf{y}}_T(h)$, i.e. W_u is the variance covariance matrix of the unconditional h -step-ahead base forecast errors. Then, the solution that minimizes (36) is given by

$$P_U = \left(S' W_u^{-1} S \right)^{-1} S' W_u^{-1}, \quad (37)$$

due to the Moore-Penrose generalized inverse. \square

A direct consequence of Lemma (7) is that, under the conditions **A1** and **C1**, \hat{P}_{MinT} in (10) and \hat{P}_{ERM} in (15) are asymptotically equivalent to P_U in (32). This is summarized in the following Corollary.

COROLLARY 8. Suppose that the regularity conditions **B1-B3** in Appendix A hold. Then, under the assumptions **A1** and **C1**, \hat{P}_{MinT} in (10) and \hat{P}_{ERM} in (15) are asymptotically equivalent to P_U in (32).

The proof is given in the Appendix B.3.

In summary, asymptotically, our ERM method includes MinT as a particular case, and reduces to MinT under the conditions **A1** and **C1**.

4.2 Regularization

Let $\tilde{\mathbf{b}}_T(h) = P^* \hat{\mathbf{y}}_T(h)$ be the bottom-level infeasible MinT forecasts where P^* is defined in (9). Using the proof of Theorem 1 in [17], we can show that $\tilde{\mathbf{b}}_T(h) = \beta$ where β is the solution of the following generalized least squares problem:

$$\min_{\beta \in \mathbb{R}^m} (\hat{\mathbf{y}}_T(h) - S\beta)' W^{-1} (\hat{\mathbf{y}}_T(h) - S\beta). \quad (38)$$

Although, the previous optimization problem does not include any explicit regularization term, the specific structure of the matrix S has an implicit regularization effect. For simplicity of exposition, let us assume that W is a block diagonal matrix, i.e. $W = \text{diag}(W_a, W_b)$. Then, since $S' = [A' \ I_m]'$ and $(\hat{\mathbf{a}}_T(h) \ \hat{\mathbf{b}}_T(h))'$, we can rewrite problem (38) as

$$\min_{\beta \in \mathbb{R}^m} (\hat{\mathbf{a}}_T(h) - A\beta)' W_a^{-1} (\hat{\mathbf{a}}_T(h) - A\beta) \quad (39)$$

$$+ (\beta - \hat{\mathbf{b}}_T(h))' W_b^{-1} (\beta - \hat{\mathbf{b}}_T(h)). \quad (40)$$

The previous optimization problem reveals that the objective function of MinT is similar to a generalized ridge regression problem, where the first term is a generalized weighted least squares criterion, and the second term is the generalized ridge penalty. A

closer look at (40) shows that MinT seeks to find the bottom-level forecasts β for which the associated bottom-up forecasts $A\beta$ are closest to the base forecasts at the aggregate levels where the importance of each difference is controlled by W_a^{-1} . Furthermore, the bottom-level forecasts β are shrunk towards the bottom-level base forecasts where the shrinkage amount is controlled by W_b^{-1} .

For MinT in high-dimensional setting, a shrinkage estimator becomes inevitable due to the accumulation of estimation errors. We can see this by noting that $|\text{Tr}(SP\hat{W}P'S') - \text{Tr}(SPW_cP'S')| \leq \|\hat{W} - W_c\|_\infty \|P'S'\|_1^2$. Related discussion can be found in [2] and [3] among many others.

Let us consider \hat{W} as defined in (11) with $\hat{W}_s = \text{diag}(\hat{W}_{s,a}, \hat{W}_{s,b})$ and $\hat{W}_d = \text{diag}(\hat{W}_{d,a}, \hat{W}_{d,b})$. Let us plug \hat{W} in (40), i.e. $W_a = (1 - \alpha)\hat{W}_{s,a} + \alpha\hat{W}_{d,a}$ and $W_b = (1 - \alpha)\hat{W}_{s,b} + \alpha\hat{W}_{d,b}$. Then, we can see that the value of α will notably affect the amount of shrinkage towards the bottom-level base forecasts. In particular, when $\alpha = 1$, which has been considered in [11], the penalty term reduces to $\sum_{j=1}^m \frac{1}{\hat{w}_j} (\beta_j - \hat{b}_{j,T}(h))^2$ where \hat{w}_j is the j th diagonal element of the matrix $\hat{W}_{d,b}$. Assuming the series are in the same scale, this shows that more shrinkage is applied when the base forecasts $\hat{b}_{j,T}(h)$ are more accurate.

Although MinT implicitly regularizes the revised forecasts, the solution is still based on the assumption **A1** and constraint **C1**. In other words, if **A1** and **C1** are not appropriate for the problem at hand, shrinking the revised forecasts towards the base forecasts can lead to poor forecast accuracy. Our method ERMreg is more robust since we try to find the best tradeoff between bias and estimation variance by explicitly adding a regularization term.

5 EXPERIMENTS

We perform multiple experiments to compare our forecasting methods with the state-of-the-art methods in different conditions. We consider both simulated and real-world datasets. For each dataset, we sample 100 different hierarchies, and the results are averaged over these 100 hierarchies.

5.1 Data

In the following, we describe the different data sets used in the experiments.

5.1.1 Simulated data (see Figure 3a). We use the data generating mechanism developed in [17] (see Sections 3.3 and 3.5) to reflect common characteristics of hierarchical time series such as correlations across series, and smoother data with more aggregation. The bottom-level series are sampled from correlated ARMA processes where the coefficients are sampled uniformly from a specific parameter space that guarantees stationarity. The random error terms have a multivariate Gaussian distribution with a covariance structure that induces a strongly positive correlation among series with the same parents, but a moderately positive correlation among series with different parents. For each series, we generate 600 observations split in training and validation sets, 2/3 and 1/3 respectively. An additional sample of 200 observations is used for testing. We consider both a small and a large two-level hierarchy where $m = 4$ and $k = 3$, and $m = 100$ and $k = 26$, respectively. More specifically,

the bottom level series were aggregated in groups of two for the next level in the small hierarchy, and groups of four in the large hierarchy. These series were then aggregated to obtain the top level series.

5.1.2 Road traffic¹ (see Figure 3b). The dataset gives the occupancy rate (between 0 and 1) of 963 car lanes of San Francisco bay area freeways. The measurements are sampled every 10 minutes from Jan. 1st 2008 to Mar. 30th 2009. This dataset has been notably used in [4] for classification tasks. We aggregate the data to 366 daily observations split in 120, 120 and 126 observations for training, validation and testing, respectively. We consider hierarchies with $m = 200$ and $k = 7$, where each hierarchy is constructed as follows. We sample 200 bottom level series from the 963 series, and compute the upper level series by aggregation. More specifically, 200 series at the bottom level were aggregated in groups of 50 for the next level, resulting in 4 series. These 4 series were then aggregated in groups of two to obtain two aggregate series and the top level series.

5.1.3 Wikipedia webpage views² (see Figure 3c). The dataset gives the number of daily views of 145,000 different Wikipedia articles starting from July, 1st, 2015 up until December 31st, 2016. We use the 366 observations for 2016 split in 86, 160 and 120 observations for training, validation and testing, respectively. We consider hierarchies where $m = 150$ and $k = 49$, where each hierarchy is constructed as follows. We sample 150 bottom level series from the 145k series, and compute the upper level series by aggregation. More specifically, the bottom level series are hierarchically aggregated into type of agent (“all-agents” and “spider”), type of access (“all-access”, “desktop” and “mobile-web”) and country codes (“en”, “fr”, “de”, “ja”, “ru”, “zh”). The hierarchy’s architecture is summarized in Figure 2, where the top level series is in the center and the most disaggregated bottom level series are on the circumference of the circle.

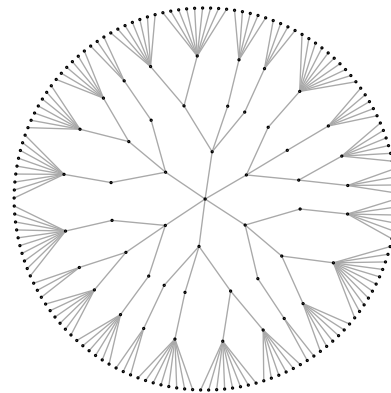


Figure 2: The hierarchical structure of the Wikipedia data.

¹<https://archive.ics.uci.edu/ml/datasets/PEMS-SF>

²<https://www.kaggle.com/c/web-traffic-time-series-forecasting/data>

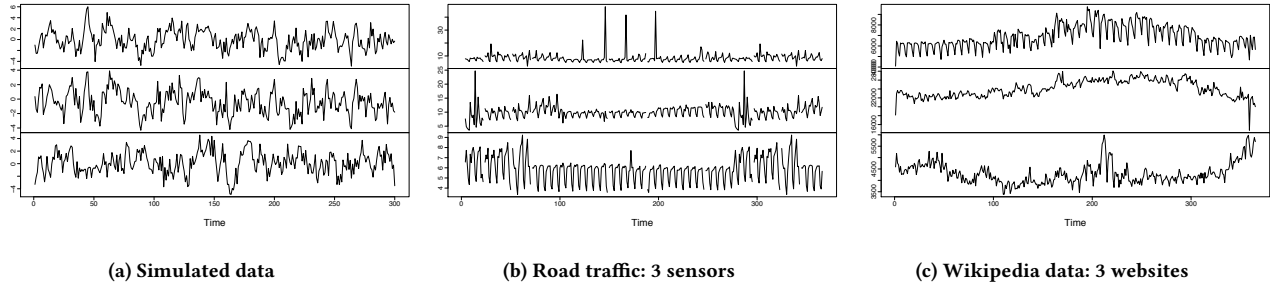


Figure 3: Visualization of the various time series considered in the experiments.

5.2 Methodology

In order to ensure stationarity, we remove the trend and seasonal component using the Seasonal and Trend decomposition using Loess (STL). We also apply a log transformation to stabilize the variance. Finally, we replace any outliers using linear interpolation, as implemented in the `tsclean` function in the forecast package for R.

For each series in each hierarchy, we produce one-step-ahead rolling-origin base forecasts for both the validation and test sets from a model fitted to the last T_1 observations. We use $T_1 = 400, 86, 120$ for the simulated, road traffic and Wikipedia data sets, respectively. We consider both ARIMA and Exponential Smoothing (ETS) forecasting algorithms fitted using the automatic model selection algorithm of [9], which is implemented in the forecast package for R. While the simulated data are generated using an ARMA process, using an ETS model allows us to simulate what happens in practice, where the true data generating processes is unknown and misspecified.

We consider the following hierarchical forecasting methods:

- BASE.** The base forecasts defined in (4) produced as described in the previous paragraph. Recall that the base forecasts are (probably) incoherent.
- BU.** The bottom-up forecasts defined in (3), i.e. where (5) is computed using (6).
- MinTsam.** The MinT forecasts given in (10) where (11) is computed with $\alpha = 1$.
- MinTols.** The MinT forecasts given in (10) where (11) is computed with $\hat{W}_d = \hat{\sigma}^2 I_n$ and $\alpha = 0$.
- MinTshr.** The MinT forecasts given in (10) where (11) is with an optimal shrinkage factor α , as described in [15].
- ERM.** The ERM forecasts where (5) is computed using (25).
- ERMreg.** The ERMreg forecasts where (5) is computed using (26). We use cross-validation to select the best value of λ , as implemented in the `glmnet` package in R.
- ERMregbu.** The ERMreg forecasts where (5) is computed using (27).

6 RESULTS

For each data set, we compute the one-step-ahead mean squared forecast errors for all series in the hierarchy, i.e.

$$\frac{1}{T_{\text{test}}} \sum_{t=T}^{T+T_{\text{test}}-1} \|\mathbf{y}_{t+1} - SP\hat{\mathbf{y}}_t(1)\|_2^2,$$

and the bottom series only, i.e.

$$\frac{1}{T_{\text{test}}} \sum_{t=T}^{T+T_{\text{test}}} \|\mathbf{b}_{t+1} - P\hat{\mathbf{y}}_t(1)\|_2^2.$$

where T_{test} is the number of observations in the test set. We have $T_{\text{test}} = 200, 120, 120$ for the simulated, road traffic and Wikipedia data, respectively. We then average the results for the 100 hierarchies.

We start by analyzing the results for the simulated data. Table 1 and 2 give the forecasting errors for the small and the large hierarchy, respectively.

	ARIMA		ETS	
	All	Bottom	All	Bottom
BASE	39.32 (0.33)	8.20 (0.05)	43.64 (0.42)	9.09 (0.07)
BU	39.11 (0.32)	8.20 (0.05)	43.59 (0.42)	9.09 (0.07)
ERM	40.23 (0.41)	8.43 (0.07)	41.26 (0.43)	8.59 (0.07)
ERMreg	39.45 (0.33)	8.29 (0.05)	40.83 (0.40)	8.53 (0.07)
ERMregbu	39.29 (0.33)	8.24 (0.05)	40.76 (0.40)	8.51 (0.07)
MINTsam	40.97 (0.48)	8.55 (0.08)	44.90 (0.71)	9.38 (0.14)
MINTols	39.14 (0.33)	8.20 (0.05)	43.63 (0.42)	9.10 (0.07)
MINTshr	39.15 (0.32)	8.20 (0.05)	43.61 (0.42)	9.09 (0.07)

Table 1: Forecast errors for the simulated data with a small hierarchy. Standard errors are given in parentheses.

Table 1 and 2 show that when using an ETS model, i.e. when the base model is misspecified, our ERMreg methods outperform all the other methods. This confirms the robustness of our methods to misspecification of the base forecasts. When an ARMA model is

	ARIMA		ETS	
	All	Bottom	All	Bottom
BASE	11.03 (0.05)	2.05 (0.00)	12.21 (0.06)	2.28 (0.00)
BU	10.87 (0.05)	2.05 (0.00)	12.16 (0.06)	2.28 (0.00)
ERM	18.41 (0.13)	3.53 (0.02)	22.01 (0.19)	4.18 (0.02)
ERMreg	11.69 (0.05)	2.28 (0.00)	12.37 (0.06)	2.37 (0.00)
ERMregbu	10.88 (0.05)	2.06 (0.00)	11.69 (0.05)	2.22 (0.00)
MINTsam	12.13 (0.10)	2.30 (0.01)	13.04 (0.07)	2.44 (0.01)
MINTols	11.00 (0.05)	2.06 (0.00)	12.21 (0.06)	2.28 (0.00)
MINTshr	10.85 (0.05)	2.05 (0.00)	12.16 (0.06)	2.28 (0.00)

Table 2: Forecast errors ($\times 10^2$) for the simulated data with a large hierarchy. Standard errors are given in parentheses.

used, i.e. the base model is well-specified, our methods still provide good forecasts although this is the ideal scenario for MINT.

In Table 2, for the large hierarchy, we can see that MINTsam and ERM have larger forecast errors than the other methods. This is not surprising since both MINTsam and ERM do not use regularization, and as a result suffer from the accumulation of errors in high dimensions. In fact, the larger number of parameters to estimate significantly increases the forecast variance which leads to poor forecast accuracy. We can also see that ERM is more affected by high dimensionality than MINTsam. With MINTsam, the accumulation of errors will indirectly affect the estimation of the combination weights. In fact, in order to compute the combination weights, MINTsam plugs \hat{W}_s given in (11) into (10). However, since ERM solves the high-dimensional regression problem given in (15), it will be directly affected by the accumulation of errors. Nevertheless, we can see that MINTshr and ERMreg significantly reduce the forecast errors thanks to regularization which reduces forecast variance. The robustness of our ERMreg methods to misspecification can also be seen in Table 2.

Table 3 and 4 give the results for the road traffic and Wikipedia data, respectively. For clarity, given that MINTsam and ERM provide very poor forecasts (due to high-dimensionality, as discussed above), we omit their results from the tables. For both data sets, our method ERMregbu provides competitive results compared to the other methods.

Figure 4 plots the weight matrix of different methods computed for one hierarchy of the road traffic data set. Each entry of the matrix gives $|P_{ij}|$, i.e. the absolute value of the i th row and the j th column of matrix P . A white cell has zero value, and a darker cell has a higher value. In addition to differences in forecast accuracy, Figure 4 shows that the different methods provide combination weights with different structures. In fact, we can see that MINTshr and ERM have a dense weight matrix while ERMreg and ERMregbu provide sparse matrices, as expected. The diagonal line for some matrices shows that the base forecast for the series under consideration has been selected in the forecast combination.

	ARIMA		ETS	
	All	Bottom	All	Bottom
BASE	34.47 (0.35)	1.72 (0.01)	36.97 (0.36)	1.82 (0.01)
BU	32.85 (0.31)	1.72 (0.01)	32.95 (0.31)	1.82 (0.01)
ERMreg	31.37 (0.53)	1.78 (0.01)	36.78 (0.78)	1.81 (0.02)
ERMregbu	31.24 (0.32)	1.72 (0.01)	36.68 (0.61)	1.88 (0.01)
MINTols	34.18 (0.34)	1.73 (0.01)	36.66 (0.36)	1.84 (0.01)
MINTshr	32.67 (0.31)	1.72 (0.01)	33.10 (0.31)	1.81 (0.01)

Table 3: Forecast errors for the road traffic data. Standard errors are given in parentheses.

	ARIMA		ETS	
	All	Bottom	All	Bottom
BASE	70.98 (0.68)	10.94 (0.09)	74.01 (0.71)	11.10 (0.10)
BU	69.71 (0.65)	10.94 (0.09)	72.43 (0.63)	11.10 (0.10)
ERMreg	72.22 (0.69)	11.19 (0.10)	73.87 (0.73)	11.21 (0.10)
ERMregbu	69.61 (0.64)	11.01 (0.09)	72.95 (0.65)	11.13 (0.10)
MINTols	70.29 (0.67)	10.96 (0.09)	73.69 (0.70)	11.11 (0.10)
MINTshr	69.31 (0.66)	10.93 (0.09)	72.78 (0.65)	11.10 (0.10)

Table 4: Forecast errors for the Wikipedia data. Standard errors are given in parentheses.

7 CONCLUSION

We considered the problem of forecasting multiple time series under hierarchical constraints. We proposed a new forecasting method which does not assume or constrain the forecasts to be unbiased but aim to find the best tradeoff between bias and forecast variance. Our method involves solving a multi-response regression problem for which we provided a closed-form solution. To deal with high-dimensional hierarchies, we also proposed a regularization method that has a theoretical justification, and proved its asymptotic properties. We also showed that our problems can be divided into multiple independent problems which can be solved efficiently in parallel. We compared our methods both theoretically and empirically with the state-of-the-art methods for hierarchical forecasting. Using both simulated and two real-world time series data, we showed our methods is competitive with the state-of-the-art methods in various conditions.

One direction for future work is the improvement of the computational complexity of our methods, as well as the development of new regularization methods which exploit the hierarchical structure of the data. Another important problem is the computation of probabilistic forecasts for high-dimensional hierarchies [1]. Finally, it is also possible to extend and adapt our methods to the problem of multi-source learning/forecasting with hierarchical structures [18–20].

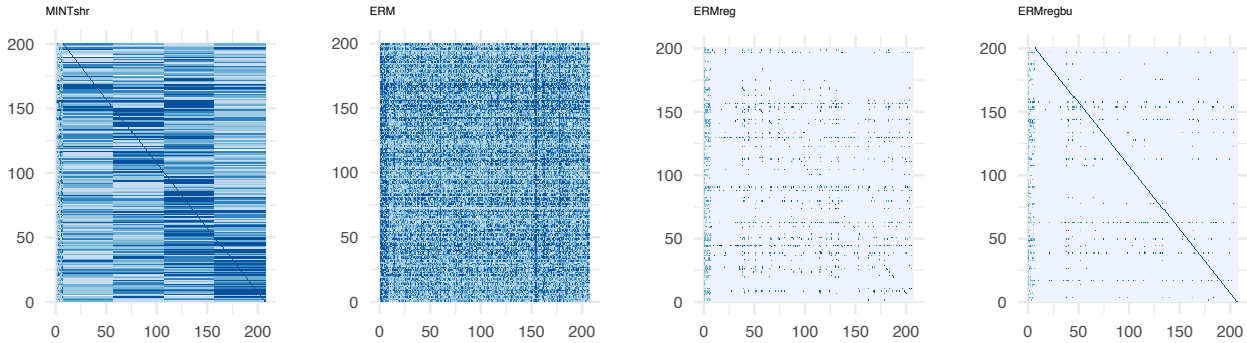


Figure 4: Matrix combination weights for different hierarchical forecasting methods.

8 ACKNOWLEDGEMENTS

This work is partially supported by the Australian Research Council under Grant No. DP150104292 and Grant No. DE170100713. We thank anonymous reviewers and Shanika Wickramasuriya for their insightful and constructive comments.

REFERENCES

- [1] Ben Taieb, James Taylor, and Rob Hyndman. 2017. Coherent Probabilistic Forecasts for Hierarchical Time Series. In *Proceedings of The 34th International Conference on Machine Learning (Proceedings of Machine Learning Research)*, Vol. 70. 3348–3357.
- [2] Peter J Bickel and Elizaveta Levina. 2008. Regularized Estimation of Large Covariance Matrices. *Annals of statistics* 36, 1 (2008), 199–227.
- [3] Peter J Bickel, Elizaveta Levina, and Others. 2008. Covariance regularization by thresholding. *Annals of statistics* 36, 6 (2008), 2577–2604.
- [4] Marco Cuturi. 2011. Fast global alignment kernels. In *Proceedings of the 28th international conference on machine learning (ICML-11)*. 929–936.
- [5] Jerome Friedman, Trevor Hastie, Holger Höfling, and Robert Tibshirani. 2007. Pathwise coordinate optimization. *The Annals of Applied Statistics* 1, 2 (Dec. 2007), 302–332.
- [6] Tilmann Gneiting. 2011. Making and evaluating point forecasts. *J. Amer. Statist. Assoc.* 106, 494 (June 2011), 746–762.
- [7] T Hastie, R Tibshirani, and M Wainwright. 2015. *Statistical Learning with Sparsity: The Lasso and Generalizations*. CRC Press.
- [8] Trevor J Hastie, Robert Tibshirani, and Jerome H Friedman. 2008. *The elements of statistical learning*. Vol. 18. Springer-Verlag, 764 pages.
- [9] Rob Hyndman and Yeasmin Khandakar. 2008. Automatic Time Series Forecasting: The forecast Package for R. *Journal of Statistical Software* 27, 1 (2008), 1–22.
- [10] Rob J Hyndman, Roman A Ahmed, George Athanasopoulos, and Han Lin Shang. 2011. Optimal combination forecasts for hierarchical time series. *Computational statistics & data analysis* 55, 9 (Sept. 2011), 2579–2589.
- [11] Rob J Hyndman, Alan J Lee, and Earo Wang. 2016. Fast computation of reconciled forecasts for hierarchical and grouped time series. *Computational Statistics & Data Analysis* 97 (May 2016), 16–32.
- [12] Keith Knight. 2008. Shrinkage estimation for nearly singular designs. *Econometric Theory* 24, 2 (April 2008), 323–337.
- [13] Mirko Kremer, Enno Siemsen, and Douglas J Thomas. 2016. The Sum and Its Parts: Judgmental Hierarchical Forecasting. *Management Science* 62, 9 (2016), 2745–2764.
- [14] Vitaly Kuznetsov and Mehryar Mohri. 2015. Learning Theory and Algorithms for Forecasting Non-stationary Time Series. In *Advances in Neural Information Processing Systems* 28, C Cortes, N D Lawrence, D D Lee, M Sugiyama, R Garnett, and R Garnett (Eds.). Curran Associates, Inc., 541–549.
- [15] Juliane Schäfer and Korbinian Strimmer. 2005. A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics. *Statistical Applications in Genetics and Molecular Biology* 4 (Nov. 2005), 32.
- [16] A Timmermann. 2006. Forecast combinations. In *Handbook of Economic Forecasting*. Vol. 1. Elsevier, 135–196.
- [17] Shanika L Wickramasuriya, George Athanasopoulos, and Rob J Hyndman. 2018. Optimal forecast reconciliation for hierarchical and grouped time series through trace minimization. *J. Amer. Statist. Assoc.* (March 2018), 1–45.
- [18] Pei Yang and Jingrui He. 2015. Model Multiple Heterogeneity via Hierarchical Multi-Latent Space Learning. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 1375–1384.
- [19] Chenwei Zhang, Sihong Xie, Yaliang Li, Jing Gao, Wei Fan, and Philip S Yu. 2016. Multi-source Hierarchical Prediction Consolidation. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*. ACM, 2251–2256.
- [20] Liang Zhao, Jieping Ye, Feng Chen, Chang-Tien Lu, and Naren Ramakrishnan. 2016. Hierarchical Incomplete Multi-source Feature Learning for Spatiotemporal Event Forecasting. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2085–2094.

A REGULARITY CONDITIONS

A.1 Definitions

In what follows, $O_p(\cdot)$ and $o_p(\cdot)$ denote the usual big O and little o in probability. \xrightarrow{p} and \xrightarrow{d} denote convergence in probability and convergence in distribution respectively.

DEFINITION 9. $U_T(s) \xrightarrow{e-d} U(s)$ if for any closed rectangles R_1, \dots, R_k with open interiors R_1^0, \dots, R_k^0 and any real numbers a_1, \dots, a_k ,

$$\begin{aligned} & P\{ \inf_{s \in R_1} U(s) > a_1, \dots, \inf_{s \in R_k} U(s) > a_k \} \\ & \leq \liminf_{T \rightarrow \infty} P\{ \inf_{s \in R_1} U(s) > a_1, \dots, \inf_{s \in R_k} U(s) > a_k \} \\ & \leq \limsup_{T \rightarrow \infty} P\{ \inf_{s \in R_1^0} U(s) > a_1, \dots, \inf_{s \in R_k^0} U(s) > a_k \} \\ & \leq P\{ \inf_{s \in R_1^0} U(s) > a_1, \dots, \inf_{s \in R_k^0} U(s) > a_k \} \end{aligned}$$

A.2 Regularity conditions for Theorem 6

D1 (near singularity)

- (i) Let $M_q = 1/q \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t'$. M_q is nonsingular for all q whereas its probability limit, M is singular.
- (ii) As $q \rightarrow \infty$, \mathbf{x}_t satisfies

$$\frac{1}{q} \max_{1 \leq t \leq q} \mathbf{x}_t M_q^{-1} \mathbf{x}_t \rightarrow 0.$$

- (iii) For some sequence of constants, $\{c_q\}$ tending to infinity,

$$c_q(M_q - M) \xrightarrow{p} D$$

where D is a positive definite matrix on the null space of M .

- D2 (Shrinkage parameter and the rate of convergence)** Let the shrinkage parameter λ be such that $\lambda/r_q \rightarrow \lambda_0$ where λ_0 is a positive constant and $r_q = \sqrt{q/c_q}$.

A.3 Regularity conditions for Corollary 8

To ensure that the sample objective function (15) converges to the limit objective function (31) and the minimizer of the former matches with that of the latter, we assume the following.

- B1 (Data)** $\{\mathbf{y}_t\}_{t=1}^T$ is strictly stationary, ergodic with exponentially decaying mixing coefficients.

- B2 (Parameter space)** The true parameter P_0 involving the Moore-Penrose generalized inverse belongs to the interior of a compact convex subset \mathcal{P} of the Euclidean space, $\mathbb{R}^{m \times n}$.

- B3 (Objective function)** The nonrandom limit objective $\mathcal{L}(P)$ exists, is minimized uniquely at P_0 involving the Moore-Penrose generalized inverse and is equicontinuous on \mathcal{P} , i.e. it is well behaved in the neighborhood around P_0 , i.e. $\|P - P_0\| \leq \varepsilon$ at some ε . For any $\eta > 0$, there exists $\varepsilon > 0$,

$$\liminf_{T \rightarrow \infty} P \left(\sup_{\|P - P_0\| \geq \eta} (L_T(P) - L_T(P_0)) \leq -\varepsilon \right) = 1,$$

and for any $\eta > 0$, there exists $\varepsilon > 0$, s.t.

$$\inf_{\|P - P_0\| > \eta} \mathcal{L}(P) - \mathcal{L}(P_0) \geq \varepsilon.$$

B PROOFS

B.1 Proof of Theorem 3

If condition. First, we show that P_{BU} is at least one of solutions of the objective function in (28). Note that $P_{\text{BU}}S = I$ and

$$\sum_{i=1}^m \sum_{j=1}^n |P_{\text{BU}[ij]}|^0 = m.$$

Therefore, in order to show that P_{BU} could be a solution, we must show that

$$\min \left(\sum_{i=1}^m \sum_{j=1}^n |P_{ij}|^0 \right) = m$$

for any P such that $PS = I$.

Note that for P such that $PS = I$, $\text{rank}(PS) = \text{rank}(I) = m \leq \min(\text{rank}(P), \text{rank}(S))$, and $\text{rank}(S) = m$. So, $\text{rank}(P) \geq m$. Since $\text{rank}(I) = m$, $\text{rank}(P) \geq m$. This implies that the L_0 norm of P must be greater or equal to m . This is because the rank of P could not be greater than or equal to m if the number of non-zero values of P is lower than m . As a result, P_{BU} is at least one of solutions of the above objective function.

Only if condition. Now, we turn our attention to showing that P_{BU} is the unique solution for (28). Suppose that P_{BU} is only one of many solutions for the above objective function. Then, there exists P^* such that $P^* = P_{\text{BU}} + C$ with a non-zero $m \times n$ sized matrix, C and $P^*S = I$. Since $\min(\sum_{i=1}^m \sum_{j=1}^n |P_{ij}^*|^0) = m$, there are two cases we need to consider.

The first case is when one of the non-zero values of P^* is different from the non-zero values of P_{BU} in the same location. This is contradictory since there is no such solution because in that case $P^*S \neq I$, i.e. it should be $P_{\text{BU}} = P^*$.

The second case is when one of the non-zero values of P_{BU} becomes zero and one of zero values of P_{BU} becomes non-zero. This should be satisfied to ensure that $\min(\sum_{i=1}^m \sum_{j=1}^n |P_{ij}^*|^0) = m$. Note that

$$\begin{aligned} P_{\text{BU}}S + CS = I &\implies CS = \mathbf{0} \implies \\ C_a A + C_b = \mathbf{0} &\implies C_b = -C_a A \end{aligned}$$

Then, one of the diagonal values of C_b should be -1 and all the others are zeros while one of the values of C_a should be non-zero and others are zeros.

Note that $C_{b,ii} = -\sum_{j=1}^n C_{a,ij} A_{ji}$ for $j = 1, \dots, n$. Suppose that $i = k$ for some fixed k . Then, $C_{b,kk} = -1$ and only one non-zero value belongs to $C_{a,kj}$ and it should be 1. However, if any of A_{jk} contains at least one non-zero value except for $j = k$, then the equality will not hold. But this is the case because A is the aggregation matrix. This is contradictory since $C_b = -C_a A$. Therefore, $C = \mathbf{0}$.

B.2 Proof of Theorem 6

Recall that for $\mathbf{y}_t = (\mathbf{a}_t \ \mathbf{b}_t)'$, $\mathbf{y}_t = S\mathbf{b}_t$ where $S = \begin{bmatrix} A' & I_m \end{bmatrix}'$. That is, \mathbf{a}_t can be expressed as linear combinations of \mathbf{b}_t . This implies that our regression problem is subject to the issue of nearly singular design. As seen from Remark 2, this is the natural consequence of hierarchical time series forecasting given that $Y'Y$ is singular whereas $\hat{Y}'\hat{Y}$ can be nonsingular due to estimation error.

Recall that

$$Q_q(\boldsymbol{\beta}) = \|\mathbf{z} - X\boldsymbol{\beta}\|_2^2/q + \lambda\|\boldsymbol{\beta}\|_1, \quad (41)$$

where $X = (S \otimes \hat{Y})$, $\mathbf{z} = \text{vec}(Y) - X \text{vec}(\hat{P}'_{BU})$, $\boldsymbol{\beta} = \text{vec}(\hat{P}' - \hat{P}'_{BU})$, $N = T - T_1 - h + 1$, and $q = Nm$.

It is worth noting that $Q_q(\cdot)$ is convex by construction. Here we borrow the concept of near collinearity or near singularity introduced in [12], i.e. for a regressor matrix X , M_q is nonsingular but its limit, M is singular where

$$M_q = \frac{1}{q} \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t'$$

$$M = \text{plim}_{q \rightarrow \infty} M_q$$

Define $\mathbf{g} = r_q(\boldsymbol{\beta} - \boldsymbol{\beta}_0)$. Note that via reparameterization, (41) can be rewritten as

$$Q_q(\mathbf{g}) = (\|\boldsymbol{\varepsilon} - X\mathbf{g}/r_q\|_2^2 - \|\boldsymbol{\varepsilon}\|_2^2)/q + \lambda(\|\boldsymbol{\beta} + \mathbf{g}/r_q\|_1 - \|\boldsymbol{\beta}\|_1) \quad (42)$$

where $\boldsymbol{\varepsilon} = \text{vec}(\mathbf{z} - X\boldsymbol{\beta})$ and λ and r_q satisfies the assumption D2. Note that if $\hat{\boldsymbol{\beta}}$ minimizes $Q_q(\boldsymbol{\beta})$, $\hat{\mathbf{g}} := r_q(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$ minimizes $Q_q(\mathbf{g})$.

Let $Q_0(\mathbf{g}) = \lim_{q \rightarrow \infty} Q_q(\mathbf{g})$. Due to convexity of $Q_q(\cdot)$, $Q_0(\cdot)$ and $Q_q(\cdot)$ is stochastically equi-lower semicontinuous, which implies that $Q_q(\mathbf{g}) \xrightarrow{e-d} Q_0(\mathbf{g})$ (see [12]). Trivially, $r_q(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) = O_p(1)$ due to the convexity of $Q_q(\mathbf{g})$ and the usual asymptotic theory related to the LASSO estimation.

Combined with $\hat{\mathbf{g}} = O_p(1)$, the epi-convergence $Q_q(\mathbf{g}) \xrightarrow{e-d} Q_0(\mathbf{g})$, implies that the limiting distribution of $\hat{\mathbf{g}}$ is equivalent to that of the minimizer of the limit objective function of the sample objective function, i.e. $\hat{\mathbf{g}} \xrightarrow{d} \text{argmin } Q_0(\mathbf{g})$.

The remainder is to obtain the finite-dimensional limit of $Q_q(\mathbf{g})$. Note that

$$(\|\boldsymbol{\varepsilon} - X\mathbf{g}/c_q\|_2^2 - \|\boldsymbol{\varepsilon}\|_2^2)/q \xrightarrow{d} -2\mathbf{g}'\mathbf{V} + \mathbf{g}'\mathbf{D}\mathbf{g}$$

$$\lambda(\|\boldsymbol{\beta} + \mathbf{g}/r_q\|_1 - \|\boldsymbol{\beta}\|_1) \xrightarrow{d} \lambda_0 \sum_{j=1}^{m \times n} \{g_j \text{sgn}(\beta_j) + |g_j|I(\beta_j = 0)\}$$

where \mathbf{V} is defined as a zero-mean multivariate normal random vector such that $\text{Var}(\mathbf{g}'\mathbf{V}) = \sigma^2 \mathbf{g}'\mathbf{D}\mathbf{g}$ for \mathbf{g} satisfying $M\mathbf{g} = \mathbf{0}$. This completes the proof.

B.3 Proof of Corollary 8

We show that \hat{P}_{ERM} computed by solving (15) and \hat{P}_{MinT} are asymptotically equivalent to P_U where P_U is defined in (37).

Consider \hat{W} defined in (11), and assume $\alpha \rightarrow 0$ when $T \rightarrow \infty$. Then, by the usual laws of large numbers for the stationary and ergodic data, $\hat{W} \xrightarrow{p} W_u$. Furthermore, using Slutsky theorem, we have

$$\hat{P}_{MinT} \xrightarrow{p} P_U.$$

In order to show that $\hat{P}_{ERM} \xrightarrow{p} P_U$, we note that

$$\sup_{P \in \mathcal{P}} |L_T(P, \mathbf{y}) - \mathcal{L}(P)| \leq \sup_{P \in \mathcal{P}} |L_T(P, \mathbf{y}) - \text{EL}_T(P, \mathbf{y})|$$

$$+ \sup_{P \in \mathcal{P}} |\text{EL}_T(P, \mathbf{y}) - \mathcal{L}(P)|$$

By construction, $L_T(P, \mathbf{y})$ is continuous in $P \in \mathcal{P}$ for all \mathbf{y} and is a measurable function of \mathbf{y} for all $P \in \mathcal{P}$. Due to the assumptions B1-B3, $L_T(P, \mathbf{y})$ is stochastic equicontinuous such that

$$\sup_{P \in \mathcal{P}} \sup_{P^* \in B(P, \epsilon)} |L_T(P, \mathbf{y}) - L_T(P^*, \mathbf{y})| = o_p(1)$$

where $B(P, \epsilon)$ is a closed ball in \mathcal{P} of radius $\epsilon \geq 0$ centered at P . Due to the pointwise convergence of $L_T(P, \mathbf{y})$ to $\mathcal{L}(P)$ on P and uniform stochastic equicontinuity of $L_T(P, \mathbf{y})$,

$$\sup_{P \in \mathcal{P}} |L_T(P, \mathbf{y}) - \mathcal{L}(P)| = o_p(1).$$

Since we define P_0 as the true parameter involving the Moore-Penrose generalized inverse in B2, P_U is the unique one for P_0 in a sense that P_U minimizes $\mathcal{L}(P)$ uniquely in the mean squared error sense. Note also that \hat{P}_{ERM} is the unique minimizer of $L_T(P, \mathbf{y})$ involving the Moore-Penrose generalized inverse and shares the same criterion as P_U . Therefore, due to the assumption B3 and the argmax theorem, $\hat{P}_{ERM} \xrightarrow{p} P_U$.