



Department of Econometrics and Business Statistics

<http://monash.edu/business/ebs/research/publications>

Forecast Linear Augmented Projection (FLAP): A free lunch to reduce forecast error variance

Yangzhuoran Fin Yang, George Athanasopoulos, Rob
J. Hyndman, Anastasios Panagiotelis

April 2024

Working Paper ??/??



AACSB
ACCREDITED



Forecast Linear Augmented Projection (FLAP): A free lunch to reduce forecast error variance

Yangzhuoran Fin Yang

Monash University
Melbourne, Australia
Email: Fin.Yang@monash.edu

George Athanasopoulos

Monash University
Melbourne, Australia

Rob J. Hyndman

Monash University
Melbourne, Australia

Anastasios Panagiotelis

University of Sydney
Sydney, Australia

21 April 2024

Forecast Linear Augmented Projection (FLAP): A free lunch to reduce forecast error variance

Abstract

A novel forecast linear augmented projection (FLAP) method is introduced. FLAP provably reduces the forecast error variance of any unbiased multivariate forecast without introducing bias. The method first constructs new series as linear combinations of the original series. Forecasts are then generated for both the original and new series. Finally, the full vector of forecasts is projected onto a linear subspace where the constraints implied by the combination weights hold. It is proven that the trace of the forecast error variance is non-increasing with the number of components, and mild conditions are established for which it is strictly decreasing. It is also shown that the proposed method achieves maximum forecast error variance reduction among linear projections. The theoretical results are validated through simulations and two empirical applications based on Australian tourism and FRED-MD data. Notably, using FLAP with Principal Component Analysis to construct the new series leads to substantial forecast error variance reduction.

1 Introduction

Multivariate forecasting arises in a number of disciplines including macroeconomics and finance; see Carriero, Galvao & Kapetanios (2019) and Tsay (2013) respectively, and references therein. We introduce a new post processing framework that (i) augments the data by constructing new series that are linear combinations of the original series, (ii) forecasts both the original and new series and (iii) recovers a new set of forecasts for the original series via projections. We refer to this method as Forecast Linear Augmented Projection (FLAP). We prove that the method reduces the forecast error variance of the original series in a way that is agnostic both with respect to the weights of the linear combinations used at step (i) and with respect to the model used to generate forecasts at step (ii). The model is inspired by the forecast reconciliation literature (Athanasopoulos et al. 2023) whereby forecasts are adjusted to cohere with known linear constraints. In contrast to that literature, the FLAP method focuses on multivariate forecasting where such constraints are not present. Indeed, the method need not only be applied to forecasting problems, but multivariate predictions in general.

It may appear puzzling that forecast accuracy can be improved, not by introducing any new information, but by simply taking linear combinations of existing time series. To give an intuition into how this puzzle can be resolved, we consider a toy example of two series y_1 and y_2 that are of concern to the forecaster and two linear combinations or *components* of the series $c_1 = 0.5z_1 + 0.5z_2$ and $c_2 = 0.5z_1 - 0.5z_2$. Denote by $\hat{y}_1, \hat{y}_2, \hat{c}_1, \hat{c}_2$ any forecasts of these original series and components, which we collectively refer to as *base forecasts*. The base forecasts may be generated by univariate methods, multivariate methods, or even based on expert judgement. When considering y_1 , there is both a *direct* forecast \hat{y}_1 and *indirect* forecast $\hat{c}_1 + \hat{c}_2$, similarly for y_2 the direct forecast is \hat{y}_2 and the indirect forecast $\hat{c}_1 - \hat{c}_2$.¹ With the exception of some pathological cases, the direct and indirect forecasts for the same variable will not, in general, be equal. Therefore forecast accuracy can be improved by combining direct and indirect forecasts, something implicitly achieved by the proposed FLAP method. The puzzle is thus resolved; while no new information is created at the data augmentation step, there is new information embedded into forecasts of the augmented series, which can be leveraged via model combination (see Wang et al. (2023) for a review of forecast combination). Something obscured by the simple toy example is the way our FLAP method differs from the usual forecast combination methods, in particular our combinations are potentially non-convex since they are obtained via projections, in a way that we now elaborate upon.

More formally and more generally, the FLAP method considers a vector of original series $\mathbf{y} \in \mathbb{R}^m$ and a vector of components $\mathbf{c} \in \mathbb{R}^p$. While $(\mathbf{y}', \mathbf{c}')'$ is a $p + m$ -vector, the construction of components as linear combinations of the original series implies that $(\mathbf{y}', \mathbf{c}')'$ lies on a linear subspace of at most dimension m . The corresponding vector of forecasts $(\hat{\mathbf{y}}', \hat{\mathbf{c}}')'$ will, in general, have support on \mathbb{R}^{m+p} . FLAP projects $(\hat{\mathbf{y}}', \hat{\mathbf{c}}')'$ onto the m -dimensional linear subspace on which $(\mathbf{y}', \mathbf{c}')'$ has support. The setup of this problem bears similarities to the well known problem of forecast reconciliation where Panagiotelis et al. (2021) provide similar geometric intuition, while Wickramasuriya, Athanasopoulos & Hyndman (2019), Athanasopoulos et al. (2017), and Di Fonzo & Girolimetto (2023) have all shown that reconciliation can reduce forecast error variance theoretically and empirically. However, we note that these papers establish that reconciliation improves forecast accuracy for the hierarchy *as a whole*. In the general multivariate setting that we consider, this would imply improvements in forecast accuracy for \mathbf{y} and \mathbf{c} taken together. This poses a problem if improvements in forecast accuracy for \mathbf{c} could be offset by a deterioration in forecast accuracy for \mathbf{y} , since the former are not of interest in and of themselves. A key insight we make in this paper in the hierarchical setting, is that reductions in forecast error variance accrue even for a subset of variables in the hierarchy. It is this contribution

¹This argument, as well as the terminology direct and indirect forecasts, is inspired by Hollyman, Petropoulos & Tipping (2021) who discuss this in the forecast reconciliation setting.

that allows us to propose a method that goes beyond the case where time series adhere to linear constraints, and that instead applies to the more general setting.

While the theoretical results apply for any linear combinations of the original series, in practice we propose to augment the data with principal components. When doing so, the FLAP method bears a resemblance to Dynamic Factor Models (DFMs), specifically those common in macroeconomic forecasting (Stock & Watson 2002b,a, 2012), their extensions in the machine learning literature (De Stefani et al. 2019) as well as the factor augmented VAR (Bernanke, Boivin & Elias 2005). The factor models assume that the multivariate time series possesses common components and the dynamics of the observed series are governed by the dynamics of these unobserved factors, often assumed to follow some parametric model. In contrast, FLAP is a post-forecasting step, indeed forecasts can even be made using a DFM and then further improved by implementing the FLAP method, something we demonstrate in Section 3 and Section 4.

In the sense that forecast accuracy can be improved without any new information, FLAP has parallels with bootstrap aggregation or “bagging” (Breiman 1996; Bergmeir, Hyndman & Benítez 2016). Bagging can reduce prediction variance without increasing bias (Hastie, Tibshirani & Friedman 2003), by mitigating model uncertainty (Petropoulos, Hyndman & Bergmeir 2018), and does so without introducing any new data, but rather resampled versions of the existing data. Our FLAP method also reduces forecast error variance without introducing new data, but using linear combinations of the existing data, rather than bootstrapping. The FLAP method (in addition to forecast reconciliation and bagging) can be viewed as contributing to the literature where forecasts are improved by combination and data augmentation methods. This includes the theta method (Assimakopoulos & Nikolopoulos 2000), temporal aggregation (Kourentzes, Petropoulos & Trapero 2014; Athanasopoulos et al. 2017), forecasting with sub-seasonal series (FOSS, Li, Petropoulos & Kang 2022) and forecast combination with multiple starting points (Disney & Petropoulos 2015); a review of all these methods can be found in Petropoulos & Spiliotis (2021) who refer to them as using “the wisdom of data”. Our FLAP method is distinct in that it aims to exploit information in the data with a focus on linear combinations of multivariate series.

The remainder of the paper is structured as follows. In Section 2, we propose the FLAP method, and highlight its theoretical properties and associated estimation methods. In Section 3, we present a simulation example demonstrating its performance and discuss the implications for sources of uncertainty. Section 4 examines the performance of FLAP in two empirical applications: forecasting Australian domestic tourism and forecasting macroeconomic variables in the FRED-MD data set. Section 5 concludes with some thoughts on future research directions. The methods introduced in

this paper are implemented in the `flap` package (Yang 2024). This paper is fully reproducible with code and documentation provided at <https://github.com/FinYang/paper-forecast-projection>.

2 Forecast Linear Augmented Projection (FLAP)

2.1 Method and theoretical properties

In the following, all vectors and matrices are denoted in bold font. We use I_n to denote the $n \times n$ identity matrix, and $O_{n \times k}$ to denote the $n \times k$ zero matrix.

Let $\mathbf{y}_t \in \mathbb{R}^m$ be a vector of m observed time series we are interested in forecasting. The FLAP method involves three steps:

1. *Form components.* Form $\mathbf{c}_t = \Phi \mathbf{y}_t \in \mathbb{R}^p$, a vector of p linear combinations of \mathbf{y}_t at time t , where $\Phi \in \mathbb{R}^{p \times m}$. We call \mathbf{c}_t the components of \mathbf{y}_t and the component weights Φ are known in the sense that they are chosen by the user of FLAP. Let $\mathbf{z}_t = [\mathbf{y}_t', \mathbf{c}_t']'$ be the concatenation of series \mathbf{y}_t and components \mathbf{c}_t . We note that \mathbf{z}_t will be constrained in the sense that $\mathbf{C} \mathbf{z}_t = \mathbf{c}_t - \Phi \mathbf{y}_t = \mathbf{0}$ for any t where $\mathbf{C} = [-\Phi \ I_p]$ is referred to as the constraint matrix.
2. *Generate forecasts.* Denote as $\hat{\mathbf{z}}_{t+h}$ the h -step-ahead base forecast of \mathbf{z}_t . The method used to generate forecasts is again selected by the user. This can be univariate or multivariate. In the more general setting where \mathbf{z}_t are not time series but cross sectional data, any prediction method can be used. In general, the constraints that hold for \mathbf{z}_t will not hold for $\hat{\mathbf{z}}_{t+h}$, i.e. $\mathbf{C} \hat{\mathbf{z}}_{t+h} \neq \mathbf{0}$
3. *Project the base forecasts.* Let $\tilde{\mathbf{z}}_{t+h}$ be a set of projected forecasts such that,

$$\tilde{\mathbf{z}}_{t+h} = \mathbf{M} \hat{\mathbf{z}}_{t+h} \quad (1)$$

with projection matrix

$$\mathbf{M} = \mathbf{I}_{m+p} - \mathbf{W}_h \mathbf{C}' (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \mathbf{C}, \quad (2)$$

where $\text{Var}(\mathbf{z}_{t+h} - \hat{\mathbf{z}}_{t+h}) = \mathbf{W}_h$ is the forecast error covariance matrix. For the proofs of this section, we will assume that \mathbf{W}_h is known, in practice a plug-in estimate can be used that will be discussed in Section 2.5.

In practice we are interested in forecasts of \mathbf{y}_t and not the full vector \mathbf{z}_t . We now introduce some notation to handle this issue. Define the selection matrix $\mathbf{J}_{n,k} = [\mathbf{I}_n \ O_{n \times k}]$, so that $\mathbf{J}_{n,k} \mathbf{A}$ selects the first n rows of a matrix \mathbf{A} . Let $\hat{\mathbf{y}}_{t+h}$ and $\tilde{\mathbf{y}}_{t+h}$ denote the first m elements of $\hat{\mathbf{z}}_{t+h}$ and $\tilde{\mathbf{z}}_{t+h}$, comprising the base and projected forecasts of \mathbf{y}_t respectively. Similarly, let $\hat{\mathbf{c}}_{t+h}$ and $\tilde{\mathbf{c}}_{t+h}$ denote the last p

elements of $\hat{\mathbf{z}}_{t+h}$ and $\tilde{\mathbf{z}}_{t+h}$, comprising the base and projected forecasts of \mathbf{c}_t respectively. Then the projected forecast of \mathbf{y}_t can be found by

$$\tilde{\mathbf{y}}_{t+h} = \mathbf{J}\tilde{\mathbf{z}}_{t+h} = \mathbf{J}\mathbf{M}\hat{\mathbf{z}}_{t+h}, \quad (3)$$

where $\mathbf{J} = \mathbf{J}_{m,p}$.

We now present some theoretical results regarding the FLAP method, with proofs provided in the Appendix. Theorem 2.1 establishes that the forecasts produced by the FLAP method, $\tilde{\mathbf{y}}_{t+h}$, dominate the base forecasts, $\hat{\mathbf{y}}_{t+h}$, in the sense that the difference between their forecast error variances is always positive definite. Theorem 2.2 establishes that the trace of the covariance of the forecast errors of $\tilde{\mathbf{y}}_{t+h}$ is non-increasing with the number of components p . The conditions needed to make the trace strictly decreasing are discussed in Theorem 2.3. In Theorem 2.4 we prove that the projection in Equation 2 achieves the minimum forecast error variance amongst the class of all projections. Finally, while the theoretical results imply that components could in principle continue to be added to improve forecasts, in practice, larger values of p will make estimation of the plug-in covariance matrix \mathbf{W}_h unreliable. We explore this issue in a simulation setting and empirically in Section 3 and Section 4, respectively.

2.2 Positive Semi-Definiteness of Error Variance Reduction

We first provide some intermediate results.

Lemma 2.1. *Matrix \mathbf{M} is a projection onto the space where the constraint $\mathbf{C}\mathbf{z}_t = \mathbf{0}$ is satisfied.*

Proof in Appendix, page 31.

Based on the attractive properties of projections, we have the following corollaries.

Corollary 2.1.

1. The projected forecast $\tilde{\mathbf{z}}_{t+h}$ satisfies the constraint $\mathbf{C}\tilde{\mathbf{z}}_{t+h} = \mathbf{0}$.
2. For \mathbf{z}_{t+h} that already satisfies the constraint, the projection does not change its value, i.e., $\mathbf{M}\mathbf{z}_{t+h} = \mathbf{z}_{t+h}$ (Rao 1974, Lemma 2.4).
3. If the base forecasts are unbiased such that $E(\hat{\mathbf{z}}_{t+h}|\mathcal{I}_t) = E(\mathbf{z}_{t+h}|\mathcal{I}_t)$, then the projected forecasts are also unbiased, i.e., $E(\tilde{\mathbf{z}}_{t+h}|\mathcal{I}_t) = E(\mathbf{z}_{t+h}|\mathcal{I}_t)$.

Proof in Appendix, page 31.

The assumption of unbiasedness of the base forecasts is not unreasonable in practice, and where it does not hold, a bias correction method can be applied. Note this is not a requirement on model

specification. We do not assume the model producing the base forecast is correctly specified like in the DFM literature (e.g., Stock & Watson 2002a). In fact, the power of FLAP manifests when the models are misspecified, as discussed in Section 3.

Lemma 2.2. *The forecast error covariance matrix of the component-augmented projected h -step-ahead forecasts $\tilde{\mathbf{z}}_{t+h}$ is*

$$\text{Var}(\mathbf{z}_{t+h} - \tilde{\mathbf{z}}_{t+h}) = \mathbf{M} \mathbf{W}_h \mathbf{M}' = \mathbf{M} \mathbf{W}_h,$$

and the forecast error covariance matrix of the projected h -step-ahead forecasts $\tilde{\mathbf{y}}_{t+h}$ is

$$\text{Var}(\mathbf{y}_{t+h} - \tilde{\mathbf{y}}_{t+h}) = \mathbf{J} \mathbf{M} \mathbf{W}_h \mathbf{J}'.$$

Proof in Appendix, page 31.

Lemma 2.2 is a well-known result in the forecast reconciliation literature (e.g., Di Fonzo & Girolimetto 2023).

Theorem 2.1 (Positive Semi-Definiteness of Error Variance Reduction). *The difference between the forecast error variances of the base and projected component-augmented forecasts,*

$$\begin{aligned} \text{Var}(\mathbf{z}_{t+h} - \hat{\mathbf{z}}_{t+h}) - \text{Var}(\mathbf{z}_{t+h} - \tilde{\mathbf{z}}_{t+h}) &= \mathbf{W}_h - \mathbf{M} \mathbf{W}_h \\ &= \mathbf{W}_h - (\mathbf{I} - \mathbf{W}_h \mathbf{C}' (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \mathbf{C}) \mathbf{W}_h \\ &= \mathbf{W}_h \mathbf{C}' (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \mathbf{C} \mathbf{W}_h, \end{aligned}$$

is positive semi-definite. The difference between the forecast error variances of the base and projected forecasts of the original series,

$$\text{Var}(\mathbf{y}_{t+h} - \hat{\mathbf{y}}_{t+h}) - \text{Var}(\mathbf{y}_{t+h} - \tilde{\mathbf{y}}_{t+h}) = \mathbf{J} \mathbf{W}_h \mathbf{C}' (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \mathbf{C} \mathbf{W}_h \mathbf{J}',$$

is therefore also positive semi-definite.

Proof in Appendix, page 32.

Theorem 2.1 is why FLAP works. The trace of $\mathbf{J} \mathbf{W}_h \mathbf{C}' (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \mathbf{C} \mathbf{W}_h \mathbf{J}'$ is the sum of the reduction in forecast error variances, and is non-negative because the matrix is positive semi-definite. It implies that the forecast error variance can be reduced by simply forecasting the components (the artificially constructed linear combinations of the original data), and mapping the forecasts using matrix \mathbf{M} . For the improvement to be zero, the trace must be zero. This implies that the entire

$\mathbf{J}\mathbf{W}_h\mathbf{C}'(\mathbf{C}\mathbf{W}_h\mathbf{C}')^{-1}\mathbf{C}\mathbf{W}_h\mathbf{J}' = \mathbf{O}_{m \times m}$ as this is a positive semi-definite matrix, something rarely observed in practice. See Theorem 2.3 for more discussion.

The following example illustrates the mechanism of the reduction in the forecast error variance.

Example 2.1. Suppose \mathbf{z}_t comprises m original series \mathbf{y}_t and p components \mathbf{c}_t . Let $\hat{\mathbf{z}}_{t+h}$ and $\tilde{\mathbf{z}}_{t+h}$ be h -step-ahead base and projected forecasts of \mathbf{z}_t . Assume that their corresponding forecast errors are uncorrelated with unit variance such that $\mathbf{W}_h = \mathbf{I}_{m+p}$, then

$$\begin{aligned} \text{Var}(\mathbf{z}_{t+h} - \hat{\mathbf{z}}_{t+h}) - \text{Var}(\mathbf{z}_{t+h} - \tilde{\mathbf{z}}_{t+h}) &= \mathbf{C}'(\mathbf{C}\mathbf{C}')^{-1}\mathbf{C} \\ &= \begin{bmatrix} -\Phi' \\ \mathbf{I}_p \end{bmatrix} (\Phi\Phi' + \mathbf{I}_p)^{-1} \begin{bmatrix} -\Phi & \mathbf{I}_p \end{bmatrix}, \end{aligned}$$

where

$$\mathbf{C} = \begin{bmatrix} -\Phi & \mathbf{I}_p \end{bmatrix}.$$

Let Φ consist of $p \leq m$ orthogonal unit vectors, for example, those obtained from Principal Component Analysis (PCA, Jolliffe 2002). In this case $\Phi\Phi' = \mathbf{I}_p$ and

$$\text{Var}(\mathbf{z}_{t+h} - \hat{\mathbf{z}}_{t+h}) - \text{Var}(\mathbf{z}_{t+h} - \tilde{\mathbf{z}}_{t+h}) = \frac{1}{2} \begin{bmatrix} \Phi'\Phi & -\Phi' \\ -\Phi & \mathbf{I}_p \end{bmatrix}.$$

Focusing on the forecast error variance reduction of the forecasts of the original series \mathbf{y}_t , i.e., $\text{tr}(\text{Var}(\mathbf{y}_{t+h} - \hat{\mathbf{y}}_{t+h}) - \text{Var}(\mathbf{y}_{t+h} - \tilde{\mathbf{y}}_{t+h})) = \frac{1}{2} \text{tr}(\Phi'\Phi)$.

- When $p < m$, since $\Phi'\Phi$ is idempotent, $\text{tr}(\Phi'\Phi) = \text{rank}(\Phi'\Phi) = p$. Hence, focusing on the original series the reduction in the total forecast error variance of the FLAP forecasts relative to the base forecasts, is $p/2$.
- When $p = m$, where all principal components are used, $\Phi'\Phi = \mathbf{I}_m$. This implies a total reduction in the error variance of $m/2$, and that for each of the m individual series the error variance is halved.

If we keep increasing the number of components beyond m , the result in Theorem 2.1 still holds, although Φ can no longer contain orthogonal vectors, and the example here becomes intractable. This is an artificial example as the forecast error variance \mathbf{W}_h can hardly be an identity in practice. It is likely that the forecast error of a linear combination of series will be correlated to the forecast error of forecasting these series directly. Nonetheless, the aim of the example is to demonstrate how the forecast error variance can be reduced as a result of the component forecasts bringing new information about the original series. The forecast error variance reduction becomes larger as we

increase the number of components p . This is not a coincidence but a desirable property of FLAP, as shown in the next section.

2.3 Monotonicity

In the results that follow, we break the base forecast error covariance matrix into smaller blocks.

$$\mathbf{W}_h = \begin{bmatrix} \mathbf{W}_{y,h} & \mathbf{W}_{yc,h} \\ \mathbf{W}'_{yc,h} & \mathbf{W}_{c,h} \end{bmatrix},$$

where $\mathbf{W}_{y,h}$ is the forecast error covariance matrix of $\hat{\mathbf{y}}_{t+h}$, $\mathbf{W}_{c,h}$ is the forecast error covariance matrix of $\hat{\mathbf{c}}_{t+h}$, and $\mathbf{W}_{yc,h}$ contains error covariances between elements of $\hat{\mathbf{y}}_{t+h}$ and $\hat{\mathbf{c}}_{t+h}$.

Theorem 2.2 (Monotonicity). *The forecast error variance reductions for each series, i.e., the diagonal elements in the matrix*

$$\text{Var}(\mathbf{y}_{t+h} - \hat{\mathbf{y}}_{t+h}) - \text{Var}(\mathbf{y}_{t+h} - \tilde{\mathbf{y}}_{t+h}) = \mathbf{J} \mathbf{W}_h \mathbf{C}' (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \mathbf{C} \mathbf{W}_h \mathbf{J}'$$

is non-decreasing as p increases. In particular, the sum of forecast error variance reductions

$$\text{tr}(\text{Var}(\mathbf{y}_{t+h} - \hat{\mathbf{y}}_{t+h}) - \text{Var}(\mathbf{y}_{t+h} - \tilde{\mathbf{y}}_{t+h})) = \text{tr}(\mathbf{J} \mathbf{W}_h \mathbf{C}' (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \mathbf{C} \mathbf{W}_h \mathbf{J}') \quad (4)$$

is non-decreasing as p increases.

Proof in Appendix, page 32.

Theorem 2.2 is the key result that demonstrates the usefulness of FLAP. It means that we can keep increasing the number of components to reduce forecast error variance, even when the number of components exceeds the number of original series. It requires \mathbf{C} to be $[-\Phi \quad \mathbf{I}_p]$ or $[-\Phi \quad \mathbf{L}]$ where \mathbf{L} is a lower triangular matrix. This implies that the components can also be constructed from existing components, not only from the original series. This has little significance since a linear combination of components of the original series, is just a linear combination of the original series. This of course assumes that the forecast error covariances are known, something we explore in Section 2.5 and Section 3.

Extending the proof of Theorem 2.2, we can outline the condition for the reduced sum of forecast error variances to be positive. Denote ϕ_i as the row vector containing the weights associated with the i th component, so that with p components, the weights matrix is $\Phi = [\phi_1' \quad \phi_2' \quad \dots \quad \phi_p']'$. Let $\mathbf{W}_{\tilde{\mathbf{y}},h}^{(i-1)}$ denote the error covariance matrix of the projected forecasts of the original series based on the first $i-1$ components, $\mathbf{w}_{c_1\tilde{\mathbf{y}},h}$ denote a vector of covariances of the first component and the base forecasts

of the original series, and $\mathbf{w}_{c_i\tilde{y},h}^{(i-1)}$ denote a vector of covariances of the projected i th component and the projected forecasts of the original series, based on the first $i - 1$ components.

Theorem 2.3 (Positive Forecast Error Variance Reduction Condition). *For the first component to achieve a guaranteed reduction of forecast error variance (for the reduced variance matrix in Theorem 2.1 to have positive trace),*

$$\phi_1 \mathbf{W}_{y,h} \neq \mathbf{w}_{c_1y,h}. \quad (5)$$

For the i th component to have a positive reduction on forecast error variance of the original series,

$$\phi_i \mathbf{W}_{\tilde{y},h}^{(i-1)} \neq \mathbf{w}_{c_i\tilde{y},h}^{(i-1)}. \quad (6)$$

Proof in Appendix, page 35.

The condition of Equation 5 demonstrates that for a new component to be beneficial, the information introduced by this new component, reflected in the error covariance, cannot be a linear combination of already existing information.

Theorem 2.3 can potentially provide insights into the selection of component weights and forecast models to satisfy the conditions. We leave this issue to a later article, as practically the conditions in Theorem 2.3 are either almost always satisfied if the weights are simulated randomly on a continuous scale, or the loss associated with the rare occasions where the conditions are not satisfied is negligible compared to the estimation error imposed by the limited sample size as the number of components increases, as discussed in Section 3 and Section 4.

2.4 Optimality of the projection

Equation 1 can be seen as a solution to the optimisation problem:

$$\arg \min_{\check{\mathbf{z}}_{T+h}} (\hat{\mathbf{z}}_{T+h} - \check{\mathbf{z}}_{T+h})' \mathbf{W}_h^{-1} (\hat{\mathbf{z}}_{T+h} - \check{\mathbf{z}}_{T+h}) \quad \text{s.t. } \mathbf{C} \check{\mathbf{z}}_{T+h} = 0.$$

If we consider the transformed space where all the vectors are first transformed via pre-multiplying by $\mathbf{W}_h^{-1/2}$, where $\mathbf{W}_h^{-1} = (\mathbf{W}_h^{-1/2})' \mathbf{W}_h^{-1/2}$, then this optimisation problem can be interpreted as finding the set of forecasts that are closest to the base forecasts on the transformed space, while satisfying the linear constraints imposed by the components.

Moreover, this is equivalent to the optimisation problem:

$$\arg \min_{\check{\mathbf{z}}_{T+h}} (\hat{\mathbf{z}}_{T+h} - \check{\mathbf{z}}_{T+h})' \mathbf{W}_h^{-1} (\hat{\mathbf{z}}_{T+h} - \check{\mathbf{z}}_{T+h}) \quad \text{s.t. } \Phi \check{\mathbf{y}}_{T+h} = \check{\mathbf{c}}_{T+h},$$

where $\check{\mathbf{z}}_{T+h}$ is the vector of the last p elements of $\hat{\mathbf{z}}_{T+h}$, corresponding to the forecast of the components as part of the solution. This equivalence is discussed in Wickramasuriya, Athanasopoulos & Hyndman (2019), where the authors find the solution by minimising the sum of forecast error variance of all series (See Ando & Narita (2022) for a simpler proof). The result is

$$\tilde{\mathbf{z}}_{t+h} = \mathbf{S}\mathbf{G}\hat{\mathbf{z}}_{t+h}, \quad (7)$$

where $\mathbf{S} = \begin{bmatrix} \mathbf{I}_m \\ \Phi \end{bmatrix}$ contains the constraints, so that $\mathbf{z}_t = \mathbf{S}\mathbf{y}_t$, and

$$\mathbf{G} = (\mathbf{S}'\mathbf{W}_h^{-1}\mathbf{S})^{-1}\mathbf{S}'\mathbf{W}_h^{-1}. \quad (8)$$

In Equation 7, $\mathbf{G}\hat{\mathbf{z}}_{t+h}$ can be viewed as mapping of all the series to a selected few. In the forecast reconciliation context, this is a mapping of all series to the “bottom level”. In our multivariate forecasting context, this is a mapping of all series including the components, to the space of the original series. This leads to the solution

$$\tilde{\mathbf{y}}_{t+h} = \mathbf{G}\hat{\mathbf{z}}_{t+h}, \quad (9)$$

as equivalent to Equation 3. Recognising that Equation 7 is equivalent to Equation 1, it is the solution that minimises the sum of forecast error variances of the original series and all the components. We go further in Theorem 2.4, and show that Equation 7 is also the solution to minimise each individual forecast error variance of the original series, and their sum. This can be viewed as a special case of Theorem 3.3 in Panagiotelis et al. (2021), or as illustrated by Ando & Narita (2022), but applied in a different context to forecast reconciliation. The earliest work we can find that noted this interpretation in a non-forecasting context is Luenberger (1969, p.85). We establish a few basic results leading to the optimality of this solution first, also to check that Lemma 2.1, Corollary 2.1 and Lemma 2.2 hold under this alternative representation.

Lemma 2.3. *The matrix $\mathbf{S}\mathbf{G}$ is a projection onto the space where the constraint $\mathbf{C}\mathbf{z}_t = \mathbf{0}$ is satisfied, provided that $\mathbf{G}\mathbf{S} = \mathbf{I}$.*

Proof in Appendix, page 36.

Corollary 2.2. *Provided that $\mathbf{G}\mathbf{S} = \mathbf{I}$, the following results hold.*

1. *The projected forecast in Equation 9 satisfies the constraint $\mathbf{C}\tilde{\mathbf{z}}_{t+h} = \mathbf{C}\mathbf{S}\tilde{\mathbf{y}}_{t+h} = \mathbf{0}$.*

2. For \mathbf{z}_{t+h} that already satisfies the constraint, the mapping does not change its value, i.e., $\mathbf{G}\mathbf{z}_{t+h} = \mathbf{y}_{t+h}$.
3. If the base forecasts are unbiased such that $E(\hat{\mathbf{z}}_{t+h}|\mathcal{I}_t) = E(\mathbf{z}_{t+h}|\mathcal{I}_t)$, then the projected forecasts in Equation 9 are also unbiased: $E(\tilde{\mathbf{y}}_{t+h}|\mathcal{I}_t) = E(\mathbf{y}_{t+h}|\mathcal{I}_t)$.

Proof in Appendix, page 37.

Lemma 2.4. The covariance matrix of the projected forecasts from Equation 9 is given by

$$\text{Var}(\mathbf{y}_{t+h} - \tilde{\mathbf{y}}_{t+h}) = \mathbf{G}\mathbf{W}_h\mathbf{G}'.$$

Proof in Appendix, page 37.

We are now ready to present the following theorem.

Theorem 2.4 (Minimum Variance Unbiased Projected Forecast). The solution to

$$\arg \min_{\mathbf{G}} \text{tr}(\mathbf{G}\mathbf{W}_h\mathbf{G}') \quad \text{s.t. } \mathbf{G}\mathbf{S} = \mathbf{I} \quad (10)$$

is Equation 8. This problem can be effectively split into independent subproblems such that $\mathbf{G} = [\mathbf{g}_1 \ \mathbf{g}_2 \ \dots \ \mathbf{g}_m]'$, where \mathbf{g}_i is the solution to the subproblem of the i th series

$$\arg \min_{\mathbf{g}_i} \mathbf{g}_i' \mathbf{W}_h \mathbf{g}_i \quad \text{s.t. } \mathbf{g}_i' \mathbf{s}_j = \delta_{ij}, \quad j = 1, 2, \dots, m, \quad (11)$$

where \mathbf{s}_j is the j th column of \mathbf{S} , and δ_{ij} is the Kronecker delta function taking value 1 if $i = j$ and 0 otherwise.

Proof in Appendix, page 38.

In other words, the forecast projection method gives optimal projected forecast for a given set of components, in the sense that the unbiased forecast of each series has minimum variance.

2.5 Estimation of \mathbf{W}_h

In practice, the base forecast error variance \mathbf{W}_h is unknown and needs to be estimated. Denote $\hat{\mathbf{e}}_{t,h} = \mathbf{z}_t - \hat{\mathbf{z}}_{t|h-h}$ as the h -step-ahead base forecast in-sample residual. The conventional forecast error variance matrix estimator

$$\widehat{\mathbf{W}}_h = \frac{1}{T-h-1} \sum_{t=h+1}^T \hat{\mathbf{e}}_{t,h} \hat{\mathbf{e}}_{t,h}',$$

albeit unbiased, is not considered a good approximation to the true forecast error variance in a finite sample when $(m+p) \approx T-h$. It is even singular when $(m+p) > T-h$, which makes the quantities

discussed in the previous sections impossible to calculate. For this reason, while other shrinkage estimators can be considered, we adopt the covariance shrinkage method of Schäfer & Strimmer (2005) and the variance shrinkage method of Opgen-Rhein & Strimmer (2007). Then, the estimated forecast error variance matrix is guaranteed to be positive definite with few numerical problems. This estimator is denoted as $\widehat{\mathbf{W}}_h^{shr} = (\hat{w}_{ij,h}^{shr})_{1 \leq i,j \leq m+p}$ with the element in row i and column j

$$\hat{w}_{ij,h}^{shr} = \hat{r}_{ij,h}^{shr} \sqrt{\hat{v}_{i,h} \hat{v}_{j,h}},$$

where $\hat{r}_{ij,h}^{shr} = (1 - \hat{\lambda}_{cor})\hat{r}_{ij,h}$ and $\hat{v}_{i,h} = \hat{\lambda}_{var}\hat{w}_{h,median} + (1 - \hat{\lambda}_{var})\hat{w}_{i,h}$, with $\hat{\lambda}_{cor}$ being the shrinkage intensity parameter for the correlation

$$\hat{\lambda}_{cor} = \min\left(1, \frac{\sum_{i \neq j} \widehat{\text{var}}(\hat{r}_{ij,h})}{\sum_{i \neq j} \hat{r}_{ij,h}^2}\right),$$

and $\hat{\lambda}_{var}$ being the shrinkage intensity parameter for the variance

$$\hat{\lambda}_{var} = \min\left(1, \frac{\sum_{i=1}^{m+p} \widehat{\text{var}}(\hat{w}_{i,h})}{\sum_{i=1}^{m+p} (\hat{w}_{i,h} - \hat{w}_{h,median})^2}\right),$$

$\hat{r}_{ij,h}$ the sample correlation of the h -step-ahead forecast error between the i th and the j th series (component) in \mathbf{z}_t , $\hat{w}_{i,h}$ the h -step-ahead sample base forecast error variance associated with the i th series (the i th diagonal element of $\widehat{\mathbf{W}}_h$), and $\hat{w}_{h,median}$ the median of the h -step-ahead sample forecast error variance of the series and components (the median of the diagonal elements of $\widehat{\mathbf{W}}_h$). The estimation of $\widehat{\mathbf{W}}_h^{shr}$ in the following sections are implemented using the package `corpcor` (Schafer et al. 2021) in R (R Core Team 2023).

Estimating $\widehat{\mathbf{W}}_h^{shr}$ for each forecast horizon h is desirable but computationally intensive. It involves the calculation of multi-step-ahead in-sample residuals of the forecast models, which is especially challenging for iterative forecasts. Because of this, in practice it is not unreasonable to assume the h -step forecast error variance is proportional to the 1-step forecast error variance by a constant η_h , as do Wickramasuriya, Athanasopoulos & Hyndman (2019):

$$\widehat{\mathbf{W}}_h^{shr} = \eta_h \widehat{\mathbf{W}}_1^{shr}.$$

Under this assumption, when $\widehat{\mathbf{W}}_h^{shr}$ is used in Equation 2, the proportionality constant η_h cancels out regardless of the value of h . We can effectively use only the one-step forecast error variance in forecast projection, if we only need point forecasts. We calculate $\widehat{\mathbf{W}}_h^{shr}$ for each h for the simulation example in Section 3, but assume this proportionality for the application in Section 4.

3 Simulation

In this section, we illustrate the performance of FLAP in a simulation setting. We generate time series of length $T = 400$ from a $m = 70$ variable VAR(3). The coefficients for the VAR DGP are estimated from the first 70 series in the Australian tourism data set used in Section 4.1. The innovations are simulated from a multivariate normal distribution with an identity covariance matrix. The estimation and simulation are done using package `tsDyn` (Fabio Di Narzo, Aznarte & Stigler 2009).

For each sample we generate $h = 1$ to 12-step-ahead base forecasts from benchmark models. We implement FLAP with a varying number of components and component construction methods. This process is repeated 220 times. The improvement of FLAP forecasts over base forecasts is assessed, and the statistical significance of these improvements is evaluated.

3.1 Benchmarks for generating base forecasts

The first benchmark is the univariate ARIMA model. For each series, we fit an ARIMA model using the `auto.arima()` function from the `forecast` package (Hyndman et al. 2023) with the default settings.

The second benchmark is the dynamic factor model (DFM). Following Stock & Watson (2002b),

$$\hat{y}_{T+h} = \hat{\alpha}_h + \sum_{j=1}^n \hat{\beta}'_{hj} \hat{F}_{T-j+1} + \sum_{j=1}^s \hat{\gamma}_{hj} y_{T-j+1},$$

where \hat{F}_t is the vector of k estimated factors, and \hat{y}_t is the target series to forecast. The factors are estimated using PCA on demeaned and scaled data. The optimal model is selected for each series based on the Bayesian information criterion (BIC) from models fitted using different combinations of meta-parameters in their corresponding range: $1 \leq k \leq 6$, $1 \leq n \leq 3$ and $1 \leq s \leq 3$. We note that the DFM produces direct forecasts in the sense that a different model is fitted for each forecast horizon h , in contrast to indirect or iterative forecasts generated by the ARIMA models.

3.2 FLAP forecasts

We use several sets of weights to construct components for FLAP. The types of components are listed in Table 1. Principal components from PCA are established using the `prcomp` function in package `stats` (R Core Team 2023). We generate random components using weights simulated from a standard normal (Norm) distribution and a uniform (Unif) distribution with range $(-1, 1)$. We normalise the weights of all randomly generated components into unit vectors to maintain some level of consistency, with $\phi_i / \sqrt{\sum_j (\phi_{ij}^2)}$ where ϕ_{ij} is the j th value in the weight vector of the i th component. The total number of components p is selected to be either $m = 70$ or 300.

Table 1: *Component construction methods for FLAP*

Component	Description
PCA	Principal components from PCA.
Norm	The weights of all components are simulated from a standard normal distribution.
Unif	The weights of all components are simulated from a uniform distribution
PCA+Norm	m principal components from PCA are complemented with random components whose weights are simulated from a standard normal distribution.
PCA+Unif	m principal components from PCA are complemented with random components whose weights are simulated from a uniform distribution.
Ortho	A random orthonormal matrix is generated using package <code>pracma</code> (Borchers 2023) as the weight matrix.
Ortho+Norm	A random orthonormal matrix is generated using package <code>pracma</code> , forming the weights of the first m components. Weights of the additional components are simulated from a standard normal distribution.

We employ the Friedman test (Friedman 1937, 1939) along with post-hoc Nemenyi tests (Nemenyi 1963; Hollander, Wolfe & Chicken 2013) using the `tsutils` package (Kourentzes 2023) to compare forecast performance between different methods. The analysis involves the use of Multiple Comparisons with the Best (MCB) plot introduced by Koning et al. (2005) to visualise the comparison. The MSE of each series over different samples is calculated, and the MSEs of all the series are treated as observations in the Nemenyi test.

The average ranks are plotted in Figure 1 for forecast horizons 1, 6 and 12. The methods using FLAP are labelled “Model – Component Weights - Number of Components”. The benchmarks are labelled “Model – Benchmark”. These points are marked with triangles. The shaded region is the confidence interval of the best-performing method. Methods outside the shaded region are significantly worse than the best model.

In Figure 2, we plot the out-of-sample MSE values as a function of the number of components p . Here we include the performance of the true data generating process (DGP) VAR model (VAR – DGP), the VAR model with the correct specification (VAR – Est.) but using estimated parameters, and their corresponding FLAP generated forecasts. We do not include methods involving a uniform distribution and random orthonormal matrices, as they are visually identical to the methods with random weights generated from a standard normal distribution. The vertical black line indicates the number of series $m = 70$.

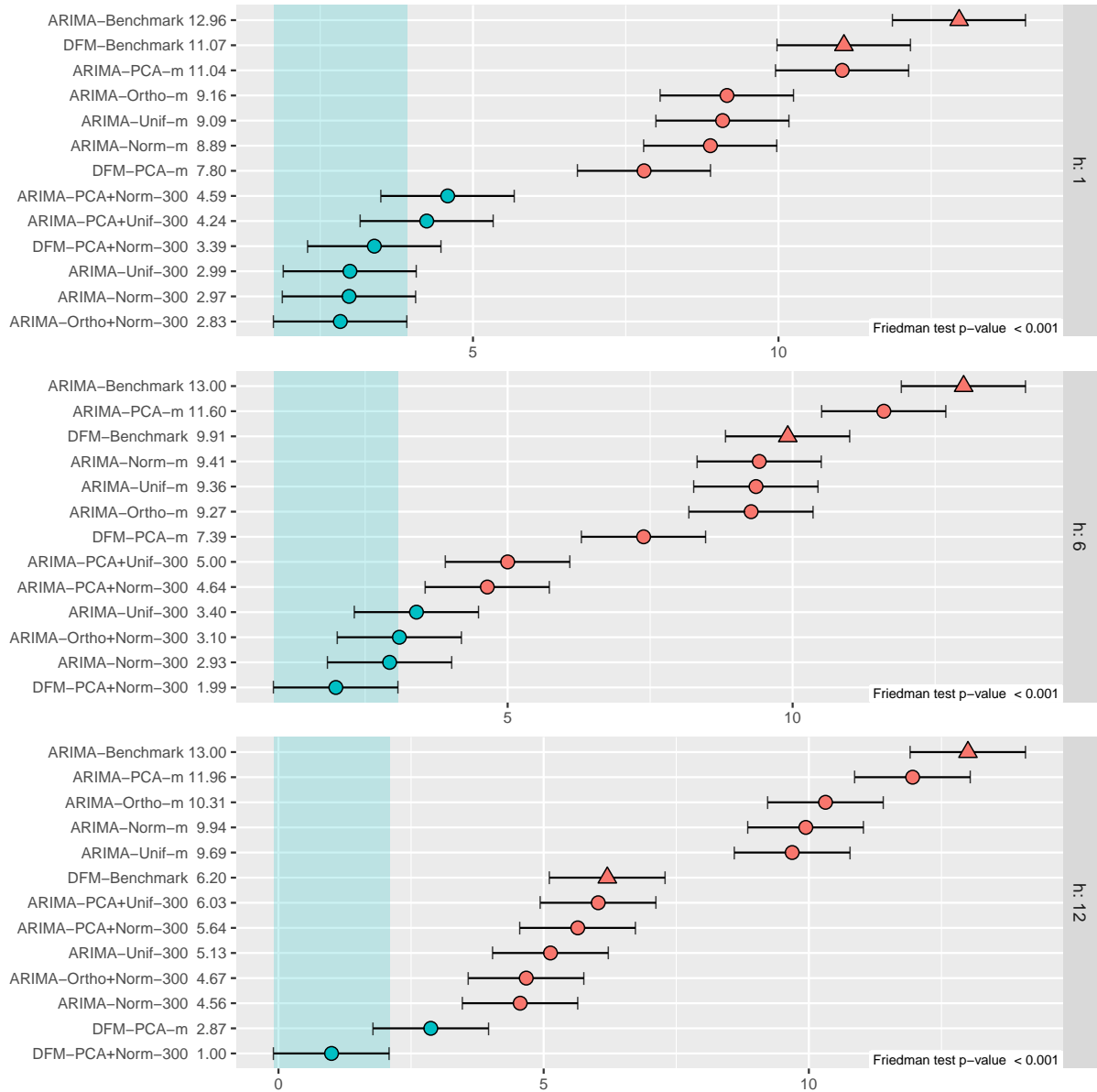


Figure 1: Average ranks of 1-, 6- and 12-step-ahead MSE of different model and component specifications in the simulation. The methods using FLAP are named as “Model – Component Weights – Number of Components”. The base models are named as “Model – Benchmark” and these points are marked with triangles. The shaded region is the confidence interval of the best performing model. Methods outside the shaded region are significantly worse than the best model.

3.3 FLAP over base forecasts

The first important observation from Figure 1 is the overall performance difference between the FLAP forecasts compared to the corresponding benchmark base forecasts. Note that we need to compare forecasts corresponding to the same model, i.e. FLAP ARIMA forecasts to base ARIMA forecasts and FLAP DFM forecasts to base DFM forecasts. The average ranks of the FLAP forecasts are better than the corresponding base forecasts for all forecast horizons, and the differences are

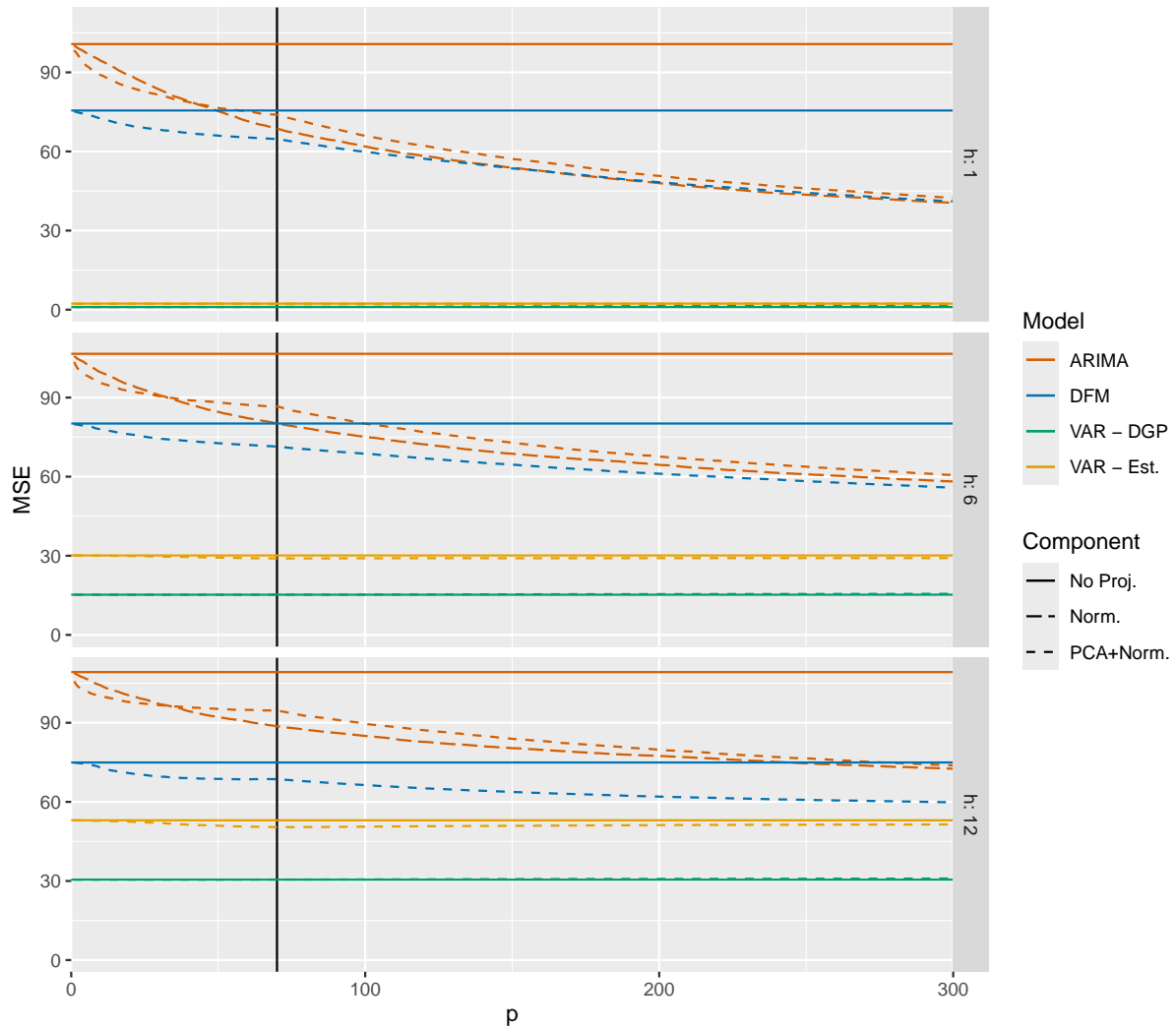


Figure 2: Out-of-sample MSE for base and FLAP forecasts as the number of components p increases, for forecast horizons 1, 6 and 12, in the simulation. “VAR – DGP” indicates the performance of the true data generating VAR model. “VAR – Est.” indicates the performance of the VAR model with the same structure as the true model but with estimated parameters. The solid horizontal lines show the MSE for the base forecasts while the dashed lines show the MSEs of the FLAP forecasts. The vertical black line indicates the location of $p = m = 70$, the number of series.

all statistically significant. The only exception is for the the PCA-related FLAP forecasts using only $m = 70$ components for forecast horizon 6 and 12. The number of components seems to be important. The best-performing models are all with 300 components. Between the one-step-ahead forecasts, the methods with 300 components are not significantly different from each other, regardless of the forecast model or how we construct the components.

Indeed, from Figure 2 we can see that the MSEs for ARIMA and DFM FALP forecasts keep decreasing from the base forecast as the number of components increases. This confirms Theorem 2.2, that the more components we include, the more forecast error variance reduction we can achieve. We should note that this is only obvious in this ideal setting where we have 400 observations in each group

while we only use at most 300 components in FLAP. This relatively large number of observations and the simple DGP can ease the challenge of estimation. This continued reduction in forecast error variance is not always achievable with real data, as we will see in Section 4, especially with FRED-MD in Section 4.2.

3.4 Base forecasts

Having simulated data from a VAR DGP, we expect the DFM to pick up the correlations between series, something not possible with univariate ARIMA. Hence we expect the DFM to perform better than the ARIMA. This is indeed the case. Figure 1 shows that the base DFM forecasts are significantly better than base ARIMA forecasts, with the exception for $h = 1$, where the difference is very close to significant. This is also observed in Figure 2. The solid horizontal line representing the MSE of the base DFM forecasts is always far below the horizontal solid line for the base ARIMA forecasts.

With the help of FLAP, a simple model like ARIMA can achieve comparable performance to more sophisticated models like the DFM. In Figure 1, all FLAP ARIMA forecasts except the one having m PCs are significantly better than base DFM at $h = 1$, and all FLAP ARIMA with 300 components are significantly better than base DFM at $h = 6$.

In Figure 2, the dashed lines of FLAP ARIMA forecasts decrease monotonically as the number of components p increases and reaches the MSE of base DFM at some point. This is because FLAP utilises shared information between series by capturing them in the components, making up for the overlooked cross-correlations in the univariate ARIMA models.

Interestingly enough, at $h = 1$, while the MSE of FLAP DFM forecasts also decreases as p increases, the MSE of the FLAP ARIMA forecasts and the MSE of the FLAP DFM forecasts seems to converge to the same value as p reaches 300, no matter how the components are constructed. Note that the same forecasts of the components, generated from univariate ARIMA models of these components, are used for both the FLAP ARIMA and the FLAP DFM methods. This implies that there is valuable information in the series that is not captured by either the ARIMA model or DFM, but is captured by the components. As the number of components increases, the information captured by the components overpowers the information captured by the base models, dominating the performance of the FLAP forecasts. Once again, this emphasises the importance of the components and FLAP. In this extreme case, the simple model is as good as the more complicated model after projection, while the forecast model itself is not as valuable as FLAP.

This observation is not as obvious as the forecast horizon increases. This is because while ARIMA produces forecasts iteratively, DFM is a direct forecast model. With this simple DGP, the performance

of DFM can be well maintained with larger h since a different model is fitted for each h . This can be seen as the MSE of the base DFM does not change much with different h , but the MSE of base ARIMA keeps increasing as h increases.

3.5 Component construction

The construction of components is obviously important in the proposed FLAP. However, the simulation results show that the method of constructing components may not be as important as one may have expected. In Figure 1, the main difference that can be observed is between using a combination of PCA and random weights, versus purely using random weights. The distribution that generates the random weights has limited effect. In Figure 1, with the same number of components and the same base ARIMA model, the MSEs are not significantly different, regardless of distribution from which the weights are simulated. The same conclusion can be drawn when PCA is used. As long as principal components are used, the performance is not different whether the additional components are generated from a normal distribution (PCA+Norm) or a uniform distribution (PCA+Unif).

When the weights are randomly simulated, the distribution is of limited importance. Therefore, in Figure 2, we look only at the inclusion of PCA and randomly generated components whose weights are from a standard normal distribution. When p is relatively small, the MSE decreases at a faster rate when PCA is used. Comparing the two dashed lines of the FLAP ARIMA forecasts, The performance of FLAP forecasts without PCA reaches and exceeds the performance with PCA before the number of components reaches m , and stays in the lead thereafter, although the gap seems to diminish with large p . This difference of PCA comes from the variances of principal components being maximised and ranked from largest to smallest. Because the performance of using random orthonormal weight matrix is the same as using only random normal weights, the difference of PCA is not from the orthogonality of the components. This might suggest the use of simple random weights if one is prepared to include a relatively large number of components in FLAP and to use PCA only when the number of components is small. However, as we will see in Section 4, this is not the case with real data. In this case, PCA seems to be the preferred approach even when the number of components is large.

Different constructions of components remain an important aspect of FLAP. One important future direction would be to find alternative and optimal components, as we do not limit the structure of the weight matrix in Section 2. This should be studied together with the selection of the forecast model since both the weight matrix Φ and the base forecast error variance \mathbf{W}_h can affect the projection simultaneously in Equation 2. This is likely to be an extension of the forecast combination literature, focusing on the properties of the base forecasts, and the diversity and robustness of the forecast model

and components. Examples of studies on this issue in the forecast combination literature include Batchelor & Dua (1995), Kang et al. (2022) and Lichtendahl & Winkler (2020).

3.6 Sources of uncertainty

At the bottom of each panel in Figure 2, the best forecasts come from the true DGP VAR model (the solid green line). The true VAR model FLAP forecasts do not improve on the base forecasts, as the uncertainty comes from the intrinsic error in the DGP that cannot be reduced. The second best forecasts come from the estimated VAR model, as the uncertainty, apart from the intrinsic error, only comes from the estimation error, not the model misspecification error like ARIMA and DFM, which are both misspecified in this simulation example. The gap between the estimated VAR and the true VAR becomes bigger for a longer forecast horizon, because VAR produces iterative forecasts, and estimation error accumulates as h becomes larger.

Forecast projection shows little, if any, improvement over the estimated VAR. This means that FLAP cannot reduce estimation error, which is termed as the parameter uncertainty in Petropoulos, Hyndman & Bergmeir (2018). On the other hand, it shows significant improvement over misspecified base models. This implies that the uncertainty FLAP can reduce is mainly the model misspecification error, referred to as the model uncertainty in Petropoulos, Hyndman & Bergmeir (2018). This is similar to bagging as bagging also reduces variance by controlling model uncertainty. Petropoulos, Hyndman & Bergmeir (2018) also examine the data uncertainty, which is not studied here. It is not clear how data uncertainty translates in a projection problem. However, the uncertainty in how the components are constructed, unique in FLAP problems, is discussed in Section 3.5 and awaits future research.

4 Empirical applications

Here we apply the FLAP method to two real data examples and draw most of the same conclusions we did from the simulations, with a few key differences.

4.1 Australian domestic tourism

The Australian Tourism Data Set compiled from the National Visitor Survey by Tourism Research Australia contains the total number of nights spent by Australians away from home. We refer to these as visitor nights. The monthly visitor nights are recorded for $m = 77$ regions around Australia, covering the period January 1998 to December 2019. To evaluate the forecast performance, we conduct expanding window time series cross-validation (Hyndman & Athanasopoulos 2021). The first $T = 84$ observations are used as the first training set. The following 12 months are used as the test set for evaluation. We repeat the evaluation for the rest of the data, by expanding each training

sample by one observation at a time. This generates 169 forecasts for each of the forecast horizons from 1 to 12 for evaluation.

The base forecasts for both the series and the components are produced using univariate ETS models selected and fitted using the `ets()` function in the `forecast` package (Hyndman et al. 2023; Hyndman & Khandakar 2008). See Hyndman & Athanasopoulos (2018, chap. 7) for more details.

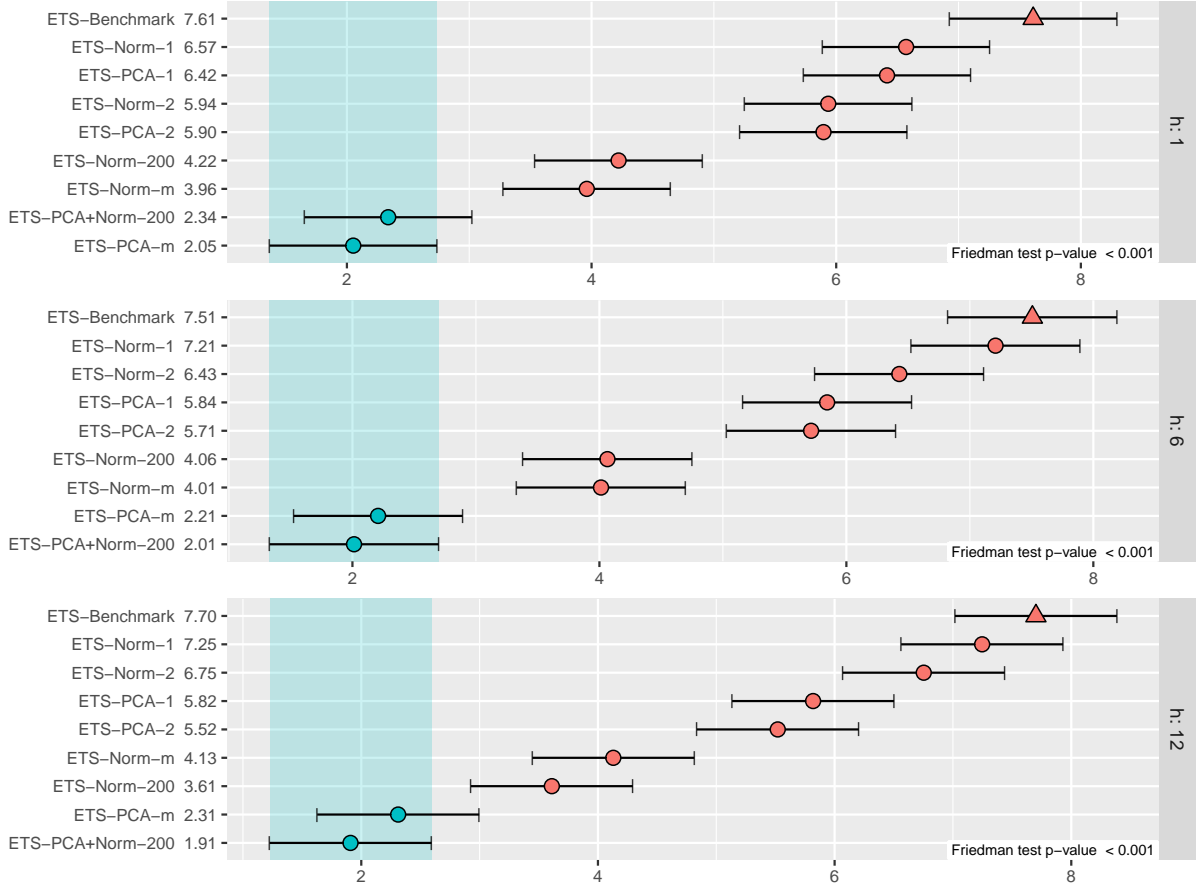


Figure 3: Average ranks of 1-, 6- and 12-step-ahead cross-validation MSE of different model and component specifications on the visitor nights data. The methods using forecast projection are named as “Model – Component Weights – Number of Components”. The base models are named as “Model – Benchmark” and these points are marked with triangles. The shaded region is the confidence interval of the best performing model. Methods outside the shaded region are significantly worse than the best model.

The MCB and MSE plots for $h = 1, 6, 12$ (as described in Section 3.2) are shown in Figure 3 and Figure 4. The FLAP forecasts in general outperforms the base forecasts, which is consistent with Section 3. In Figure 3, the base forecasts are always ranked last. Even FLAP with only one component are significantly better than the base forecasts for $h = 6$ and 12.

We highlight two different observations from Section 3. First, in Figure 4, the MSE of FLAP forecasts does not always decrease as the number of components increases — see for example for $h = 1$ with $p > 150$. This can also be seen from Figure 3, where the two best methods are not significantly different,

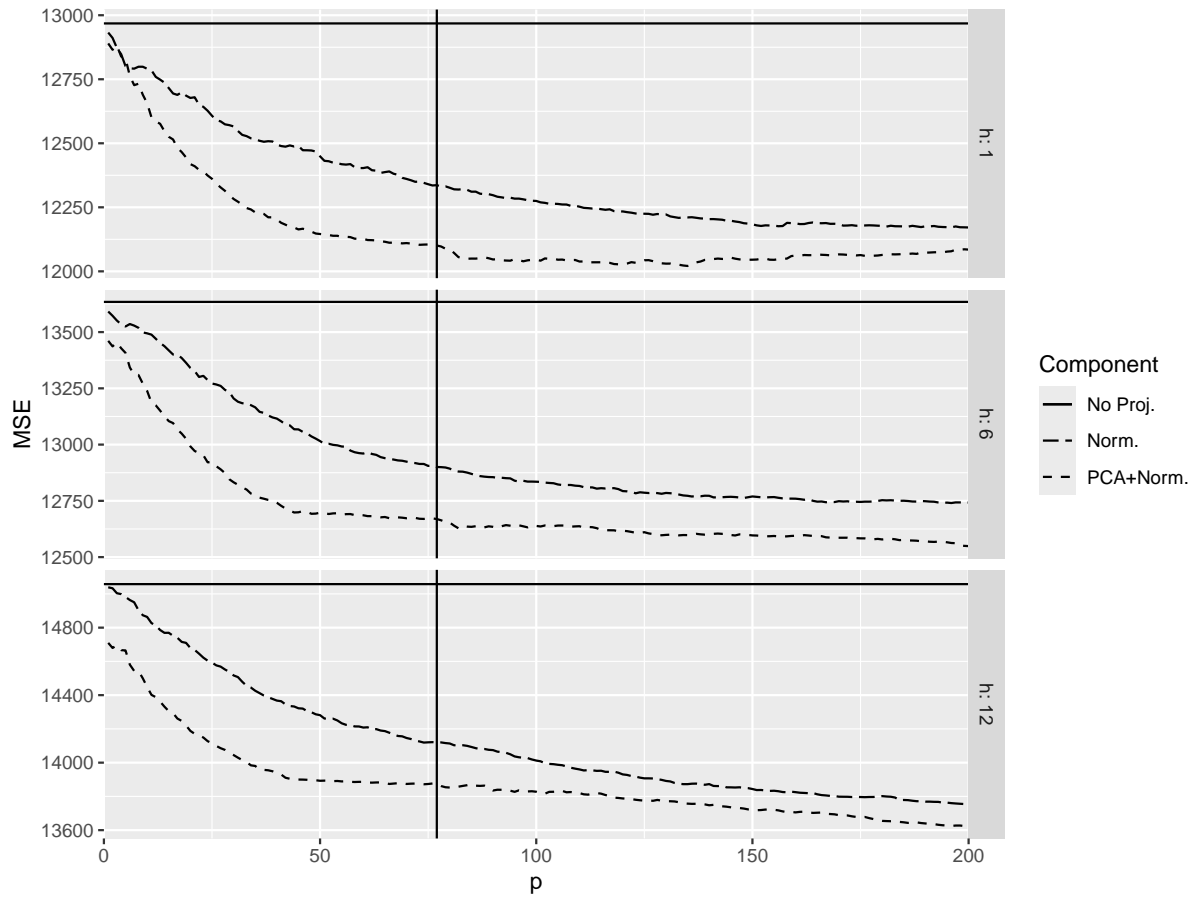


Figure 4: Out-of-sample MSE of base and FLAP forecasts as the number of components p increases, for forecast horizons 1, 6 and 12, using the visitor nights data. The solid horizontal lines show the MSE for the base forecasts while the dashed lines show the MSEs of the FLAP forecasts. The vertical black line indicates the location of $p = m$, the number of series.

even though they have very different numbers of components ($p = 77$ and $p = 200$). Choosing the number of components is a tradeoff between the increasing estimation error as the dimension of forecast error variance \mathbf{W}_h increases, and the additional benefit brought by the information embedded in the new components, depending on the complexity of the DGP. For the visitor nights data set, the benefit of components above the estimation error diminishes after the number of components reaches $p = m = 77$.

Second, unlike Section 3, using principal components is significantly better than simply using random components with normal weights. In Figure 3, this is observed by comparing the performance of PCA related FLAP and the performance of FLAP with only random components, with a relatively large number of components, for example, $m = 77$ or 300. This is also clear from Figure 4, where the reduction in terms of sample MSE from using PCA is always in the lead, even after the m principal components are exhausted and random components with weights generated from a normal distribution are added.

PCA aims to find components along the direction where the data varies the most. The first principal components accounts for the maximum variance in the data. Subsequent principal components are orthogonal to the previous ones and capture the maximum remaining variance. Based on this variance maximisation and ranking of PCA, we propose two potential explanations for its superior performance.

1. **Optimality:** By maximising and ranking the variance of PCs from largest to smallest, we ensure that the projection utilises components containing significantly more information (as measured by variance) compared to randomly weighted components.
2. **Diversity:** Actively seeking PCs with the highest variance results in the incorporation of a more diverse set of components into the projection.

4.2 FRED-MD

The FRED-MD (McCracken & Ng 2016) is a popular monthly data set of macroeconomic variables, and shares similar properties with the Stock & Watson (2002a) data. We downloaded and transformed the data set using the `fbi` package (Chen, Ng & Bai 2023). The period we use for this exercise is from January 1959 to September 2023, containing 777 observations. Following McCracken & Ng (2016), we replace observations that deviate from the sample median by more than 10 interquartile ranges (which are recognised as outliers), with missing values. We then drop any series with more than 5% observations missing. This left us with $m = 122$ series. We fill in the missing values using the expectation-maximization (EM) algorithm described in Stock & Watson (2002b) with 8 factors. The number 8 is identified by McCracken & Ng (2016), albeit with a different time span. In order to relate to the theoretical forecast error variance reduction, we use MSE as the error measure, instead of other scaled or percentage error measures. To reliably calculate MSE over series with different scales, we demean the series and scale them to have variance 1. The MSEs are calculated on this standardised scale without back-transformation.

We evaluate the performance of forecasts using time series cross-validation. Starting with 300 observations in the first training set and the following 12 observations as the test set, we repeat the evaluation for the rest of the data with the size of the training set increasing by 1 in each iteration. This provides us with 466 forecasts for each of the forecast horizons from 1 to 12 for evaluation. We generate base forecasts from ARIMA models using `auto.arima()` function from the `forecast` package (Hyndman et al. 2023) with the default settings, and DFMs. The ranges of the meta-parameters in the DFMs are $1 \leq k \leq 8$ (since 8 factors are identified and used to fill in the missing values), $1 \leq n \leq 3$ and $0 \leq s \leq 6$. For more details see Section 3.1. The MCB and MSE plots for $h = 1, 6, 12$ are shown in Figure 5 and Figure 6.

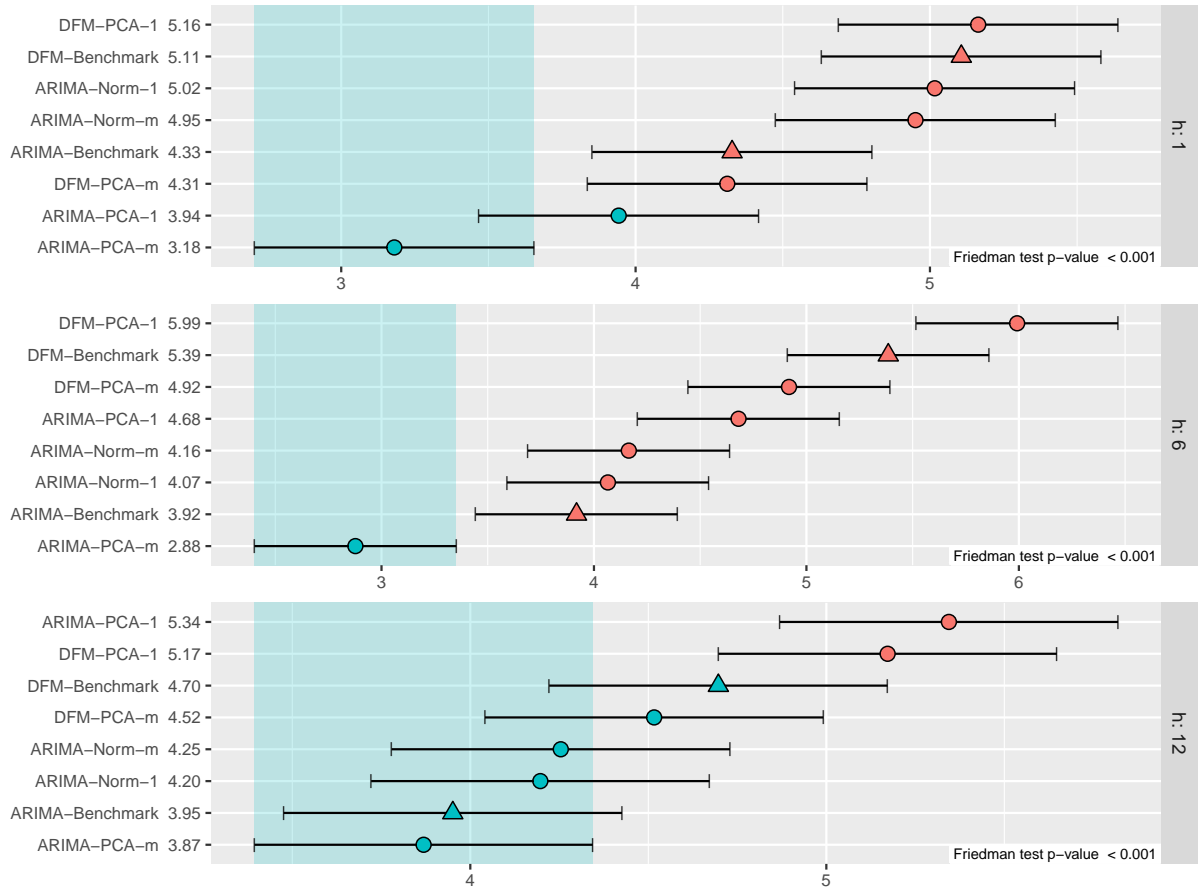


Figure 5: Average ranks of 1-, 6- and 12-step-ahead cross-validation MSE of different model and component specifications on the FRED-MD data. The methods using forecast projection are named as “Model – Component Weights – Number of Components”. The base models are named as “Model – Benchmark” and these points are marked with triangles. The shaded region is the confidence interval of the best performing model. Methods outside the shaded region are significantly worse than the best model.

The first observation worth noting in Figure 5 is the performance of base ARIMA forecasts exceeds that of the base DFM. This difference is statistically significant for $h = 6$. The best models at all forecast horizons are once again FLAP forecasts using components from PCA. These are significantly better than the base forecasts at $h = 1$ and 6. The fact that the FLAP with PCA is statistically significantly better than with components with random weights, which can be seen from both the rankings in Figure 5 and the MSE in Figure 6, reaffirms our findings from the tourism data set discussed in Section 4.1.

In Figure 6, FLAP forecasts for both ARIMA and DFM seem to be worse than the base forecasts at the beginning, when the number of components included in FLAP is small. However, this improves as p becomes larger and FLAP forecasts outperform the base forecasts gradually. The tradeoff between the benefit of including more components and the difficulty of estimation in a large dimension is once again observed, however seems to be more extreme compared to the tourism example. In this

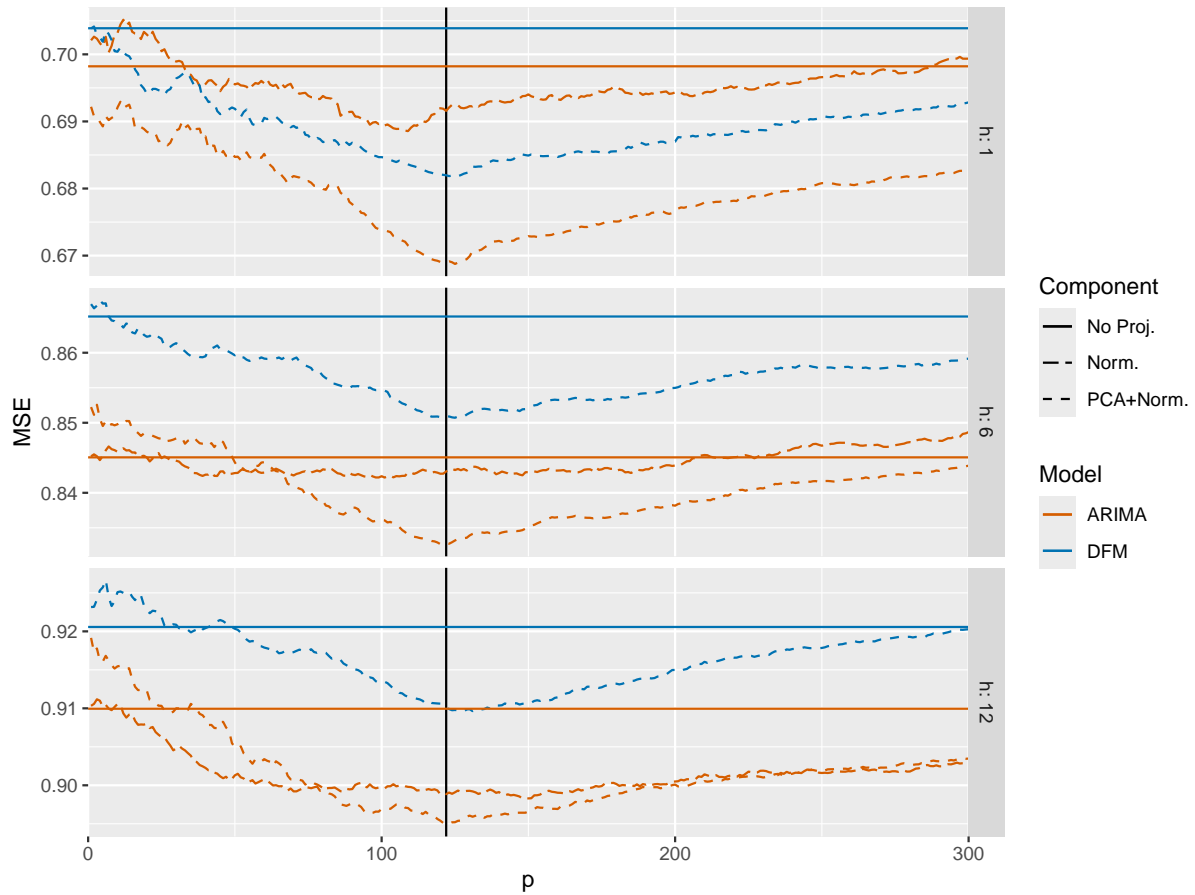


Figure 6: Out-of-sample MSE for base and FLAP forecasts as the number of components p increases, for forecast horizons 1, 6 and 12, using the FRED-MD data. The solid horizontal lines show the MSE for the base forecasts while the dashed lines show the MSEs of the FLAP forecasts. The vertical black line indicates the location of $p = m$ the number of series.

case, the MSEs start to increase once p becomes larger than m . As we have seen in Section 3 and Section 4.1, m is not always a clear cut-off point. Where the performance of FLAP turns should be jointly determined by the number of series m , the sample size T , the component construction method, and the DGP. In the case of FRED-MD, it signals the importance of PCA, as m is the point that the component changes from PCA to random normal weighted linear combinations, implying PCA can exploit the information in the data while random weights cannot. This is more obvious for $h = 1$ and $h = 6$, as PCA works when $p < m$, but random normal weights do not seem to work from the beginning.

5 Conclusion

The proposed forecast linear augmented projection (FLAP) method has been shown to be a simple but effective way to reduce forecast error variance of any multivariate forecasting problem. It simply involves augmenting the data with linear combinations, forecasting these and then projecting the

augmented vector of forecasts. We have shown theoretically that FLAP will continue to reduce forecast error covariance as more components are added, assuming we that the forecast error covariance matrix is known. In practice, a plug in estimate of this covariance matrix can be used, and in both simulated and empirical data we demonstrate that a simple shrinkage estimator does indeed lead to improvements in forecast accuracy. Regarding the construction of components we find that PCA in practice achieves significant improvements in forecast accuracy. Another appealing property of FLAP is that the projection step can even compensate for a poor choice of base forecasting model. This is particularly attractive since it makes FLAP robust against model misspecification in the base forecasting step.

One outstanding issue is to find alternatives to PCA to select component weights. For example, Goerg (2013) proposed “forecastable components” that are optimal in the sense of minimising the forecast error variance of the components, while Matteson & Tsay (2011) proposed “dynamic orthogonal components” that reduce a multivariate time series to a set of uncorrelated univariate time series. It would be interesting to explore whether these components (or other similar suggestions) can be used effectively in FLAP. Also, another route to improving FLAP may be found by optimizing Equation 10 over Φ and G rather than just G .

Finally, while FLAP is motivated by the forecast reconciliation literature, the focus is very much on multivariate time series with no constraints. However, it would be possible to use both forecast reconciliation and forecast projection together. This may be particularly useful when there are relationships between series that are not captured in the known hierarchical structure.

Acknowledgements

We thank Daniele Girolimetto for contributing to the initial proof of Theorem 2.2 in his unpublished work.

References

- Ando, S & F Narita (2022). “An alternative proof of minimum trace reconciliation”. <https://doi.org/10.2139/ssrn.4199143>.
- Assimakopoulos, V & K Nikolopoulos (2000). The theta model: a decomposition approach to forecasting. *International Journal of Forecasting* **16** (4), 521–530.
- Athanasopoulos, G, RJ Hyndman, N Kourentzes & A Panagiotelis (2023). Forecast reconciliation: A review. *International Journal of Forecasting*.
- Athanasopoulos, G, RJ Hyndman, N Kourentzes & F Petropoulos (2017). Forecasting with temporal hierarchies. *European Journal of Operational Research* **262** (1), 60–74.
- Batchelor, R & P Dua (1995). Forecaster diversity and the benefits of combining forecasts. *Management Science* **41** (1), 68–75.
- Bergmeir, C, RJ Hyndman & JM Benítez (2016). Bagging exponential smoothing methods using STL decomposition and Box–Cox transformation. *International Journal of Forecasting* **32** (2), 303–312.
- Bernanke, BS, J Boivin & P Elias (2005). Measuring the effects of monetary policy: a factor-augmented vector autoregressive (FAVAR) approach. *The Quarterly journal of economics* **120**(1), 387–422.
- Borchers, HW (2023). *pracma: Practical Numerical Math Functions*. R package version 2.4.4. <https://CRAN.R-project.org/package=pracma>.
- Breiman, L (1996). Bagging predictors. *Machine learning* **24** (2), 123–140.
- Carriero, A, AB Galvao & G Kapetanios (2019). A comprehensive evaluation of macroeconomic forecasting methods. *International Journal of Forecasting* **35**(4), 1226–1239.
- Chen, Y, S Ng & J Bai (2023). *fbi: Factor-Based Imputation and FRED-MD/QD Data Set*. R package version 0.7.0. <https://github.com/cykbennie/fbi>.
- De Stefani, J, YA Le Borgne, O Caelen, D Hattab & G Bontempi (2019). Batch and incremental dynamic factor machine learning for multivariate and multi-step-ahead forecasting. *International Journal of Data Science and Analytics* **7** (4), 311–329.
- Di Fonzo, T & D Girolimetto (2023). Cross-temporal forecast reconciliation: Optimal combination method and heuristic alternatives. *International Journal of Forecasting* **39** (1), 39–57.
- Disney, SM & F Petropoulos (2015). Forecast combinations using multiple starting points. In: Logistics and Operations Management Section Annual Conference (Cardiff, Jan. 9, 2015).
- Fabio Di Narzo, A, JL Aznarte & M Stigler (2009). *tsDyn: Time series analysis based on dynamical systems theory*. R package version 0.7. <https://CRAN.R-project.org/package=tsDyn>.
- Friedman, M (1937). The Use of Ranks to Avoid the Assumption of Normality Implicit in the Analysis of Variance. *Journal of the American Statistical Association* **32** (200), 675–701.

- Friedman, M (1939). A correction: The use of ranks to avoid the assumption of normality implicit in the analysis of variance. *Journal of the American Statistical Association*.
- Goerg, G (2013). Forecastable Component Analysis. In: *Proceedings of The 30th International Conference on Machine Learning*, pp.64–72. <http://jmlr.org/proceedings/papers/v28/goerg13.pdf>.
- Hastie, T, R Tibshirani & J Friedman (2003). *The elements of statistical learning: Data mining, inference, and prediction*. 1st ed. Springer series in statistics. New York, NY: Springer. 536 pp.
- Hollander, M, DA Wolfe & E Chicken (2013). *Nonparametric Statistical Methods*. John Wiley & Sons.
- Hollyman, R, F Petropoulos & ME Tipping (2021). Understanding forecast reconciliation. *European Journal of Operational Research* **294** (1), 149–160.
- Hyndman, R, G Athanasopoulos, C Bergmeir, G Caceres, L Chhay, M O'Hara-Wild, F Petropoulos, S Razbash, E Wang & F Yasmeeen (2023). *forecast: Forecasting functions for time series and linear models*. R package version 8.21.1. <https://pkg.robjhyndman.com/forecast/>.
- Hyndman, RJ & G Athanasopoulos (2018). *Forecasting: principles and practice*. 2nd ed. Melbourne, Australia: OTexts. <http://OTexts.org/fpp2>.
- Hyndman, RJ & G Athanasopoulos (2021). *Forecasting: principles and practice*. 3rd ed. Melbourne, Australia: OTexts. <http://OTexts.org/fpp3>.
- Hyndman, RJ & Y Khandakar (2008). Automatic time series forecasting: The forecast package for R. *Journal of Statistical Software* **27** (3), 1–22.
- Jolliffe, IT (2002). *Principal Component Analysis*. Springer, New York, NY.
- Kang, Y, W Cao, F Petropoulos & F Li (2022). Forecast with forecasts: Diversity matters. *European Journal of Operational Research* **301** (1), 180–190.
- Koning, AJ, PH Franses, M Hibon & HO Stekler (2005). The M3 competition: Statistical tests of the results. *International Journal of Forecasting* **21** (3), 397–409.
- Kourentzes, N (2023). *tsutils: Time Series Exploration, Modelling and Forecasting*. R package version 0.9.4. <https://CRAN.R-project.org/package=tsutils>.
- Kourentzes, N, F Petropoulos & JR Trapero (2014). Improving forecasting by estimating time series structural components across multiple frequencies. *International Journal of Forecasting* **30** (2), 291–302.
- Li, X, F Petropoulos & Y Kang (2022). Improving forecasting by subsampling seasonal time series. *International Journal of Production Research* **61** (3), 1–17.
- Lichtendahl Jr, KC & RL Winkler (2020). Why do some combinations perform better than others? *International Journal of Forecasting* **36** (1), 142–149.
- Luenberger, DG (1969). *Optimization by vector space methods*. Nashville, TN: John Wiley & Sons.

- Matteson, DS & RS Tsay (2011). Dynamic Orthogonal Components for Multivariate Time Series. *Journal of the American Statistical Association* **106**(496), 1450–1463. <http://dx.doi.org/10.1198/jasa.2011.tm10616>.
- McCracken, MW & S Ng (2016). FRED-MD: A Monthly Database for Macroeconomic Research. *Journal of Business & Economic Statistics* **34** (4), 574–589.
- Nemenyi, PB (1963). “Distribution-free multiple comparisons”. PhD thesis. Princeton University.
- Opgen-Rhein, R & K Strimmer (2007). Accurate ranking of differentially expressed genes by a distribution-free shrinkage approach. *Statistical Applications in Genetics and Molecular Biology* **6** (1).
- Panagiotelis, A, G Athanasopoulos, P Gamakumara & RJ Hyndman (2021). Forecast reconciliation: A geometric view with new insights on bias correction. *International Journal of Forecasting* **37** (1), 343–359.
- Petropoulos, F, RJ Hyndman & C Bergmeir (2018). Exploring the sources of uncertainty: Why does bagging for time series forecasting work? *European Journal of Operational Research* **268** (2), 545–554.
- Petropoulos, F & E Spiliotis (2021). The wisdom of the data: Getting the most out of univariate time series forecasting. *Forecasting* **3** (3), 478–497.
- R Core Team (2023). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. Vienna, Austria. <https://www.R-project.org/>.
- Rao, CR (1974). Projectors, generalized inverses and the Blue’s. *Journal of the Royal Statistical Society* **36** (3), 442–448.
- Schafer, J, R Opgen-Rhein, V Zuber, M Ahdesmaki, APD Silva & K Strimmer. (2021). *corpcor: Efficient Estimation of Covariance and (Partial) Correlation*. R package version 1.6.10. <https://CRAN.R-project.org/package=corpcor>.
- Schäfer, J & K Strimmer (2005). A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics. *Statistical Applications in Genetics and Molecular Biology* **4** (1).
- Stock, JH & MW Watson (2002a). Forecasting Using Principal Components From a Large Number of Predictors. *Journal of the American Statistical Association* **97** (460), 1167–1179.
- Stock, JH & MW Watson (2002b). Macroeconomic Forecasting Using Diffusion Indexes. *Journal of Business & Economic Statistics* **20** (2), 147–162.
- Stock, JH & MW Watson (2012). Generalized shrinkage methods for forecasting using many predictors. *Journal of Business & Economic Statistics* **30** (4), 481–493.
- Tsay, RS (2013). *Multivariate time series analysis: with R and financial applications*. John Wiley & Sons.

- Wang, X, RJ Hyndman, F Li & Y Kang (2023). Forecast combinations: An over 50-year review. *International Journal of Forecasting* **39** (4), 1518–1547.
- Wickramasuriya, SL, G Athanasopoulos & RJ Hyndman (2019). Optimal forecast reconciliation for hierarchical and grouped time series through trace minimization. *Journal of the American Statistical Association* **114** (526), 804–819.
- Yang, YF (2024). *flap: Forecast Linear Augmented Projection*. R package version 0.1.0. <https://CRAN.R-project.org/package=flap>.

A Proofs for Section 2 (Forecast Linear Augmented Projection (FLAP))

Proof of Lemma 2.1

We have

$$\begin{aligned}
 MM &= I_{m+p} - 2W_h C' (C W_h C')^{-1} C \\
 &\quad + W_h C' (C W_h C')^{-1} C W_h C' (C W_h C')^{-1} C \\
 &= I_{m+p} - W_h C' (C W_h C')^{-1} C \\
 &= M,
 \end{aligned}$$

so M is a projection matrix. For any \mathbf{z} such that $M\mathbf{z} = \mathbf{y}$ for some \mathbf{y} , we have

$$C\mathbf{y} = CM\mathbf{z} = C\mathbf{z} - C W_h C' (C W_h C')^{-1} C\mathbf{z} = \mathbf{0}.$$

Thus, M projects any vector onto the space where the constraint $C\mathbf{y} = \mathbf{0}$ is satisfied.

Proof of Corollary 2.1

Items 1 and 2 are trivial application of Lemma 2.1. To prove 3, we have

$$E(\tilde{\mathbf{z}}_{t+h} | \mathcal{I}_t) = E(M \hat{\mathbf{z}}_{t+h} | \mathcal{I}_t) = M E(\hat{\mathbf{z}}_{t+h} | \mathcal{I}_t) = M E(\mathbf{z}_{t+h} | \mathcal{I}_t) = E(M \mathbf{z}_{t+h} | \mathcal{I}_t) = E(\mathbf{z}_{t+h} | \mathcal{I}_t).$$

Proof of Lemma 2.2

$$\text{Var}(\tilde{\mathbf{z}}_{t+h} - \mathbf{z}_{t+h}) = \text{Var}(M \hat{\mathbf{z}}_{t+h} - M \mathbf{z}_{t+h}) = M \text{Var}(\hat{\mathbf{z}}_{t+h} - \mathbf{z}_{t+h}) M' = M W_h M'.$$

If we simplify it further, we have

$$\begin{aligned}
 M W_h M' &= (I - W_h C' (C W_h C')^{-1} C) W_h (I - W_h C' (C W_h C')^{-1} C)' \\
 &= W_h - W_h C' (C W_h C')^{-1} C W_h - W_h C' (C W_h C')^{-1} C W_h \\
 &\quad + W_h C' (C W_h C')^{-1} C W_h C' (C W_h C')^{-1} C W_h \\
 &= W_h - W_h C' (C W_h C')^{-1} C W_h \\
 &= M W_h.
 \end{aligned}$$

To get $\text{Var}(\tilde{\mathbf{y}}_{t+h} - \mathbf{y}_{t+h})$, we just need to recognise that it is the first $m \times m$ leading principal submatrix of $\text{Var}(\tilde{\mathbf{z}}_{t+h} - \mathbf{z}_{t+h})$.

Proof of Theorem 2.1

Trivially, $W_h C' (C W_h C')^{-1} C W_h$ and $J W_h C' (C W_h C')^{-1} C W_h J'$ are positive semi-definite. Note that $\text{Var}(\hat{y}_{t+h} - z_{t+h}) - \text{Var}(\tilde{y}_{t+h} - z_{t+h})$ is the leading principal submatrix of $W_h C' (C W_h C')^{-1} C W_h$, and the leading principal submatrix of a positive semi-definite matrix is positive semi-definite.

Proof of Theorem 2.2

Suppose now that we want to include q more components $c_t^* = \Phi^* y_t$ in the projection. We define

$z_t^* = \begin{bmatrix} z_t \\ c_t^* \end{bmatrix}$, the constraint matrix

$$C^* = \begin{bmatrix} C & \mathbf{0}_{p \times q} \\ -\Phi^* & \mathbf{0}_{q \times p} \end{bmatrix} = \begin{bmatrix} -\Phi & I_p & \mathbf{0}_{p \times q} \\ -\Phi^* & \mathbf{0}_{q \times p} & I_q \end{bmatrix} = \begin{bmatrix} \overline{C} \\ \underline{C} \end{bmatrix} \quad (12)$$

where \overline{C} contains the first p rows of C^* and \underline{C} contains the remaining q rows of C^* , the forecast error variance matrix

$$\text{Var}(\hat{z}_{t+h}^* - z_{t+h}^*) = W_h^* = \begin{bmatrix} W_h & W_{yc,h}^* \\ W_{cy,h}^* & W_{c,h}^* \end{bmatrix}.$$

where \hat{z}_{t+h}^* is the h -step-ahead base forecasts of z_t^* :

$$\hat{z}_{t+h}^* = \begin{bmatrix} \hat{z}_{t+h} \\ \hat{c}_{t+h}^* \end{bmatrix},$$

and the corresponding

$$M^* = I - W_h^* C^{*'} (C^* W_h^* C^{*'})^{-1} C^*.$$

Proving Theorem 2.2 requires proving the following two items.

1. Including additional components in the mapping without including corresponding component constraints is equivalent to not including these additional components at all.
2. For a fixed set of components to be included in the mapping, adding constraints will reduce forecast error variance.

We start by proving the first statement. Consider the case where we include the additional series c_t^* without using the additional constraint Φ^* . Defining M^+ only with \overline{C} :

$$M^+ = I_{m+p+q} - W_h^* \overline{C}' (\overline{C} W_h^* \overline{C}')^{-1} \overline{C}, \quad (13)$$

we have $\tilde{\mathbf{z}}_{t+h}^+ = \mathbf{M}^+ \hat{\mathbf{z}}_{t+h}^*$. Further, we obtain

$$\mathbf{W}_h^* \bar{\mathbf{C}}' = \begin{bmatrix} \mathbf{W}_h & \mathbf{W}_{yc,h}^* \\ \mathbf{W}_{cy,h}^* & \mathbf{W}_{c,h}^* \end{bmatrix} \begin{bmatrix} \mathbf{C}' \\ \mathbf{0}_{q \times p} \end{bmatrix} = \begin{bmatrix} \mathbf{W}_h \mathbf{C}' \\ \mathbf{W}_{cy,h}^* \mathbf{C}' \end{bmatrix}$$

and

$$\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}' = \begin{bmatrix} \mathbf{C} & \mathbf{0}_{p \times q} \end{bmatrix} \begin{bmatrix} \mathbf{W}_h \mathbf{C}' \\ \mathbf{W}_{cy,h}^* \mathbf{C}' \end{bmatrix} = \mathbf{C} \mathbf{W}_h \mathbf{C}',$$

which gives

$$\begin{aligned} \mathbf{M}^+ &= \mathbf{I}_{m+p+q} - \begin{bmatrix} \mathbf{W}_h \mathbf{C}' \\ \mathbf{W}_{cy,h}^* \mathbf{C}' \end{bmatrix} (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \begin{bmatrix} \mathbf{C} & \mathbf{0}_{p \times q} \end{bmatrix} \\ &= \mathbf{I}_{m+p+q} - \begin{bmatrix} \mathbf{W}_h \mathbf{C}' (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \mathbf{C} & \mathbf{0} \\ \mathbf{W}_{cy,h}^* \mathbf{C}' (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \mathbf{C} & \mathbf{0} \end{bmatrix}, \end{aligned}$$

and

$$\begin{aligned} \tilde{\mathbf{z}}_{t+h}^+ &= \mathbf{M}^+ \hat{\mathbf{z}}_{t+h}^* \\ &= \left(\mathbf{I}_{m+p+q} - \begin{bmatrix} \mathbf{W}_h \mathbf{C}' (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \mathbf{C} & \mathbf{0} \\ \mathbf{W}_{cy,h}^* \mathbf{C}' (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \mathbf{C} & \mathbf{0} \end{bmatrix} \right) \begin{bmatrix} \hat{\mathbf{z}}_{t+h} \\ \hat{\mathbf{c}}_{t+h}^* \end{bmatrix} \\ &= \begin{bmatrix} (\mathbf{I}_{n+p} - \mathbf{W}_h \mathbf{C}' (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \mathbf{C}) \hat{\mathbf{z}}_{t+h} \\ \hat{\mathbf{c}}_{t+h}^* - \mathbf{W}_{cy,h}^* \mathbf{C}' (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \mathbf{C} \hat{\mathbf{z}}_{t+h} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{M} \hat{\mathbf{z}}_{t+h} \\ \hat{\mathbf{c}}_{t+h}^* - \mathbf{W}_{cy,h}^* \mathbf{C}' (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \mathbf{C} \hat{\mathbf{z}}_{t+h} \end{bmatrix} \\ &= \begin{bmatrix} \tilde{\mathbf{z}}_{t+h} \\ \hat{\mathbf{c}}_{t+h}^* - \mathbf{W}_{cy,h}^* \mathbf{C}' (\mathbf{C} \mathbf{W}_h \mathbf{C}')^{-1} \mathbf{C} \hat{\mathbf{z}}_{t+h} \end{bmatrix}. \end{aligned}$$

If we only consider the forecast performance relevant to \mathbf{y}_{t+h} , and define $\mathbf{J}^* = \mathbf{J}_{m,p+q} = \begin{bmatrix} \mathbf{I}_m & \mathbf{0}_{m \times (p+q)} \end{bmatrix}$, we have

$$\tilde{\mathbf{y}}_{t+h}^+ = \mathbf{J}^* \tilde{\mathbf{z}}_{t+h}^+ = \mathbf{J} \tilde{\mathbf{z}}_{t+h} = \tilde{\mathbf{y}}_{t+h}.$$

This means adding additional components without imposing the corresponding constraints will yield the same projected forecasts as if these additional components are not added, which implies that the forecast error variance stays the same:

$$\text{Var}(\tilde{\mathbf{y}}_{t+h}^+ - \mathbf{y}_{t+h}) = \text{Var}(\tilde{\mathbf{y}}_{t+h} - \mathbf{y}_{t+h}) = \mathbf{J} \mathbf{M} \mathbf{W}_h \mathbf{J}'. \quad (14)$$

This finishes the proof of the first statement. Now we move on to proving the second statement. We have the forecast error variance matrices

$$\begin{aligned} \text{Var}(\tilde{\mathbf{z}}_{t+h}^+ - \mathbf{z}_{t+h}^*) &= \mathbf{M}^+ \mathbf{W}_h^* = (\mathbf{I}_{m+p+q} - \mathbf{W}_h^* \bar{\mathbf{C}}' (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} \bar{\mathbf{C}}) \mathbf{W}_h^* \\ \text{and} \quad \text{Var}(\tilde{\mathbf{z}}_{t+h}^* - \mathbf{z}_{t+h}^*) &= \mathbf{M}^* \mathbf{W}_h^* = (\mathbf{I}_{m+p+q} - \mathbf{W}_h^* \mathbf{C}^{*'} (\mathbf{C}^* \mathbf{W}_h^* \mathbf{C}^{*'})^{-1} \mathbf{C}^*) \mathbf{W}_h^*. \end{aligned}$$

Taking the difference, we have

$$\begin{aligned} \text{Var}(\tilde{\mathbf{z}}_{t+h}^+ - \mathbf{z}_{t+h}^*) - \text{Var}(\tilde{\mathbf{z}}_{t+h}^* - \mathbf{z}_{t+h}^*) &= (\mathbf{W}_h^* \mathbf{C}^{*'} (\mathbf{C}^* \mathbf{W}_h^* \mathbf{C}^{*'})^{-1} \mathbf{C}^* - \mathbf{W}_h^* \bar{\mathbf{C}}' (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} \bar{\mathbf{C}}) \mathbf{W}_h^* \\ &= \mathbf{W}_h^* (\mathbf{C}^{*'} (\mathbf{C}^* \mathbf{W}_h^* \mathbf{C}^{*'})^{-1} \mathbf{C}^* - \bar{\mathbf{C}}' (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} \bar{\mathbf{C}}) \mathbf{W}_h^*. \end{aligned}$$

Using block matrix inversion, we have

$$\begin{aligned} \mathbf{C}^{*'} (\mathbf{C}^* \mathbf{W}_h^* \mathbf{C}^{*'})^{-1} \mathbf{C}^* &= \begin{bmatrix} \bar{\mathbf{C}}' & \underline{\mathbf{C}}' \end{bmatrix} \begin{bmatrix} \bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}' & \bar{\mathbf{C}} \mathbf{W}_h^* \underline{\mathbf{C}}' \\ \underline{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}' & \underline{\mathbf{C}} \mathbf{W}_h^* \underline{\mathbf{C}}' \end{bmatrix}^{-1} \begin{bmatrix} \bar{\mathbf{C}} \\ \underline{\mathbf{C}} \end{bmatrix} \\ &= \begin{bmatrix} \bar{\mathbf{C}}' & \underline{\mathbf{C}}' \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} \bar{\mathbf{C}} \\ \underline{\mathbf{C}} \end{bmatrix} \\ &= \bar{\mathbf{C}}' a \bar{\mathbf{C}} + \bar{\mathbf{C}}' b \underline{\mathbf{C}} + \underline{\mathbf{C}}' c \bar{\mathbf{C}} + \underline{\mathbf{C}}' d \underline{\mathbf{C}}, \end{aligned}$$

where

$$\begin{aligned} a &= (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} + (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} \bar{\mathbf{C}} \mathbf{W}_h^* \underline{\mathbf{C}}' \\ &\quad (\underline{\mathbf{C}} \mathbf{W}_h^* \underline{\mathbf{C}}' - \underline{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}' (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} \bar{\mathbf{C}} \mathbf{W}_h^* \underline{\mathbf{C}}')^{-1} \underline{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}' (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} \\ &= (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} + (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} \bar{\mathbf{C}} \mathbf{W}_h^* \underline{\mathbf{C}}' (\underline{\mathbf{C}} \mathbf{M}^+ \mathbf{W}_h^* \underline{\mathbf{C}}')^{-1} \underline{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}' (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1}, \\ b &= -(\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} \bar{\mathbf{C}} \mathbf{W}_h^* \underline{\mathbf{C}}' (\underline{\mathbf{C}} \mathbf{W}_h^* \underline{\mathbf{C}}' - \underline{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}' (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} \bar{\mathbf{C}} \mathbf{W}_h^* \underline{\mathbf{C}}')^{-1} \\ &= -(\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} \bar{\mathbf{C}} \mathbf{W}_h^* \underline{\mathbf{C}}' (\underline{\mathbf{C}} \mathbf{M}^+ \mathbf{W}_h^* \underline{\mathbf{C}}')^{-1}, \\ c &= -(\underline{\mathbf{C}} \mathbf{M}^+ \mathbf{W}_h^* \underline{\mathbf{C}}')^{-1} \underline{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}' (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1}, \\ d &= (\underline{\mathbf{C}} \mathbf{M}^+ \mathbf{W}_h^* \underline{\mathbf{C}}')^{-1}. \end{aligned}$$

Thus,

$$\begin{aligned} \mathbf{C}^{*'} (\mathbf{C}^* \mathbf{W}_h^* \mathbf{C}^{*'})^{-1} \mathbf{C}^* &= \bar{\mathbf{C}}' (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} \bar{\mathbf{C}} \\ &\quad + \bar{\mathbf{C}}' (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} \bar{\mathbf{C}} \mathbf{W}_h^* \underline{\mathbf{C}}' (\underline{\mathbf{C}} \mathbf{M}^+ \mathbf{W}_h^* \underline{\mathbf{C}}')^{-1} \underline{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}' (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} \bar{\mathbf{C}} \\ &\quad - \bar{\mathbf{C}}' (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} \bar{\mathbf{C}} \mathbf{W}_h^* \underline{\mathbf{C}}' (\underline{\mathbf{C}} \mathbf{M}^+ \mathbf{W}_h^* \underline{\mathbf{C}}')^{-1} \underline{\mathbf{C}} \\ &\quad - \underline{\mathbf{C}}' (\underline{\mathbf{C}} \mathbf{M}^+ \mathbf{W}_h^* \underline{\mathbf{C}}')^{-1} \underline{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}' (\bar{\mathbf{C}} \mathbf{W}_h^* \bar{\mathbf{C}}')^{-1} \bar{\mathbf{C}} \\ &\quad + \underline{\mathbf{C}}' (\underline{\mathbf{C}} \mathbf{M}^+ \mathbf{W}_h^* \underline{\mathbf{C}}')^{-1} \underline{\mathbf{C}} \end{aligned}$$

$$\begin{aligned}
 &= \bar{C}'(\bar{C}W_h^*\bar{C}')^{-1}\bar{C} \\
 &\quad - \bar{C}(\bar{C}W_h^*\bar{C}')^{-1}\bar{C}W_h^*\underline{C}'(\underline{C}M^+W_h^*\underline{C}')^{-1}\underline{C}M^+ \\
 &\quad + \underline{C}'(\underline{C}M^+W_h^*\underline{C}')^{-1}\underline{C}M^+ \\
 &= \bar{C}'(\bar{C}W_h^*\bar{C}')^{-1}\bar{C} + M^{+'}\underline{C}'(\underline{C}M^+W_h^*\underline{C}')^{-1}\underline{C}M^+.
 \end{aligned}$$

Therefore,

$$\begin{aligned}
 &\text{Var}(\tilde{z}_{t+h}^+ - z_{t+h}^*) - \text{Var}(\tilde{z}_{t+h}^* - z_{t+h}^*) \\
 &= W_h^*(\bar{C}'(\bar{C}W_h^*\bar{C}')^{-1}\bar{C} + M^{+'}\underline{C}'(\underline{C}M^+W_h^*\underline{C}')^{-1}\underline{C}M^+ - \bar{C}'(\bar{C}W_h^*\bar{C}')^{-1}\bar{C})W_h^* \\
 &= W_h^*(M^{+'}\underline{C}'(\underline{C}M^+W_h^*\underline{C}')^{-1}\underline{C}M^+)W_h^*
 \end{aligned}$$

is positive semi-definite. This concludes the proof of the second statement. Combining the results above, we have

$$\begin{aligned}
 \text{Var}(\tilde{y}_{t+h} - y_{t+h}) - \text{Var}(\tilde{y}_{t+h}^* - y_{t+h}) &= \text{Var}(\tilde{y}_{t+h}^+ - y_{t+h}) - \text{Var}(\tilde{y}_{t+h}^* - y_{t+h}) \\
 &= J^* \text{Var}(\tilde{z}_{t+h}^+ - z_{t+h}^*)J^{*'} - J^* \text{Var}(\tilde{z}_{t+h}^* - z_{t+h}^*)J^{*'} \quad (15) \\
 &= J^*W_h^*M^{+'}\underline{C}'(\underline{C}M^+W_h^*\underline{C}')^{-1}\underline{C}M^+W_h^*J^{*'}
 \end{aligned}$$

being positive semi-definite. Finally, we have

$$\begin{aligned}
 (\text{Var}(\hat{y}_{t+h} - y_{t+h}) - \text{Var}(\tilde{y}_{t+h}^* - y_{t+h})) - (\text{Var}(\hat{y}_{t+h} - y_{t+h}) - \text{Var}(\tilde{y}_{t+h} - y_{t+h})) \\
 = \text{Var}(\tilde{y}_{t+h} - y_{t+h}) - \text{Var}(\tilde{y}_{t+h}^* - y_{t+h})
 \end{aligned}$$

being a positive semi-definite matrix where the diagonal terms are non-negative, whose trace, therefore, is non-negative. This means using a larger number of components in the mapping achieves lower forecast error variances, giving Theorem 2.2.

Proof of Theorem 2.3

Denote $\psi_i = \begin{bmatrix} -\phi_i & \mathbf{0}_{1 \times (i-1)} & 1 \end{bmatrix}$ and $W_h^{(i)}$ to be the base forecast error variance of the original series and the first i components. Starting with the first component, Equation 4 becomes

$$\text{tr}(J_{m,1}W_h^{(1)}\psi_1'(\psi_1W_h^{(1)}\psi_1')^{-1}\psi_1W_h^{(1)}J_{m,1}') = (\psi_1W_h^{(1)}\psi_1')^{-1}\psi_1W_h^{(1)}J_{m,1}'J_{m,1}W_h^{(1)}\psi_1', \quad (16)$$

$$\begin{aligned}
 \text{where} \quad \psi_1W_h^{(1)}J_{m,1}' &= \begin{bmatrix} -\phi_1 & 1 \end{bmatrix} \begin{bmatrix} W_{z,h} & w_{c_1z,h}' \\ w_{c_1z,h} & W_{c_1,h} \end{bmatrix} \begin{bmatrix} I_m \\ 0 \end{bmatrix} \\
 &= -\phi_1 W_{z,h} + w_{c_1z,h}.
 \end{aligned}$$

Equation 16 is obviously non-negative. For it to be larger than 0, we need $\psi_1 \mathbf{W}_h^{(1)} \mathbf{J}_{m,1}' \neq 0$, which gives $\phi_1 \mathbf{W}_{z,h} \neq \mathbf{w}_{c_1 z, h}$.

When it comes to adding the i th component on top of the first $i - 1$ components, we define

$$\bar{\mathbf{C}}_i = \begin{bmatrix} \psi_1 & \mathbf{0}_{1 \times i} \\ \psi_2 & \mathbf{0}_{1 \times (i-1)} \\ \vdots & \vdots \\ \psi_i & 0 \end{bmatrix}$$

and

$$\mathbf{M}_i^+ = \mathbf{I}_{m+i} - \mathbf{W}_h^{(i)} \bar{\mathbf{C}}_{i-1}' (\bar{\mathbf{C}}_{i-1} \mathbf{W}_h^{(i)} \bar{\mathbf{C}}_{i-1}')^{-1} \bar{\mathbf{C}}_{i-1}$$

analogously to Equation 12 and Equation 13. Following Equation 15, the additional reduction of forecast error variance when adding the i th component becomes

$$\mathbf{J}_{m,i} \mathbf{W}_h^{(i)} \mathbf{M}_i^{+'} \psi_i' (\psi_i \mathbf{M}_i^+ \mathbf{W}_h^{(i)} \psi_i')^{-1} \psi_i \mathbf{M}_i^+ \mathbf{W}_h^{(i)} \mathbf{J}_{m,i}' = (\psi_i \mathbf{M}_i^+ \mathbf{W}_h^{(i)} \psi_i')^{-1} \psi_i \mathbf{M}_i^+ \mathbf{W}_h^{(i)} \mathbf{J}_{m,i}' \mathbf{J}_{m,i} \mathbf{W}_h^{(i)} \mathbf{M}_i^{+'} \psi_i'.$$

Similar to before, we would want $\psi_i \mathbf{M}_i^+ \mathbf{W}_h^{(i)} \mathbf{J}_{m,i}' \neq \mathbf{0}$. Note that ψ_i concerns the first m rows and the last row of $\mathbf{M}_i^+ \mathbf{W}_h^{(i)}$, and $\mathbf{J}_{m,i}'$ concerns the first m columns. Combined with the implication from Equation 14 that the $m \times m$ leading principal submatrix in equation $\mathbf{J}_{m,i} \mathbf{M}_i^+ \mathbf{W}_h^{(i)} \mathbf{J}_{m,i}' = \mathbf{J}_{m,i-1} \mathbf{M}_{i-1} \mathbf{W}_h^{(i-1)} \mathbf{J}_{m,i-1}'$ is the same, we suppress the straightforward yet tiresome details, and obtain

$$\phi_i \mathbf{W}_{z,h}^{(i-1)} \neq [\mathbf{0}_{1 \times m+i-1} \quad 1] \mathbf{M}_i^+ \mathbf{W}_h^{(i)} \mathbf{J}_{m,i}',$$

where $\mathbf{W}_{z,h}^{(i-1)} = \mathbf{J}_{m,i-1} \mathbf{M}_{i-1} \mathbf{W}_h^{(i-1)} \mathbf{J}_{m,i-1}'$ is the projected forecast error variance of the original series using the first $i - 1$ components, and the right hand side of the inequality is simply a one-row matrix consisting of the first m elements in the last row of $\mathbf{M}_i^+ \mathbf{W}_h^{(i)}$, which can be denoted as $\mathbf{w}_{\tilde{c}_i z, h}^{(i-1)}$ and interpreted as the covariance between the projected forecast of the original series using the first $i - 1$ components, and the projected forecast of the i th component using the first $i - 1$ components.

Proof of Lemma 2.3

If $\mathbf{G}\mathbf{S} = \mathbf{I}$, $\mathbf{S}\mathbf{G}$ is a projection matrix: $\mathbf{S}\mathbf{G}\mathbf{S}\mathbf{G} = \mathbf{S}\mathbf{G}$.

For any \mathbf{z} such that $\mathbf{S}\mathbf{G}\mathbf{z} = \mathbf{y}$ for some \mathbf{y} , we have $\mathbf{C}\mathbf{y} = \mathbf{C}\mathbf{S}\mathbf{G}\mathbf{z} = \mathbf{0}$ because $\mathbf{C}\mathbf{S} = [-\Phi \quad \mathbf{I}][\mathbf{I} \quad \Phi']' = \mathbf{0}$. Similarly to \mathbf{M} , $\mathbf{S}\mathbf{G}$ projects a vector to the same space where \mathbf{C} is satisfied.

Proof of Corollary 2.2

Item 1 is an direct application of Lemma 2.3. From Lemma 2.3 and Lemma 2.4 in Rao (1974), we have

$$SG\mathbf{z}_{t+h} = \mathbf{z}_{t+h} = S\mathbf{y}_{t+h}.$$

Left multiplying by G on both sides, we have $G\mathbf{z}_{t+h} = \mathbf{y}_{t+h}$ and item 2 is proven. To prove Item 3, we have

$$E(\tilde{\mathbf{y}}_{t+h}|\mathcal{I}_t) = E(G\hat{\mathbf{z}}_{t+h}|\mathcal{I}_t) = G E(\hat{\mathbf{z}}_{t+h}|\mathcal{I}_t) = G E(\mathbf{z}_{t+h}|\mathcal{I}_t) = E(G\mathbf{z}_{t+h}|\mathcal{I}_t) = E(\mathbf{y}_{t+h}|\mathcal{I}_t).$$

Proof of Lemma 2.4

Let the base and projected forecast errors be given as

$$\hat{\mathbf{e}}_{y,t+h} = \mathbf{y}_{t+h} - \hat{\mathbf{y}}_{t+h},$$

$$\hat{\mathbf{e}}_{z,t+h} = \mathbf{z}_{t+h} - \hat{\mathbf{z}}_{t+h},$$

$$\tilde{\mathbf{e}}_{y,t+h} = \mathbf{y}_{t+h} - \tilde{\mathbf{y}}_{t+h},$$

$$\text{and } \tilde{\mathbf{e}}_{z,t+h} = \mathbf{z}_{t+h} - \tilde{\mathbf{z}}_{t+h} = S\mathbf{y}_{t+h} - S\tilde{\mathbf{y}}_{t+h} = S\tilde{\mathbf{e}}_{y,t+h}.$$

$$\begin{aligned} \text{Then we have } \tilde{\mathbf{e}}_{z,t+h} &= \hat{\mathbf{e}}_{z,t+h} + \hat{\mathbf{z}}_{t+h} - \tilde{\mathbf{z}}_{t+h} \\ &= \hat{\mathbf{e}}_{z,t+h} + \hat{\mathbf{z}}_{t+h} - SG\hat{\mathbf{z}}_{t+h} \\ &= \hat{\mathbf{e}}_{z,t+h} + (I - SG)(\mathbf{z}_{t+h} - \hat{\mathbf{e}}_{z,t+h}) \end{aligned}$$

$$\text{and } = SG\hat{\mathbf{e}}_{z,t+h} + (I - SG)S\mathbf{y}_{t+h}$$

$$S\tilde{\mathbf{e}}_{y,t+h} = SG\hat{\mathbf{e}}_{z,t+h},$$

where the last line comes from $GS = I$. Left multiplying by G on both sides, we have

$$GS\tilde{\mathbf{e}}_{y,t+h} = GSG\hat{\mathbf{e}}_{z,t+h} \quad \text{and} \quad \tilde{\mathbf{e}}_{y,t+h} = G\hat{\mathbf{e}}_{z,t+h},$$

and therefore

$$\text{Var}(\tilde{\mathbf{y}}_{t+h} - \mathbf{y}_{t+h}) = \text{Var}(\tilde{\mathbf{e}}_{y,t+h}) = \text{Var}(G\hat{\mathbf{e}}_{z,t+h}) = G \text{Var}(\hat{\mathbf{e}}_{z,t+h})G' = GW_hG'.$$

Proof of Theorem 2.4

This can be proved in a few different ways. We adopt the approach of Ando & Narita (2022) to obtain the solution to Equation 10, but the procedure from Luenberger (1969, p. 85) can also be used, where the problem is divided to Equation 11 and reconstructed to find the solution to Equation 10.

There exists a Lagrange multiplier Λ such that

$$L(\mathbf{G}) = \text{tr}(\mathbf{G}\mathbf{W}_h\mathbf{G}') + \text{tr}(\Lambda'(I - \mathbf{G}\mathbf{S}))$$

is stationary at an extremum \mathbf{G} (Luenberger 1969, p. 243, Theorem 1). We set the Gateaux differential (Luenberger 1969, p. 171) to zero for any matrix \mathbf{H} :

$$\begin{aligned} \lim_{\alpha \rightarrow 0} \frac{L(\mathbf{G} + \alpha\mathbf{H}) - L(\mathbf{G})}{\alpha} &= 0 \\ \text{tr}(\mathbf{G}\mathbf{W}_h\mathbf{H}') + \text{tr}(\mathbf{H}\mathbf{W}_h\mathbf{G}') - \text{tr}(\Lambda'(\mathbf{H}\mathbf{S})) &= \text{tr}(2\mathbf{H}\mathbf{W}_h\mathbf{G}' - \Lambda'\mathbf{H}\mathbf{S}) \\ &= \text{tr}(\mathbf{H}(2\mathbf{W}_h\mathbf{G}' - \mathbf{S}\Lambda')) \\ &= 0 \\ 2\mathbf{W}_h\mathbf{G} &= \mathbf{S}\Lambda' \\ \mathbf{G}' &= \frac{1}{2}\mathbf{W}_h^{-1}\mathbf{S}\Lambda'. \end{aligned}$$

Multiplying \mathbf{S}' to the left of both sides. we have

$$\mathbf{S}'\mathbf{G}' = \mathbf{I} = \frac{1}{2}\mathbf{S}'\mathbf{W}_h^{-1}\mathbf{S}\Lambda' \quad \text{and} \quad \Lambda' = 2(\mathbf{S}'\mathbf{W}_h^{-1}\mathbf{S})^{-1}$$

because $\mathbf{G}\mathbf{S} = \mathbf{I}$. Putting it back in, we have

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