

To appear in *Quantitative Finance*, Vol. 00, No. 00, Month 20XX, 1–34

Improving out-of-sample forecasts of stock price indexes with forecast reconciliation and clustering

Raffaele Mattera[†], George Athanasopoulos^{*‡} and Rob Hyndman[‡]

[†]Department of Social and Economic Sciences, Sapienza University of Rome, Italy

[‡]Department of Econometrics and Business Statistics, Monash University, Australia

(Received 00 Month 20XX; in final form 00 Month 20XX)

In this paper we propose a novel approach to improving forecasts of stock market indexes by considering common stock prices as hierarchical time series, combining clustering with forecast reconciliation. We propose to group the individual stock price series in various ways including via metadata and using unsupervised learning techniques. The proposed approach is applied to the Dow Jones Industrial Average Index and the Standard & Poor 500 Index and their component stocks, and the results obtained with different grouping approaches are compared. The results demonstrate empirically that the combined use of clustering and reconciliation improves the forecast accuracy of the stock market indexes and their constituents.

Keywords: financial time series, hierarchical forecasting, clustering, unsupervised learning, prediction, machine learning, finance

JEL Classification: G17, C53, C58

1. Introduction

Accurate stock price forecasts are crucial in finance: they enable investors to make informed decisions, they help traders construct profitable investment strategies, and they assist policy-makers in monitoring the evolution of financial markets (e.g. see Marquering and Verbeek 2004, Kong et al. 2011, Allen et al. 2019). Chartists are traders looking for price patterns, and making investment decisions based on their price forecasts. The use of inaccurate forecasts by these traders represents a major driver of fluctuations in market sentiment, significantly contributing to mispricing in the stock market (Giamattei et al. 2020). A simple forecast-based trading strategy is, for example, to go long if the market index price is predicted to increase in the next h time periods, and to go short otherwise (Anatolyev and Gerko 2005, Blaskowitz and Herwartz 2011). This type of strategy is directly investable via deeply liquid derivatives markets, such as futures contracts on the market index (Alexander 2008). For example, Trippi and DeSieno (1992) built a trading strategy on the market index features using neural network-based forecasts and showed that this outperformed a passive investment strategy. Similar strategies are discussed in Sutcliffe (2018) and de Prado (2020).

Stock price forecasting is notoriously difficult, as financial markets are complex and turbulent systems. Indeed, stocks can be characterized by volatility clustering, non-linear relationships, long memory and hierarchical structures (e.g., see Tumminello et al. 2010). These difficulties, along with the efficient market hypothesis, have led many authors to conclude that it is impossible to

*Corresponding author. Email: George.Athanasopoulos@monash.edu

obtain better stock price forecasts than those from a random walk model (Fama 1995). Those that have proposed more complex approaches (Kumbure et al. 2022) have usually had limited empirical success. The success of random walk forecasts suggests that future values of the time series cannot be predicted with the currently available information only. For this reason, researchers and practitioners often use additional variables, known as factors or predictive signals, to aid in predicting prices and return (e.g. see Green et al. 2013). In this paper, we provide evidence that, if the hierarchical structure of the stock market is taken into consideration via forecast reconciliation, we can obtain more accurate forecasts for the market index than the random walk without reverting to sophisticated methods for forecasting the common stock prices. For this aim, approaches for forecasting hierarchical time series need to be adopted.

Forecast reconciliation is a statistical technique dealing with multivariate time series following a hierarchical structure (Hyndman et al. 2011), or that more generally adhere to linear constraints (Hyndman et al. 2016, Panagiotelis et al. 2021). Hierarchical time series are common in many application domains. For example, national level retail sales (Makridakis et al. 2022), tourism flows (Athanasopoulos et al. 2009), or electricity demand (Panagiotelis et al. 2023) can be disaggregated by state, regional or an even finer geographic grids. GDP is constructed using income or expense components (Bisaglia et al. 2020).

The motivation behind forecast reconciliation is that the forecasts for disaggregated and aggregated series are not necessarily coherent with the hierarchical structure. In other words, while time series data naturally aggregate according to a hierarchical structure, forecasts usually do not. This issue has been traditionally addressed by using forecasts at only one level of the hierarchy, from which the remaining forecasts are computed.

For example, in the bottom-up approach, forecasts are computed for the most disaggregated, bottom-level, series first. These are then aggregated to obtain forecasts for aggregate series (Dunn et al. 1976). However, the resulting forecasts of the aggregated series are often not as accurate as forecasting the aggregated series directly. The opposite approach is top-down forecasting, where forecasts of the most aggregated series are computed first, and these are then disaggregated down to obtain forecasts of series at lower levels of the hierarchy (e.g. see Gross and Sohl 1990). However, this often results in poor forecast accuracy at the disaggregated levels, and will always produce biased forecasts even if the original forecasts are unbiased (Hyndman et al. 2011).

Hyndman et al. (2011) proposed a least squares reconciliation approach providing an ex post adjustment to a set of “base” forecasts. That is, base forecasts (an initial set of forecasts) are produced for all disaggregate and aggregate series. These are then adjusted to ensure they are coherent. Wickramasuriya et al. (2019) clarified and generalized the results leading to *Minimum Trace (MinT)* reconciliation, guaranteeing minimum variance unbiased forecasts. Panagiotelis et al. (2021) shows that the predictive accuracy of reconciled forecasts cannot be worse than unreconciled forecasts in the mean squared error sense.

In a recent study, Hollyman et al. (2021) discuss the connection between forecast reconciliation and forecast combination, which is a widely used technique in financial forecasting (Rapach et al. 2010). Therefore, forecast reconciliation offers two crucial advantages to the forecasting process. First, it ensures that forecasts are coherent with the hierarchical structure. Second, it allows getting more accurate predictions, according to a mechanism similar to the one behind the combination of alternative forecasts.

Successful applications of reconciliation techniques have been proposed in tourism (Athanasopoulos et al. 2009), macroeconomics (Athanasopoulos et al. 2020, Eckert et al. 2021, Lila et al. 2022), demography (Yang et al. 2022), and energy (Jeon et al. 2019, Di Fonzo and Girolimetto 2023). Some early papers that discuss hierarchical forecasting for stock price indexes have been Lee and Swaminathan (1999) who used a bottom-up approach to forecast the Dow Jones Industrial Average Index, and Darrough and Russell (2002) who provided a comparison between bottom-up and top-down approaches for this aim. However, to the best of our knowledge, forecast reconciliation has not been used when forecasting stock prices. In the context of financial forecasting, Li and Tang (2019) recently adopted MinT reconciliation, but for predicting mortality bond indexes rather than

stock market indexes, while Caporin et al. (2024) proposed a reconciliation procedure for realized volatility. (Hyndman and Athanasopoulos 2021, Chapter 11) provide an introductory exposition to hierarchical forecasting and forecast reconciliation. Athanasopoulos et al. (2023) provide a comprehensive literature review.

Common stock prices form hierarchical time series because they linearly aggregate to the market index (Lee and Swaminathan 1999, Darrough and Russell 2002). In line with the aforementioned literature, we expect that more accurate forecasts of the time series at both the top of the hierarchy (i.e. the market index) and the bottom of the hierarchy (i.e. the common stocks included in the index) can be obtained by adjusting forecasts so that they become coherent. That is, adjust forecasts so that they aggregate in a coherent manner. The main issue in applying hierarchical forecasting to the stock market is that the structure of the hierarchy is unknown in advance, so it needs to be estimated.

Indeed, there is evidence that stocks are characterized by hierarchical (or clustering) structures of unknown form (e.g. see Brown and Goetzmann 1997, Raffinot 2017, Tumminello et al. 2010). Therefore, it may be beneficial to first identify, estimate and forecast such hierarchies, and then apply forecast reconciliation to exploit the hierarchical structures of the series. Multiple hierarchical configurations are possible, based on different clustering tools and different dissimilarity measures applied to the stock price data. We discuss how multiple hierarchies can be used in reconciling stock price and index forecasts.

The contribution of this paper is twofold. First, we apply optimal forecast reconciliation to a new domain, that is financial markets. Second, we develop a novel forecasting framework that combines reconciliation and clustering, by implementing optimal ex post adjustment to the forecasts with the aim of making them coherent with the underlying hierarchical structure, based on the clustering of individual stock price time series. We note that clustering has been adopted by previous studies both for portfolio selection (e.g. see Tola et al. 2008, Raffinot 2017) and as a tool for improving forecasting accuracy (e.g. Marsili 2002, Sáenz et al. 2023). Hence, our paper provides additional evidence of the usefulness of clustering in financial practice.

As an empirical experiment, we first consider the Dow Jones Industrial Average (DJIA) index and its constituents. The DJIA is the oldest stock market index in the United States and, for this reason, it is an established benchmark for tracking the overall market performance (Brown et al. 1998, Lee et al. 1999, Kim et al. 2011). Indeed, although it includes thirty stocks, the index has a broad market coverage with different sectors, so shocks affecting specific sectors are reflected by the index. Moreover, the DJIA represents a classical hierarchical time series because, as happens with any equally weighted index, the top-level series is obtained as a simple linear function of the bottom series. To further illustrate our proposed method, we also consider an application to the price-weighted version of the Standard & Poor 500 (S&P 500) Index.

Clustering is introduced using both available meta-data (such as industry group or stock exchange), and empirical clustering algorithms based on several dissimilarity measures. Forecasts are obtained and evaluated for both the bottom-level series (common stocks) and the top-level (index) series. The results demonstrate empirically that forecast reconciliation, and furthermore forecast reconciliation based on hierarchical structures identified by clustering, is useful for predicting the stock market index and its constituents. Moreover, further results of an investment exercise on the top-level series based on the reconciled forecasts show that these have economic significance.

The rest of the paper is structured as follows. Section 2 discusses the methodology employed in the paper, considering the combination of optimal forecast reconciliation and clustering. The data adopted for the main empirical experiment, based on the Dow Jones Industrial Average Index during the period 2020-2022, and the forecasting set-up are described in Section 3. Section 4 shows in detail the results of cluster analysis used for the estimation of the alternative hierarchies, while Section 5 provides a discussion on the forecasting results and the usefulness of the proposed approach for investment purposes. Section 5.4 provides further analyses in order to enhance the robustness of our findings. We consider the Dow Jones Industrial Average Index over a different period, 2015-2017, as well as the S&P 500 Index over the period 2022-2024, considering both

forecast accuracy and investment strategies. We conclude with some final remarks in the Section 6.

2. Forecast reconciliation with clustering structure

Let $\mathbf{p}_i = [p_{i,1}, \dots, p_{i,T}]'$ be the daily closing price time series of the i th stock. The DJIA stock price index at time t , denoted by y_t , is obtained by the sum of its N constituents:

$$y_t = \frac{\sum_{i=1}^N p_{i,t}}{\gamma_t}, \quad (1)$$

discounted by a factor γ_t — equal for all the stocks — which accounts for market operations such as changes in the index composition and stock splits. The adjusted price series are given by $b_{i,t} = p_{i,t}/\gamma_t$, and so we have the linear constraint:

$$y_t = \sum_{i=1}^N b_{i,t}. \quad (2)$$

2.1. Groups and clusters of stocks

Let \mathbf{b}_t be the vector of all N stocks of interest observed at time t , and let \mathbf{a}_t be a corresponding vector of n_a aggregated time series:

$$\mathbf{a}_t = \mathbf{A}\mathbf{b}_t. \quad (3)$$

The first element of \mathbf{a}_t is the stock index y_t . Other elements of \mathbf{a}_t are aggregations based on subsets of stocks. For example, suppose we aggregate the prices for each of the n_1 exchanges on which they are traded. Let $c_{i,j} = 1$ if stock j is traded on exchange i , and 0 otherwise, and define \mathbf{C}_1 to be the $n_1 \times N$ matrix with element $c_{i,j}$ in row i and column j . Then $\mathbf{C}_1\mathbf{b}_t$ gives the aggregated prices for all exchanges at time t . We can similarly define \mathbf{C}_2 to denote the grouping of stocks based on industries, where each row corresponds to a different industry group. Other groups or clusters of stocks can also be defined. This leads to the aggregation matrix

$$\mathbf{A} = \begin{bmatrix} \mathbf{1}' \\ \mathbf{C}_1 \\ \vdots \\ \mathbf{C}_L \end{bmatrix}, \quad (4)$$

where each \mathbf{C}_ℓ denotes a grouping or clustering of stocks. The “aggregation” matrix \mathbf{A} has dimension $n_a \times N$ and specifies how the stock price series \mathbf{b}_t aggregate to form \mathbf{a}_t . The first row of \mathbf{A} defines the sum of all stocks making up the stock index, while other rows define sums of different groups of stocks.

Some components of \mathbf{a}_t , and therefore \mathbf{A} , can be formed using a data-driven hierarchical (or partitioning) clustering method. Time series clustering can be divided into three well-known classes (Maharaj et al. 2019): observation-based, feature-based and model-based. Observation-based approaches group time series according to their observed values. Given a pair of bottom time series $\mathbf{b}_i = [b_{i1}, \dots, b_{iT}]'$ and $\mathbf{b}_j = [b_{j1}, \dots, b_{jT}]'$, a simple observation-based approach involves the use of

the standard Euclidean distance:

$$d_{\text{EUC}}(\mathbf{b}_i, \mathbf{b}_j) = \sqrt{\sum_{t=1}^T (b_{it} - b_{jt})^2}. \quad (5)$$

Observation-based approaches can be useful for clustering short time series, although they impose strong stationarity conditions on the original series. Notice that returns-based clustering approaches are widespread in the asset pricing literature (Brown and Goetzmann 1997, Brown et al. 2012). Correlation-based approaches (Mantegna 1999, Tola et al. 2008, Raffinot 2017), instead, use pairwise returns' correlations for computing dissimilarities across stocks. We therefore consider the correlation-based distance:

$$d_{\text{COR}}(\mathbf{b}_i, \mathbf{b}_j) = \sqrt{2(1 - \rho_{ij})},$$

where ρ_{ij} is the sample Pearson correlation coefficient between log-returns, defined as $r_{it} = \log(b_{it}) - \log(b_{it-1})$. This correlation-based distance fulfills the three axioms of a metric distance. Although other clustering methods based on correlations can be used (e.g. see Giada and Marsili 2001, Marsili 2002), this correlation distance is easy to compute and can be directly used in k -means and hierarchical clustering procedures. Moreover, from the financial viewpoint, this approach is linked with the Hierarchical Risk Parity, and it is widely used to cluster stocks to build portfolios in practice (e.g. see de Prado 2016, Lohre et al. 2020, Lee et al. 2023, Raffinot 2017).

Model-based approaches define distances based on parameter estimates from statistical models. A well-known example is the ARIMA-based distance. Given two time series \mathbf{b}_i and \mathbf{b}_j , Piccolo (1990) defined the distance between two invertible ARIMA processes as the Euclidean distance between the $\text{AR}(\infty)$ representation of the two series, i.e.:

$$d_{\text{ARIMA}}(\mathbf{b}_i, \mathbf{b}_j) = \sqrt{\sum_{k=1}^K (\pi_{i,k} - \pi_{j,k})^2}, \quad (6)$$

where $\pi_{i,k}$ denotes the k th “ π weight” (Box et al. 2016, p51) for the i th stock.

Given the distance measure, we also need to choose the clustering algorithm to be used. In this paper, we use the Partition Around Medoids algorithm (PAM, Kaufman and Rousseeuw 1990), which provides an iterative solution to the following minimization problem:

$$\min : \sum_{i=1}^N \sum_{c=1}^C d^2(\mathbf{b}_i, \mathbf{b}_c), \quad (7)$$

where $d^2(\mathbf{b}_i, \mathbf{b}_c)$ is the squared distance between the i th unit and the c th cluster centroid time series. (Any of the above distances measures may be used.) In what follows we focus on the PAM algorithm rather than C -means, because it is more interpretable (using a real time series as the centroid rather than an average), and is more robust to outliers.

The main drawback of the PAM clustering approach lies in the *a priori* selection of the number of clusters C . To address this issue, we follow Arbelaiz et al. (2013) and Batool and Hennig (2021) and use the Average Silhouette Width (ASW), a well-known cluster validity index for evaluating the quality of a partition, measuring the within-cluster cohesion and inter-cluster dispersion. The Silhouette for the i th object can be computed as

$$SW_i = \frac{g_i - f_i}{\max\{g_i, f_i\}} \quad (8)$$

where f_i is the average distance of the i th unit to the other units belonging to the same cluster, and g_i is the average distance of the same unit to others belonging to the closest different cluster. The number of clusters is commonly selected as the number C maximizing the Silhouette in Equation 8.

Once the number of clusters, C , has been selected, the result of cluster analysis on the individual stock time series is the membership matrix \mathcal{C} of dimension $C \times N$, with element $\mathcal{C}_{c,j} = 1$ when \mathbf{b}_j belongs to cluster c , and 0 otherwise.

We note that financial time series can exhibit various behaviors that a single clustering approach might not fully capture. In other words, by relying on a single clustering \mathcal{C}_ℓ , there is a risk of overfitting the cluster structure. Alternatively, combining different clustering structures enhances robustness and allows for uncovering distinct patterns in the data. In this regard, we note that the aggregation matrix defined in Equation 4 can be thought of as an ensemble of the different clustering $\mathcal{C}_1, \dots, \mathcal{C}_\ell$. This provides a more comprehensive understanding of the underlying dynamics of the data, and can help mitigate the risk of overfitting (de Prado 2020).

2.2. Forecast reconciliation

The full vector of time series at time t is given by

$$\mathbf{y}_t = \begin{bmatrix} \mathbf{a}_t \\ \mathbf{b}_t \end{bmatrix} = \mathbf{S}\mathbf{b}_t, \quad (9)$$

where $\mathbf{S} = \begin{bmatrix} \mathbf{A} \\ \mathbf{I}_N \end{bmatrix}$ denotes the “summation” matrix of dimension $n \times N$, where $n = n_a + N$. Let $\hat{\mathbf{y}}_h$ be the vector of h -step-ahead forecasts obtained with a generic forecasting model. Base forecasts $\hat{\mathbf{y}}_h$ generally do not sum up to the top levels, so we say they are not “coherent”. Forecast reconciliation methods aim at making forecasts coherent across the aggregation structure. We denote coherent forecasts as $\tilde{\mathbf{y}}_h$. Linear reconciliation can be written as follows:

$$\tilde{\mathbf{y}}_h = \mathbf{M}\hat{\mathbf{y}}_h, \quad (10)$$

where $\mathbf{M} = \mathbf{S}\mathbf{G}_h$ is a $n \times n$ mapping matrix, whose role is to project the base forecasts $\hat{\mathbf{y}}_h$ onto a coherent subspace (Panagiotelis et al. 2021). For example, in the bottom-up approach, we define $\mathbf{G}_h = [\mathbf{0} \quad \mathbf{I}_N]$, with $\mathbf{0}$ denoting a vector of zeros. The optimal least squares approach (known as MinT for Minimum Trace) is obtained (Wickramasuriya et al. 2019) with

$$\mathbf{G}_h = (\mathbf{S}'\mathbf{W}_h^{-1}\mathbf{S})^{-1}\mathbf{S}'\mathbf{W}_h^{-1}, \quad (11)$$

where \mathbf{W} is the $n \times n$ covariance matrix of the h -step base forecast errors.

Following previous studies (Tola et al. 2008, Raffinot 2017), it is interesting to highlight that, in this framework, each cluster can be seen as an equally weighted portfolio of the stocks included therein. Moreover, we note that there were no negative reconciled forecasts in the empirical experiment. Hence, there was no need to implement non-negativity constraints in the MinT reconciliation procedure to ensure positive prices. The proposed procedure could easily be extended to account for non-negativity constraints, following Wickramasuriya et al. (2020).

3. Stock market data and experimental set up

We evaluate the usefulness of reconciliation in forecasting the prices of the Dow Jones Industrial Average (DJIA) and of its constituents. We use clustering to determine some possible grouping structures as explained in the Section 2.

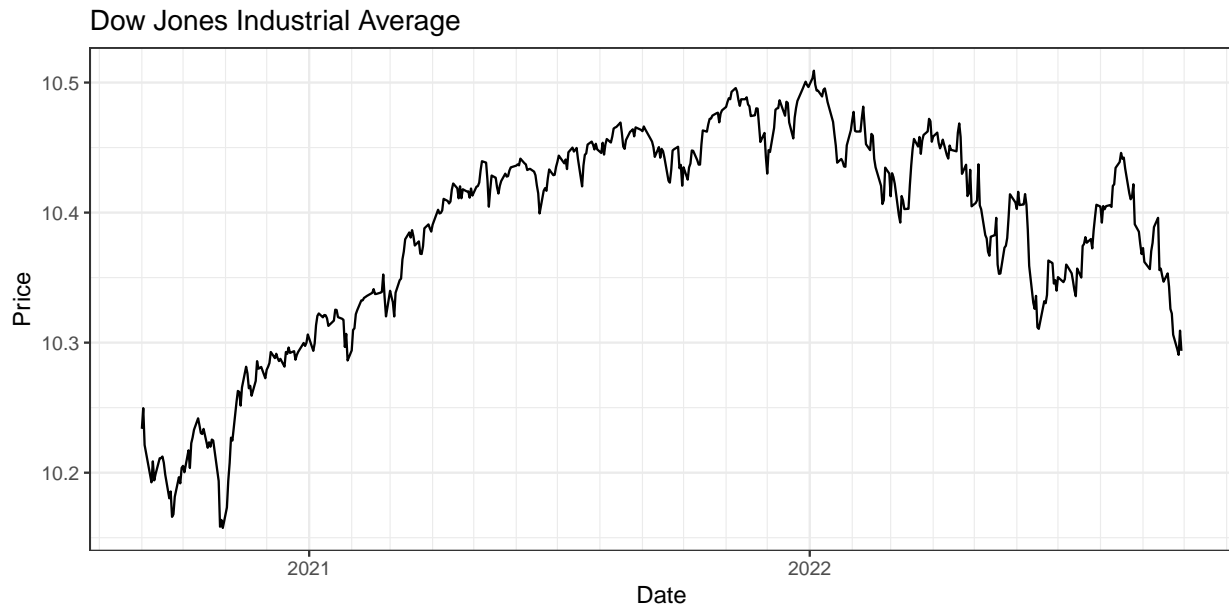


Figure 1. Dow Jones Industrial Average (DJIA) index: price time series. The time period spans from September 1, 2020 to September 30, 2022. Prices are shown on a logarithmic scale.

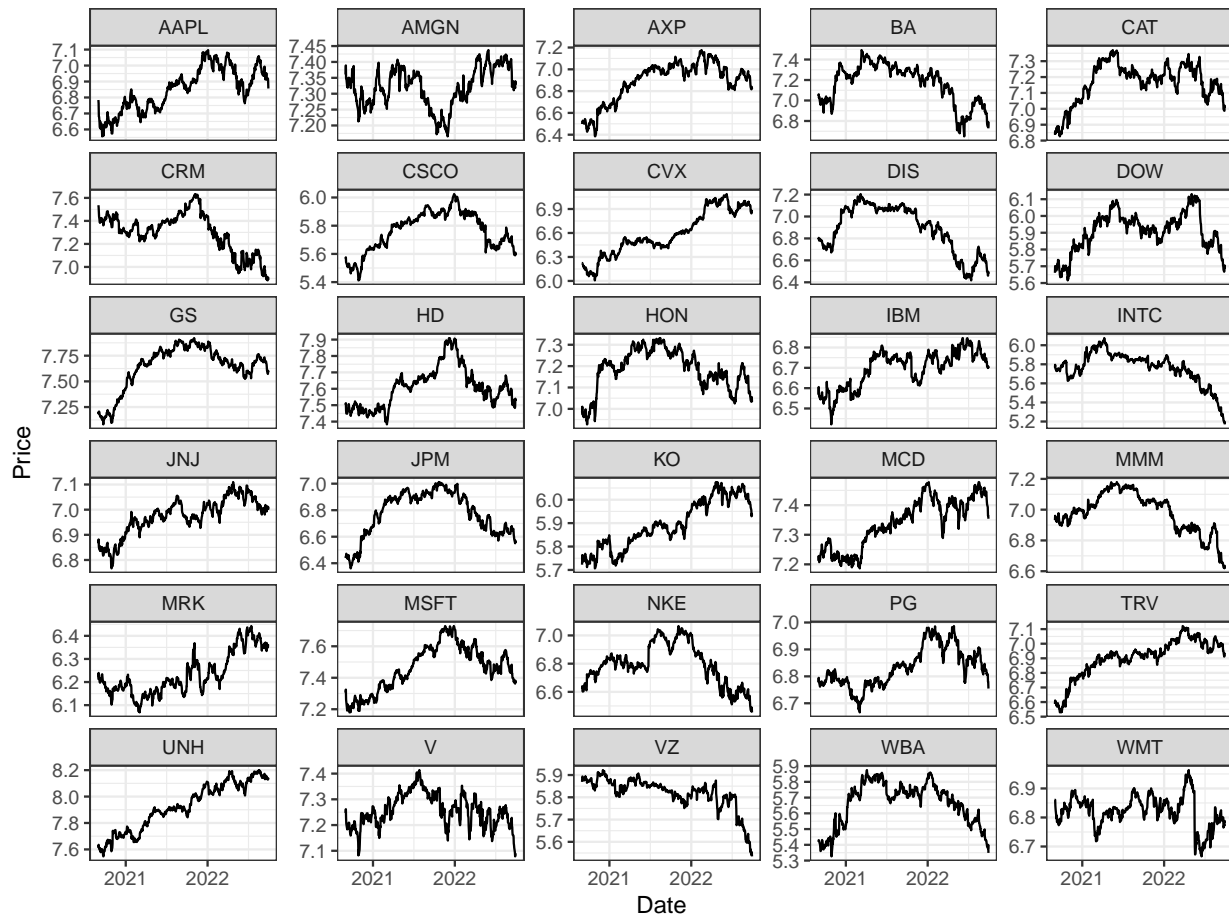


Figure 2. DJIA constituents: price time series. The time period spans from September 1, 2020 to September 30, 2022. Prices are shown on a logarithmic scale.

The dataset used for the empirical analysis consists of the daily time series associated with the DJIA index from 1 September 2020 to 30 September 2022. We made this choice because the DJIA composition (see Table A1 in Appendix A) changed on 31 August 2020 and has not been revised since; the divisor takes the constant value $d_t = 0.15$ during this period. The DJIA price time series is shown in Figure 1, while Figure 2 shows the time series of its $N = 30$ constituents.

Since we are working with non-stationary integrated price time series, we forecast the series using ARIMA models. A rolling-window procedure is used to obtain the forecasts, where at each step of the recursion we choose the best ARIMA model by means of the automatic procedure described in Hyndman and Khandakar (2008). We use the \mathbf{G} matrix resulting from the MinT approach of Wickramasuriya et al. (2019), shown in Equation 11.

The experimental design can be outlined as follows. The time series have length $T = 524$, and we leave the last $R = 124$ observations for out-of-sample testing. The clustering structures are estimated within the training set, i.e., considering only the first 400 observations. Forecasts at $h = \{1, 3, 6, 12\}$ steps ahead are produced, so the out-of-sample length is equal to $R - h$ ($r = 1, \dots, R - h$). At each r recursion, an estimation window of length 400 is considered for the model selection procedure and to make the h -step-ahead forecasts with the selected ARIMA model. We evaluate if the different MinT reconciliation approaches improve with respect to base forecasts, random walk and bottom-up reconciliation.

Let us define $e_t = \hat{y}_{t+h} - y_{t+h}$ the h step ahead forecasting error. Forecast accuracy is evaluated in terms of both absolute and squared errors. For the top level and bottom level series MAE and RMSE are considered, that is:

$$\text{RMSE} = \sqrt{\frac{1}{R-h} \sum_{t=T-R+h}^T e_t^2} \quad (12)$$

and:

$$\text{MAE} = \frac{1}{R-h} \sum_{t=T-R+h}^T |e_t| \quad (13)$$

Then, predictive accuracy tests of Diebold and Mariano (2002) and van Dijk and Franses (2003) are considered to evaluate the statistical significance of the forecasting error differences. Let us define $d_t = g(e_{1,t}) - g(e_{2,t})$ the error differential between two forecasting approaches up to some transformation $g(\cdot)$; in this paper, we use squares $g(e_{1,t}) = e_{1,t}^2$ and absolute values $g(e_{1,t}) = |e_{1,t}|$. Assuming covariance stationarity of the loss differential series d_t , Diebold and Mariano (2002) show that the sample mean of the loss differential,

$$\bar{d} \equiv \frac{1}{R-h} \sum_{t=T-R+h}^T d_t, \quad (14)$$

follows an asymptotically standard Normal distribution. Therefore, testing the null hypothesis of equal forecast accuracy can be obtained by calculating the following statistic:

$$\text{DM} = \bar{d}[V(\bar{d})]^{-1/2}, \quad (15)$$

where $V(\bar{d})$ is consistently estimated assuming a certain autocorrelation structure of the forecast errors. We might expect that forecasting models provide poor forecasts in the presence of trend reversions or changes of regimes. Therefore, we also consider an alternative testing approach, proposed in van Dijk and Franses (2003), that is based on the introduction of a weighting scheme

providing lower weights to "less relevant" information. The resulting statistic is given by:

$$\bar{d}_w \equiv \frac{1}{R-h} \sum_{t=T-R+h}^T \frac{\hat{f}(y_t)}{\max \hat{f}(y_t)} d_t, \quad (16)$$

where $\hat{f}(\cdot)$ denotes the density function of y_t estimated by means of a standard Nadaraya-Watson kernel estimator with a Gaussian kernel. In doing so, we provide less weight to the forecast errors occurring when the realized value y_{t+h} is in the tails of the distribution. The resulting weighted Diebold-Mariano statistic follows a standard Normal distribution asymptotically. In this case the covariance $V(\bar{d}_w)$ is computed considering the weighted loss differentials and their autocovariances.

4. Finding clustering structures: results

In this section, we explore various possible sets of clusters to be used in the forecast reconciliation.

4.1. *No clustering structure*

A natural starting point in our setting is the use of the MinT forecast reconciliation approach which employs no clustering structure at all. In other words, in this case we use Equation 11 considering $\mathbf{A} = \mathbf{1}'$.

4.2. *Industry-based clustering*

There is a long tradition in financial economics of considering the presence of industry-based clusters of stocks. The rationale behind this clustering approach is that stocks belonging to the same industry sector are affected by common shocks (e.g. King 1966, Livingston 1977). In the case of the Dow-Jones stocks in our sample, we have $C = 20$ Industry-based clusters (see Table A1, Appendix A). Several groups are based on singletons, due to the relatively large number of industries relative to the number of stocks included in the DJIA. The implied hierarchy is shown in Figure B1, Appendix B. The hierarchy shows the bottom-level time series (stocks) and the middle-level time series (clusters). We do not show the top-level time series, the Dow Jones Industrial Average Index, since it is already presented in Figure 1. The cluster structure is highlighted with different colors. In the case of industry-based clustering, we have that Financial Services and IT industries include 4 and 5 stocks respectively, more than any other industry groups, which is why their aggregate values are larger than for the other industries.

4.3. *Exchange-based clustering*

This clustering approach involves $C = 2$ groups, because DJIA stocks are traded at NYSE and NASDAQ only. The groups are quite unbalanced because most of the stocks included in the DJIA index are traded on the NYSE exchange. The NYSE is known to have a higher average market capitalization of listed companies compared to the NASDAQ so that, on average, the companies listed on the NYSE are larger and more established than those listed on the NASDAQ. Moreover, NASDAQ is known for its focus on IT companies, while NYSE lists a wider range of industries. The exchange-based clusters are shown in Table A1, Appendix A. The implied hierarchy is shown in Figure B2, Appendix B.

4.4. Observation-based clustering

We then consider the unsupervised approaches discussed in Section 2 for constructing hierarchies. The resulting clusters are presented in Table 1. We first construct clusters of stocks considering the standard Euclidean distance between log returns, that is:

$$x_{i,t} = \log(p_{i,t}) - \log(p_{i,t-1}). \quad (17)$$

We use log returns rather than prices, so that we consider stationary series. As explained in Section 2, we choose the number of clusters C maximizing ASW. Figure B6, Appendix B, shows the ASW associated with $C \in \{1, \dots, 10\}$ different number of clusters, with $C = 3$ giving the maximum ASW. The resulting clusters are shown in Panel A of Table 1, and the aggregated series are shown in Figure B3, Appendix B.

Cluster	Stocks
Panel A: Observation-based	
A	AAPL, CRM, CSCO, HD, INTC, MCD, MSFT, NKE, V
B	AMGN, JNJ, KO, MMM, MRK, PG, UNH, VZ, WMT
C	AXP, BA, CAT, CVX, DIS, DOW, GS, HON, IBM, JPM, TRV, WBA
Panel B: Correlation-based	
A	AAPL, CRM, INTC, MSFT, NKE
B	AMGN, CSCO, HD, HON, IBM, JNJ, KO, MCD, MMM, MRK, PG, UNH, V, VZ, WBA, WMT
C	AXP, BA, CAT, CVX, DIS, DOW, GS, JPM, TRV
Panel C: Model-based	
A	AAPL, AMGN, DOW, JNJ, KO, MRK, MSFT, PG
B	CRM, CSCO, HON, IBM, MMM
C	GS, HD, MCD, WMT
D	TRV, VZ
E	AXP, BA, CAT, CVX, DIS, INTC, JPM, NKE, UNH, V, WBA

Table 1. Cluster assignment of stocks according to three different methods. Panel A shows the clusters obtained using Euclidean distance on returns; Panel B shows the clusters obtained using correlation-based distance on returns; and, Panel C shows the clusters obtained using ARIMA-based distance on prices.

4.5. Correlation-based clustering

We also construct clusters of stocks based on the correlation between log-returns. The correlation-based approach is one of the most widely adopted for financial time series clustering. Figure B7 in Appendix B shows the ASW associated with different number of clusters, again showing $C = 3$ as the maximum value (although the ASW value is relatively lower compared to the previous observation-based clustering). The resulting clusters are shown in Panel B of Table 1 and the aggregated series are shown in Figure B4, Appendix B.

4.6. Model-based clustering

The first step of ARIMA-based clustering requires estimating the best fitting models within the training set. The resulting models identified using the automatic ARIMA algorithm of Hyndman and Khandakar (2008) are shown in Table A2, Appendix A, for readers' convenience.

The stocks with random walk models (shown as ARIMA(0,1,0)) are grouped in a separate cluster, and we apply cluster analysis on the remaining stocks using the Piccolo (1990) distance across their $AR(\infty)$ coefficients. The ASW values are shown in Figure B8, with the maximum given by $C = 4$. The resulting clusters are given in Panel C of Table 1, along with the random walk cluster labeled E.

Figure B9, Appendix B, shows the $AR(\infty)$ weights for each stock, colored by cluster. Cluster E includes stocks with zero coefficients (random walk processes), while Cluster D contains two stocks with persistent $AR(\infty)$ coefficients. Clusters A and C are characterized by similar patterns of the coefficients, but the parameters of stocks in Cluster A decay to zero faster than those included in Cluster C. Figure B5, Appendix B, shows the aggregated series resulting from these clusters.

5. Forecast reconciliation: results

In this section we provide details about the forecasting experiments. We compare the forecast accuracy of implementing forecast reconciliation using the various hierarchical structures based on clustering, versus various benchmarks. We consider as benchmarks: Base (unreconciled) forecasts; random walk (RW) forecasts, bottom-up (BU) forecasts; and MinT reconciled forecasts without using any clustering (MinT). Base and random walk forecasts provide natural benchmarks without reconciliation. BU shows the usefulness of MinT reconciliation versus using a traditional the single level approach, while MinT without clustering allows us to evaluate the usefulness of clustering within the MinT framework. In Section 5.1 we discuss the results in terms of the market index, while in Section 5.2 we analyze the results for the common individual stocks included in the index. Section 5.3 discusses the results of an investment strategy based on reconciled forecasts and Section 5.4 presents the results of a robustness analysis.

5.1. Forecasting top series: Dow Jones Index Average

Table 2 shows the out-of-sample average errors at different forecasting horizons. Panel A shows the accuracy metric in terms of absolute errors, while Panel B is in terms of squared errors. In each case we are always able to find a reconciliation approach that is more accurate than the considered benchmarks, although the rankings of the reconciliation approaches vary. This indicates that reconciliation can be successfully employed to improve forecasts of the stock market index. In terms of MAE loss, the reconciliation approaches combining all the different clustering structures —

Panel A: MAE loss	$h = 1$	$h = 3$	$h = 6$	$h = 12$
Base (unreconciled)	333.65	625.31	925.13	1436.55
RW	320.15	588.88	889.49	1344.02
BU	321.12	588.62	873.97	1309.89
MinT	320.62	588.39	873.20	1307.99
MinT: IND	320.59	589.20	871.64	1299.28
MinT: EXCH	320.02	588.39	872.82	1308.80
MinT: EUCL	321.00	588.50	874.69	1302.19
MinT: COR	320.88	589.20	872.10	1302.65
MinT: ARMA	319.40	590.07	873.83	1307.92
MinT: ALL	319.34	590.92	868.85	1288.51
Panel B: RMSE loss	$h = 1$	$h = 3$	$h = 6$	$h = 12$
Base (unreconciled)	439.87	794.42	1212.65	1831.28
RW	424.80	749.47	1091.46	1570.48
BU	426.09	747.00	1076.13	1517.53
MinT	425.86	746.92	1074.50	1513.70
MinT: IND	426.24	750.88	1079.38	1512.85
MinT: EXCH	424.37	745.96	1073.09	1513.67
MinT: EUCL	426.32	748.78	1076.69	1510.22
MinT: COR	426.33	748.72	1075.55	1508.20
MinT: ARMA	424.86	747.45	1074.01	1511.86
MinT: ALL	424.93	752.71	1078.47	1508.26

Table 2. Accuracy evaluation for the top-level series (the DJIA Index). MAE and RMSE are presented for different forecast horizons h . The best model is highlighted in bold.

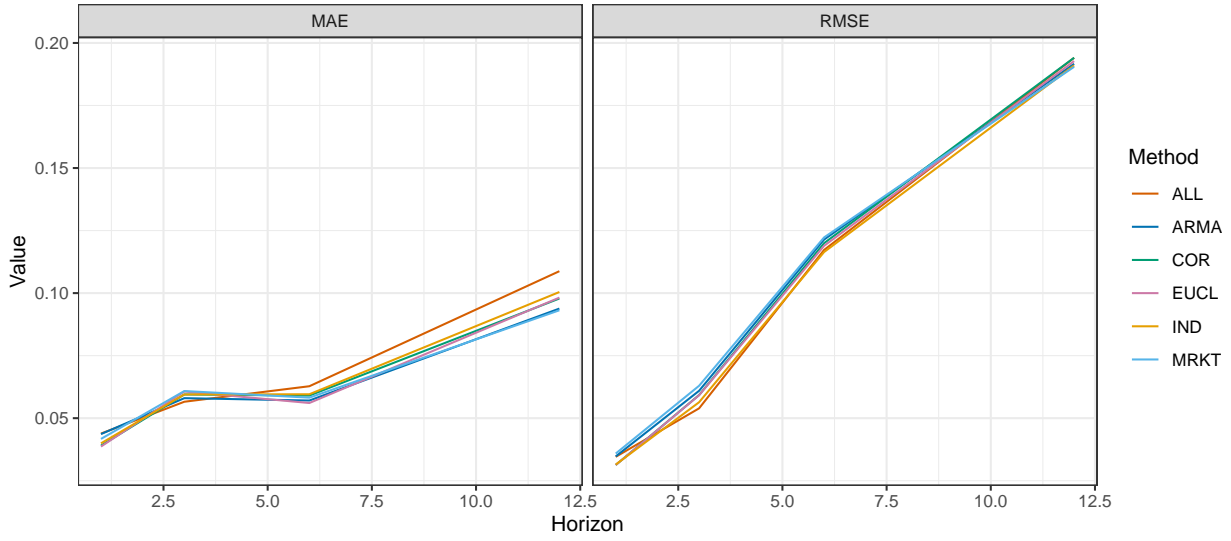


Figure 3. Average relative accuracy of reconciliation methods compared to the base unreconciled forecasts. Values above zero indicate higher (an improvement in) forecast accuracy. Left: Forecast accuracy is given by the log of the Mean Absolute Errors (MAE) of the reconciled forecasts relative to the base forecasts. Right: Forecast accuracy is given by the log of the Root Mean Squared Errors (RMSE) of the reconciled forecasts relative to the base forecasts.

denoted by MinT: ALL — provides the most accurate forecasts for $h = 1, 6, 12$ horizons. Exchange-based clustering dominates the alternatives for $h = 3$ in terms of MAE loss and is the best approach for $h = 3, 6$ in terms of RMSE. ARIMA-based clustering provides the most accurate out-of-sample forecasts considering RMSE for $h = 1$ step ahead. Correlation-based clustering seems to be the best approach for $h = 12$ under RMSE loss. Hence, there is no clear clustering structure dominating the others across all forecasting horizons and error measures. However, MinT: ALL provides the most consistent results in terms of MAE loss. As this approach also represents a simple way of hedging against the wrong specification of the clustering structure, we consider this approach as the best among the considered alternatives.

The most important result, overall, is that reconciliation provides more accurate forecasts than unreconciled approaches. Figure 3 shows the average relative accuracy of reconciliation methods in terms of MAE and RMSE, respectively. Values greater than zero indicate higher forecast accuracy relative to the base unreconciled forecasts. In particular, we observe how the benefit of reconciliation becomes larger with increasing forecasting horizon h . For $h = 1$ step forecasts, the benefit of reconciliation is around 5%, while for $h = 12$ it is larger than 10% for MAE loss and around 20% under RMSE loss. Moreover, we notice that with increasing forecasting horizons, all the reconciliation based forecasts provide similar improvements in accuracy compared with base forecasts.

Table 3 shows the results of predictive accuracy tests, with absolute errors and squared errors, respectively. The p-values of the modified Diebold and Mariano (2002) test of van Dijk and Franses (2003) are reported. Under the null hypothesis, the average error of the benchmark, shown across the columns, is equal to the average error from the reconciliation approach, shown down the rows. The results in terms of standard Diebold and Mariano (2002), which provide similar results, are shown in Appendix A (see Table A3).

Base. The first important result to highlight is that for most forecast horizons the reconciliation approaches provide statistically more accurate forecasts than the base benchmark, regardless of the loss employed for the tests. This suggests that, whatever the forecasting horizon is, forecast reconciliation should be used to improve forecasting of the market index.

Random Walk. Not all the reconciliation approaches provide statistically significant improvements in forecasting the top series compared with the naïve approach. Nevertheless, under absolute error loss for $h = 1, 6, 12$, MinT: ALL reconciliation, results to statistically more accurate forecasts

Absolute errors	Base	RW	BU	MinT	Squared errors	Base	RW	BU	MinT
$h = 1$ step ahead									
MinT: IND	0.02	0.96	0.09	0.27	MinT: IND	0.02	0.99	0.46	0.51
MinT: EXCH	0.03	0.60	0.04	0.10	MinT: EXCH	0.02	0.22	0.08	0.09
MinT: EUCL	0.03	0.95	0.03	0.11	MinT: EUCL	0.02	0.98	0.16	0.16
MinT: COR	0.03	0.99	0.39	0.80	MinT: COR	0.02	1.00	0.93	0.96
MinT: ARMA	0.02	0.13	0.00	0.00	MinT: ARMA	0.02	0.65	0.00	0.00
MinT: ALL	0.02	0.06	0.01	0.03	MinT: ALL	0.02	0.12	0.04	0.05
$h = 3$ steps ahead									
MinT: IND	0.02	0.19	0.12	0.27	MinT: IND	0.02	0.32	0.79	0.87
MinT: EXCH	0.02	0.33	0.04	0.29	MinT: EXCH	0.02	0.15	0.07	0.13
MinT: EUCL	0.02	0.28	0.02	0.12	MinT: EUCL	0.02	0.21	0.36	0.72
MinT: COR	0.02	0.31	0.06	0.40	MinT: COR	0.02	0.20	0.39	0.74
MinT: ARMA	0.02	0.43	0.56	0.73	MinT: ARMA	0.02	0.13	0.47	0.66
MinT: ALL	0.02	0.14	0.15	0.23	MinT: ALL	0.02	0.23	0.36	0.43
$h = 6$ steps ahead									
MinT: IND	0.02	0.01	0.06	0.05	MinT: IND	0.02	0.02	0.09	0.27
MinT: EXCH	0.03	0.06	0.11	0.08	MinT: EXCH	0.02	0.04	0.03	0.08
MinT: EUCL	0.03	0.05	0.53	0.73	MinT: EUCL	0.02	0.04	0.40	0.90
MinT: COR	0.02	0.04	0.10	0.09	MinT: COR	0.02	0.03	0.06	0.26
MinT: ARMA	0.02	0.03	0.46	0.57	MinT: ARMA	0.02	0.02	0.20	0.54
MinT: ALL	0.02	0.01	0.06	0.06	MinT: ALL	0.02	0.02	0.09	0.15
$h = 12$ steps ahead									
MinT: IND	0.01	0.04	0.10	0.10	MinT: IND	0.01	0.01	0.19	0.32
MinT: EXCH	0.01	0.14	0.57	0.96	MinT: EXCH	0.01	0.03	0.19	0.63
MinT: EUCL	0.01	0.11	0.01	0.00	MinT: EUCL	0.01	0.03	0.01	0.00
MinT: COR	0.01	0.09	0.09	0.06	MinT: COR	0.01	0.02	0.02	0.01
MinT: ARMA	0.01	0.11	0.52	0.61	MinT: ARMA	0.01	0.02	0.20	0.33
MinT: ALL	0.01	0.01	0.04	0.03	MinT: ALL	0.01	0.00	0.12	0.17

Table 3. P-values from the modified Diebold and Mariano (2002) test of van Dijk and Franses (2003) for the top-level series (the DJIA Index). Under the null, the difference in the out of sample forecast errors between the model in the row and the model in the column is equal to zero. Entries in bold indicate a rejection of the null at a 5% level of significance.

compared to the random walk. Some statistically significant improvements are also observed for other approaches, especially for the longer forecast horizons. The results show that random walk forecasts remain a valid (but less accurate) alternative, but only for shorter forecast horizons and under squared loss. In the other settings, reconciliation provides statistically more accurate forecasts.

Bottom-up. For all scenarios we are able to find a clustering-based MinT reconciliation approach which performs statistically better than bottom-up. The strongest results are shown for MinT: ALL, which provides statistically more accurate forecasts than bottom-up for $h = 1, 6, 12$ for absolute loss and for $h = 1, 6$ for squared loss. The overall evidence suggests that the clustering-based reconciliation approaches provide more accurate forecasts than bottom-up.

MinT without clustering. In general, we observe that clustering-based approaches are more accurate than when no clustering is considered, with many statistically significant entries in the MinT columns. In particular, under absolute error loss for $h = 1, 6, 12$ the MinT: ALL approach provides the more accurate forecasts (as shown in Table 2) which are also statistically significantly better than all the benchmarks, including MinT with no clustering. Using squared loss, MinT: ARMA and MinT: ALL provide statistically lower forecast errors than MinT with no clustering, while for $h = 3$, the best model is MinT: EXCH which gives statistically lower forecast errors. For $h = 12$ the MinT: COR is the most accurate (as shown in Table 2) and statistically better than all the benchmarks including MinT with no clustering.

These findings provide strong evidence that not only is MinT reconciliation useful for financial forecasting, but considering and exploring clustering structures in stocks, allows further improvements for forecasting the market index.

5.2. Forecasting bottom series: common stocks

In this section we evaluate the effect of forecast reconciliation in improving the accuracy of forecasting the individual stock prices. As these form the bottom-level series we remove from the benchmarks the bottom-up forecasts, as these are identical to the Base (unreconciled) forecasts.

Table 4 shows the results in terms of average loss across the 30 stocks included the index. Regardless of the loss and the forecast horizon, we are always able to find a reconciliation approach that performs better than the base forecasts. This indicates that forecast reconciliation can be successfully employed for more accurately forecasting both the index and its constituents. Most clustering approaches improve on MinT with no clustering. Considering squared loss, the random walk provides the most accurate forecasts for short horizons, $h = 1, 3$, while MinT: ARMA provides the most accurate forecasts for the longer horizons $h = 6, 12$. Considering absolute error loss, the MinT reconciliation approaches are generally more accurate than the random walk.

Table 5 shows the results of predictive accuracy tests, with absolute errors and squared errors, respectively. Under the null hypothesis, the average error of the benchmark, shown across the columns, is equal to the average error from the reconciliation approach, shown down the rows. In contrast to the tests applied to the market index, here we report the results of the Diebold and Mariano (2002) test, because average errors from different time series are considered. The results of the modified test van Dijk and Franses (2003) are reported in Appendix A (see Table A4) for the sake of consistency.

Base. In general, the MinT reconciliation approaches including clustering, provide statistically significant improvements over the base forecasts for both absolute and squared error losses. The number of statistically significant improvements seems to increase as the forecast horizon increases. In particular, for $h = 12$ all the clustering-based MinT approaches, with the exception of MinT: IND for absolute loss, provide statistically more accurate forecasts than base forecasts, which do not involve reconciliation. Hence, the results send a strong signal that forecast reconciliation improves the forecast accuracy for individual stock price series.

Random Walk. Most reconciliation procedures allow reducing the forecast error compared with a random walk, especially for longer forecast horizons. However, in contrast to forecasting the market index, we do not find any statistically significant differences in forecast accuracy between the reconciliation approaches and the random walk for forecasting the individual stock price series.

Panel A: average MAE loss	$h = 1$	$h = 3$	$h = 6$	$h = 12$
Base (unreconciled)	15.40	27.65	39.81	57.11
RW	15.35	27.55	39.77	57.38
MinT	15.40	27.64	39.77	57.01
MinT: IND	15.40	27.64	39.77	56.85
MinT: EXCH	15.38	27.61	39.70	56.89
MinT: EUCL	15.42	27.65	39.77	56.86
MinT: COR	15.41	27.66	39.75	56.86
MinT: ARMA	15.37	27.61	39.72	56.80
MinT: ALL	17.78	28.50	42.72	58.52
Panel B: average RMSE loss	$h = 1$	$h = 3$	$h = 6$	$h = 12$
Base (unreconciled)	20.23	35.72	50.49	70.57
RW	20.14	35.59	50.37	70.73
MinT	20.23	35.71	50.42	70.38
MinT: IND	20.23	35.77	50.38	69.99
MinT: EXCH	20.18	35.65	50.25	70.09
MinT: EUCL	20.24	35.73	50.37	70.09
MinT: COR	20.24	35.75	50.34	70.06
MinT: ARMA	20.20	35.65	50.21	69.90
MinT: ALL	22.79	35.74	51.30	77.37

Table 4. Accuracy evaluation for the bottom-level series (the DJIA Index constituents). Average MAE and RMSE are calculated across all individual stocks for different forecast horizons h . The best model is highlighted in bold.

Absolute errors	Base	RW	MinT	Squared errors	Base	RW	MinT
$h = 1$ step ahead							
MinT: IND	0.08	1.00	0.19	MinT: IND	0.42	1.00	0.64
MinT: EXCH	0.05	1.00	0.12	MinT: EXCH	0.08	1.00	0.13
MinT: EUCL	0.66	1.00	0.79	MinT: EUCL	0.64	1.00	0.84
MinT: COR	0.34	1.00	0.59	MinT: COR	0.70	1.00	0.96
MinT: ARMA	0.00	0.99	0.00	MinT: ARMA	0.00	1.00	0.00
MinT: ALL	0.15	0.88	0.19	MinT: ALL	0.34	0.99	0.42
$h = 3$ steps ahead							
MinT: IND	0.27	1.00	0.36	MinT: IND	0.81	1.00	0.85
MinT: EXCH	0.06	0.99	0.12	MinT: EXCH	0.06	0.95	0.10
MinT: EUCL	0.45	0.99	0.53	MinT: EUCL	0.47	0.98	0.57
MinT: COR	0.68	1.00	0.80	MinT: COR	0.77	1.00	0.87
MinT: ARMA	0.10	1.00	0.12	MinT: ARMA	0.01	0.97	0.01
MinT: ALL	0.47	0.95	0.50	MinT: ALL	0.81	0.99	0.84
$h = 6$ steps ahead							
MinT: IND	0.16	0.80	0.34	MinT: IND	0.31	0.71	0.51
MinT: EXCH	0.02	0.70	0.06	MinT: EXCH	0.03	0.56	0.06
MinT: EUCL	0.14	0.80	0.37	MinT: EUCL	0.01	0.62	0.05
MinT: COR	0.12	0.77	0.25	MinT: COR	0.03	0.61	0.08
MinT: ARMA	0.01	0.72	0.02	MinT: ARMA	0.00	0.48	0.00
MinT: ALL	0.10	0.61	0.14	MinT: ALL	0.16	0.59	0.26
$h = 12$ steps ahead							
MinT: IND	0.11	0.28	0.34	MinT: IND	0.01	0.32	0.09
MinT: EXCH	0.00	0.29	0.02	MinT: EXCH	0.01	0.35	0.03
MinT: EUCL	0.00	0.26	0.00	MinT: EUCL	0.00	0.31	0.00
MinT: COR	0.02	0.25	0.06	MinT: COR	0.00	0.29	0.00
MinT: ARMA	0.00	0.22	0.00	MinT: ARMA	0.00	0.26	0.00
MinT: ALL	0.03	0.17	0.06	MinT: ALL	0.01	0.26	0.03

Table 5. P-values from the Diebold and Mariano (2002) test for the bottom-level series (the DJIA Index constituents). Under the null, the difference in the (average) out-of-sample forecast errors between the model in the row and the model in the column is equal to zero. Entries in bold indicate a rejection of the null at a 5% level of significance.

That is, we are not able to reject the null hypothesis of the test for any of the forecasting horizons h and error losses.

MinT without clustering. In all cases, there exist clustering-based reconciliation approaches that generate statistically significant forecast accuracy improvements compared to reconciliation with no clustering. In particular, MinT: ARMA and MinT: EXCH consistently provide statistically more accurate forecasts for all forecast horizons and error losses. Similar to the comparison with the base forecasts, the number of statistically significant improvements increases as the forecast horizon increases. Hence, coupling forecast reconciliation with clustering structures seems to also be beneficial for improving on average the forecasting accuracy for the bottom-level series.

5.3. Investing with forecast reconciliation

Forecasts are commonly used for constructing profitable trading strategies in finance. For this reason, an alternative approach of comparing forecast methods is by studying the profitability of investment strategies based on the obtained forecasts. While the previous sections highlight the advantageous properties of forecast reconciliation and clustering in reducing out-of-sample forecast errors, in this section we evaluate the performance of the developed approach from a financial investment viewpoint.

We compare the performance of two investment strategies that involve buying the market index if the price is predicted to increase in the next h time periods, and selling it otherwise. The two strategies differ in the forecast method employed. In particular, we compare the performance of

alternative reconciliation and clustering approaches against the base ARIMA model, bottom-up reconciliation and MinT without clustering.

Let us define as \hat{y}_{t+h} the market index price forecast at time $t + h$ and as y_t the actual price at time t . The predicted return at time $t + h$ implied by the price forecast is computed as,

$$\hat{x}_{t+h} = \frac{\hat{y}_{t+h} - y_t}{y_t}, \quad (18)$$

while the actual return observed at time $t + h$ is given by,

$$x_{t+h} = \frac{y_{t+h} - y_t}{y_t}. \quad (19)$$

Following Anatolyev and Gerko (2005), the realized return of this investment strategy at time $t + h$ is given by,

$$r_{t+h} = \text{sign}(\hat{x}_{t+h})x_{t+h}, \quad (20)$$

where $\text{sign}(\hat{x}_{t+h})$ is a function taking the value of 1 if $\hat{x}_{t+h} \geq 0$ and -1 otherwise. Therefore, assuming that short selling is allowed, investors realize positive returns if they buy the market index at time t and its price increases at time $t + h$, but also if they sell it at time t and the price is lower at time $t + h$.

Given two alternative forecasting methods, A and B , we construct two alternative investment strategies with different ex-post realized returns, $r_{A,t+h}$ and $r_{B,t+h}$. We compare the forecast methods in terms of their implied financial performance, measured by the Sharpe ratio (Sharpe 1963), which is a commonly used metric for evaluating the performance of an investment strategy in terms of its return/risk trade-off. The Sharpe ratio of the strategy employing forecasts from method A is given by,

$$\text{SR}_A = \frac{\hat{\mu}_A}{\hat{\sigma}_A}, \quad (21)$$

where $\hat{\mu}_A$ is the average return of $r_{A,t+h}$ and $\hat{\sigma}_A$ is its risk, computed with the standard deviation. The Sharpe ratio of the strategy using forecast from method B is similarly defined. We then test if the two strategies lead to statistically different Sharpe ratios using the procedure proposed by Jobson and Korkie (1981) and Memmel (2003).

Table 6 shows the difference in Sharpe ratio of investment strategies based on alternative forecasting methods. We report the results for the top-level series for the out-of-sample forecasting evaluation as defined in Section 3. The results are evaluated at forecast horizons $h = 1, 3, 6, 12$. We consider base ARIMA, bottom-up, and MinT reconciliation without clustering as the benchmarks. The alternative forecasting methods based on MinT reconciliation under different clustering structures are presented in the rows. Positive entries indicate that the investment strategy based on the forecasting approach in the row provides a higher Sharpe ratio than the benchmark considered in the column.

First, we highlight that for $h = 1$ and $h = 6$ step-ahead forecasts, the investment strategy based on MinT: ALL forecasts is the only one that provides a statistically higher Sharpe ratio compared with base ARIMA forecasts. The Sharpe ratio is 14% higher than the one obtained using ARIMA-based forecasts. For the other forecast horizons, we do not reject the null hypothesis of the test. The results obtained for $h = 1$ step-ahead forecasts are arguably the most interesting from a financial viewpoint, as trading strategies are usually constructed and evaluated considering short-term forecasts. Forecasts about financial prices are indeed less accurate for longer horizons, as we also have highlighted in the previous section.

$h = 1$	Base	BU	MinT	$h = 3$	Base	BU	MinT
MinT: IND	-4.16	14.30***	11.40***	MinT: IND	-2.40	-2.07	-2.07
MinT: EXCH	-8.87	9.59***	6.69***	MinT: EXCH	-2.43	-2.43	-2.43
MinT: EUCL	-11.47	6.99***	4.09***	MinT: EUCL	-0.32	0.00	0.00
MinT: COR	-15.43	3.03**	0.13	MinT: COR	-0.32	0.00	0.00
MinT: ARMA	-0.12	18.34***	15.44***	MinT: ARMA	0.36	0.68**	0.68**
MinT: ALL	13.97***	32.44***	29.54***	MinT: ALL	-2.40	-2.07	-2.07
$h = 6$	Base	BU	MinT	$h = 12$	Base	BU	MinT
MinT: IND	-2.37	-0.87	-0.87	MinT: IND	-3.57	0.00	0.00
MinT: EXCH	0.22	0.00	0.00	MinT: EXCH	0.52	0.00	0.00
MinT: EUCL	-1.70	-0.20	-0.20	MinT: EUCL	-1.39	2.18***	2.18***
MinT: COR	0.22	1.72***	1.72***	MinT: COR	-3.57	0.00	0.00
MinT: ARMA	-1.50	0.00	0.00	MinT: ARMA	-3.57	0.00	0.00
MinT: ALL	1.53***	3.03***	3.03***	MinT: ALL	-1.39	2.18***	2.18***

Table 6. Difference of Sharpe ratios (%) between two forecast-based investment strategies on the DJIA Index. The forecasts of the models in the rows are compared with those of the models in the column. Positive values, highlighted in bold, indicate that the forecast method in the row provides a higher Sharpe ratio than the benchmark method in the column. Note: *** indicate the rejection of the null hypothesis at a 1% level of significance, ** at 5% level, and * at 10% level.

Another interesting result in Table 6 is that, for $h = 1$, all the reconciliation approaches provide statistically higher Sharpe ratios compared with both bottom-up and MinT reconciliation without clustering. Therefore, the combined use of clustering and MinT reconciliation is also useful in constructing more profitable trading strategies. Specifically, MinT: ALL has a Sharpe ratio 32.4% higher than the one obtained with bottom-up forecasts, and 29.5% higher than forecasts based on MinT reconciliation without clustering.

Considering longer horizons h , we find a reconciliation procedure outperforming the bottom-up and MinT benchmarks in all cases. In particular, MinT: ALL is the best model in most scenarios, namely $h = 1, 6$ and 12 . For $h = 3$ forecasting horizon, MinT: ARMA is the only model providing a statistically higher Sharpe ratio than bottom-up and MinT reconciliation without clustering.

The above results, however, do not consider the presence of transaction costs. In the realm of investment analysis, transaction costs need to be included in the assessment of portfolio performance. In particular, the inclusion of fixed transaction costs encompasses brokerage fees, taxes, and other non-variable expenses. Ignoring these costs could lead to a distorted view of a portfolio's true performance, potentially undermining the accuracy of the investment decisions. In our setting, we assume that the investor using the forecast-based investment strategies pays a fixed amount of money c , proportional to the price of the index, for making each trade. Table 7 shows the comparison in terms of Sharpe ratios computed on net returns, where the net return of a strategy is the realized return at time $t + h$ given a fixed transaction cost $c = 0.5\%$ paid by the investor at time t . More precisely, the net return of the forecast-based strategy at time $t + h$ can be computed as

$$r_{t+h} = \text{sign}(\hat{x}_{t+h}) \tilde{x}_{t+h} \quad (22)$$

where $\tilde{x}_{t+h} = \left(\frac{y_{t+h}}{(1+c)y_t} - 1 \right)$ is the net return with $(1+c)y_t$ being the cost of the transaction made at time t and $\text{sign}(\hat{x}_{t+h})$ the indication about the trade's direction. It is easy to see that if $c = 0$ no transaction cost is involved, and the Equation 22 coincides with the Equation 19.

In sum, the results in Table 7 provide the same indications as those shown in Table 6, suggesting that investment strategies based on the index and constructed by the use of forecast reconciliation provide better performance also from an investment perspective, especially for investment decisions based on one-step ahead forecasts. It is interesting to highlight, however, that the investment strategy based on the MinT:ALL reconciliation under $h = 1$ forecasts increases its performance considerably compared to its benchmark. For example, the difference between the

$h = 1$	Base	BU	MinT	$h = 3$	Base	BU	MinT
MinT: IND	12.89***	5.21***	2.48	MinT: IND	-3.90	-3.04	-3.04
MinT: EXCH	17.95***	10.27***	7.54***	MinT: EXCH	-3.65	-2.79	-2.79
MinT: EUCL	11.65***	3.96**	1.23	MinT: EUCL	-0.87	0.00	0.00
MinT: COR	9.87***	2.19	-0.54	MinT: COR	-0.87	0.00	0.00
MinT: ARMA	27.25***	19.57***	16.84***	MinT: ARMA	-0.25	0.62***	0.62***
MinT: ALL	33.27***	25.59***	22.86***	MinT: ALL	-3.90	-3.04	-3.04
$h = 6$	Base	BU	MinT	$h = 12$	Base	BU	MinT
MinT: IND	-0.99	-1.38	-1.38	MinT: IND	-3.27	0.00	0.00
MinT: EXCH	0.39**	0.00	0.00	MinT: EXCH	-3.27	0.00	0.00
MinT: EUCL	-1.89	1.38***	1.38***	MinT: EUCL	-1.89	1.38*	1.38*
MinT: COR	1.37***	0.98***	0.98*	MinT: COR	-3.27	0.00	0.00
MinT: ARMA	0.39**	0.00	0.00	MinT: ARMA	-3.27	0.00	0.00
MinT: ALL	0.80***	0.41***	0.41***	MinT: ALL	-1.89	1.38*	1.38*

Table 7. Difference of Sharpe ratios (%) between two forecast-based investment strategies on the DJIA Index, given a fixed transaction cost equal to $c = 0.5\%$. The forecasts of the models in the rows are compared with those of the models in the column. Positive values, highlighted in bold, indicate that the forecast method in the row provides a higher Sharpe ratio than the benchmark method in the column, given the presence of fixed transaction costs. Note: *** indicate the rejection of the null hypothesis at a 1% level of significance, ** at 5% level, and * at 10% level.

Sharpe ratios of the MinT:ALL compared to the baseline forecasts increases from 14% to 33%. Moreover, other reconciliation approaches provide with worse results under the assumption of no transaction costs, now provide much better Sharpe ratios compared to the baseline strategy. Then, considering investment strategies based on longer forecasting horizons h , no considerable differences can be found compared with the case of no transaction costs, even if it seems that the differences between the Sharpe ratios of alternative strategies become lower.

5.4. Robustness analysis

In what follows we aim to further demonstrate the usefulness of reconciliation by providing additional empirical evidence. We consider two different applications: the DJIA in a different time period, and the Standard and Poor 500 (S&P 500) Index in the last two years. For the DJIA, we consider the period August 2015 – August 2017, since the index composition remained constant and the divisor took a constant value equal to $d_t = 0.146$. For the S&P 500 index, we consider its price-weighted version. Financial products based on price-weighted indexes are directly investable in the financial market and usually offer better performances (DeMiguel et al. 2009, Yuan and Zhou 2023). We focus on top-level series forecast accuracy and its economic significance by investing with the obtained forecasts. We therefore show that the benefit of reconciliation is not limited to the specific period considered in the main application or by the selected index.

In line with the main application, we evaluate the forecast reconciliation performance under the Industry-based clustering and the unsupervised clustering approaches discussed in Section 2. We excluded the exchange-based clustering because few stocks were included in the DJIA Index in the period 2015–2017. We also include the MinT: ALL approach that combines different cluster structures. For the sake of brevity, we do not show the clustering results here, but the results are available upon request. In both applications, we leave the last six months for out-of-sample testing and consider an initial estimation window of one and a half years. The main results in terms of forecasting accuracy are presented in Table 8, while the results of the investment analysis are shown in Table 9.

We first study the accuracy results for the top-level series forecasting. Table 8 shows MAE and RMSE losses for the forecast horizons $h = 1, 3, 6, 12$. The results of the pairwise predictive accuracy tests are shown in Table A5 for the DJIA, and in Table A6 for the S&P 500. In summary, the results are consistent with the ones obtained in the main application for DJIA in 2020–2022, in that unreconciled base forecasts perform worse than the reconciled forecasts, and improvements in

accuracy are more evident at larger forecasting horizons. The best method for the DJIA application is the MinT: ALL, while for the S&P500 we find the MinT: ARMA to perform best.

Next, we evaluate if reconciled forecasts have economic value. To this aim, we study the same forecast-based strategy on the index discussed in Section 5.3. Table 9 shows the net returns for different forecast horizons, assuming a fixed transaction cost of $c = 0.5\%$. Consistent with the main results, we find that the MinT: ALL reconciliation provides the best results.

In the case of the DJIA application, we find a larger economic value of the reconciled forecasts compared to the analysis conducted in the period 2020–2022. For $h = 1$ the investment strategy employing MinT: ALL forecasts provides a Sharpe ratio 12.8% larger than the base approach, and 20.3% larger compared to both bottom-up and MinT without clustering. Considering longer horizons, MinT: ALL still provides the best results. For instance, it provides a Sharpe ratio that is 6.2% larger compared to the base approach for $h = 3$, 12% larger for $h = 6$, and 2.3% larger for $h = 12$. MinT: ALL provides better results also when compared to other reconciliation approaches. For the S&P500 application, we also find the net returns are overall higher for reconciled forecasts compared to the base approach obtained from the prediction of the index, especially for a short forecast horizon. In particular, considering $h = 1$, the investment strategy employing MinT: ALL forecasts provides a Sharpe ratio 7% larger than the base approach.

DJIA (2015-2017)					S&P 500 (2022-2024)				
Panel A: MAE loss	$h = 1$	$h = 3$	$h = 6$	$h = 12$	Panel A: MAE loss	$h = 1$	$h = 3$	$h = 6$	$h = 12$
Base (unreconciled)	55.06	95.51	130.98	193.18	Base (unreconciled)	607.43	929.54	1309.48	2325.93
RW	54.44	96.22	132.53	177.56	RW	610.49	899.57	1181.61	1941.92
BU	53.87	95.38	129.50	179.06	BU	605.17	894.79	1202.30	1930.68
MinT	54.03	95.29	129.36	178.60	MinT	603.46	890.15	1181.63	1866.98
MinT: IND	54.31	96.21	130.02	174.83	MinT: IND	599.66	861.59	1143.18	1785.85
MinT: EUCL	54.03	95.26	129.04	178.21	MinT: EUCL	600.30	889.70	1168.81	1799.93
MinT: COR	54.10	95.36	129.01	177.16	MinT: COR	615.50	960.74	1380.84	2359.57
MinT: ARMA	54.14	95.44	128.82	176.98	MinT: ARMA	599.56	872.05	1132.51	1728.89
MinT: ALL	54.54	96.50	129.59	171.98	MinT: ALL	609.96	909.82	1285.67	2053.36
Panel B: RMSE loss	$h = 1$	$h = 3$	$h = 6$	$h = 12$	Panel B: RMSE loss	$h = 1$	$h = 3$	$h = 6$	$h = 12$
Base (unreconciled)	76.19	125.56	164.17	224.76	Base (unreconciled)	745.88	1136.33	1489.94	2533.27
RW	75.57	125.71	163.67	210.86	RW	744.61	1103.06	1351.90	2202.12
BU	75.41	125.15	164.15	212.91	BU	743.16	1098.95	1372.58	2214.09
MinT	75.59	125.16	163.93	212.17	MinT	741.72	1090.20	1352.15	2163.31
MinT: IND	75.71	125.78	164.08	208.29	MinT: IND	738.29	1069.85	1315.20	2122.75
MinT: EUCL	75.62	125.23	163.87	211.83	MinT: EUCL	739.40	1085.39	1334.29	2106.88
MinT: COR	75.74	125.23	163.66	210.96	MinT: COR	755.46	1159.44	1571.30	2656.32
MinT: ARMA	75.98	125.27	163.61	210.08	MinT: ARMA	737.64	1064.44	1298.95	2028.17
MinT: ALL	76.20	126.27	163.74	205.58	MinT: ALL	748.11	1109.50	1467.07	2440.44

Table 8. Accuracy evaluation for the top-level series: MAE and RMSE results for different forecasting horizons h .

DJIA (2015–2017) and S&P 500 (2022–2024) are considered. The best model is highlighted in bold.

6. Conclusions

The objective of this research has been to explore the potential benefits of employing forecast reconciliation for forecasting stock market indexes and their underlying constituents. Both meta-data groups and empirical clustering techniques have been used to determine the underlying structure of the price time series. The study makes two contributions. First, to the best of our knowledge, it applies forecast reconciliation to the financial domain for the first time. Second, it combines cluster analysis with forecast reconciliation. This approach offers insights into the efficacy of reconciliation within the context of latent hierarchical structures.

To evaluate our proposed approach, we apply it to the Dow Jones Industrial Average index and its constituents in different time periods, and to the Standard & Poor 500 index. ARIMA models are used to generate forecasts for the time series using a rolling-window procedure. The reconciliation approach combines MinT reconciliation (Wickramasuriya et al. 2019) with cluster analysis. Three different dissimilarity measures are utilized for PAM-based clustering: raw returns, return

DJIA (2015-2017)				S&P 500 (2022-2024)			
	Base	BU	MinT		Base	BU	MinT
$h = 1$							
MinT: IND	7.40*	14.91***	14.91***	MinT: IND	-0.45	6.26***	6.26***
MinT: EUCL	-4.32	3.19	3.19	MinT: EUCL	-15.14	-8.44	-8.44
MinT: COR	-9.01	-1.51	-1.51	MinT: COR	-3.98	2.72***	2.72***
MinT: ARMA	-3.46	4.04	4.04	MinT: ARMA	-21.11	-14.40	-14.40
MinT: ALL	12.80***	20.30***	20.30***	MinT: ALL	7.02***	13.73***	13.73***
$h = 3$							
MinT: IND	6.21***	0.54*	0.54*	MinT: IND	-11.63	-2.96	-2.96
MinT: EUCL	5.67***	0.00	0.00	MinT: EUCL	-13.18	-4.52	-4.52
MinT: COR	5.67***	0.00	0.00	MinT: COR	-8.40	0.27	0.27
MinT: ARMA	6.21***	0.54*	0.54*	MinT: ARMA	-13.18	-4.52	-4.52
MinT: ALL	6.21***	0.54*	0.54*	MinT: ALL	-12.18	-3.51	-3.51
$h = 6$							
MinT: IND	9.17***	2.35***	2.35***	MinT: IND	-10.12	-8.58	-4.41
MinT: EUCL	6.82***	0.00	0.00	MinT: EUCL	-5.71	-4.17	0.00
MinT: COR	6.82***	0.00	0.00	MinT: COR	-8.13	-6.60	-2.42
MinT: ARMA	12.54***	5.71***	5.71***	MinT: ARMA	-1.03	0.50	4.68***
MinT: ALL	12.06***	5.23***	5.23***	MinT: ALL	-7.62	-6.09	-1.91
$h = 12$							
MinT: IND	0.34	-3.89	-2.31	MinT: IND	-9.35	-11.14	-4.97
MinT: EUCL	2.65***	-1.57	0.00	MinT: EUCL	-4.37	-6.17	0.00
MinT: COR	2.65***	-1.57	0.00	MinT: COR	-6.06	-7.85	-1.69
MinT: ARMA	2.65***	-1.57	0.00	MinT: ARMA	-4.37	-6.17	0.00
MinT: ALL	2.28***	-1.94	-0.37	MinT: ALL	-13.52	-15.32	-9.15

Table 9. Difference of Sharpe ratios (%) between two forecast-based investment strategies on the DJIA Index (2015-2017) and S&P500 Index (2022-2024), given a fixed transaction cost equal to $c = 0.5\%$. The forecasts of the models in the rows are compared with those of the models in the column. Positive values, highlighted in bold, indicate that the forecast method in the row provides a higher Sharpe ratio than the benchmark method in the column, given the presence of fixed transaction costs. Note: *** indicate the rejection of the null hypothesis at a 1% level of significance, ** at 5% level, and * at 10% level.

correlation, and ARIMA distance, as proposed by Piccolo (1990). We show how all the clusters can be used simultaneously within a forecast reconciliation context. Furthermore, we investigate the usefulness of clustering by considering a reconciliation approach without a clustering structure, where stocks aggregate directly to the market index.

We evaluate the usefulness of reconciliation and clustering in terms of out-of-sample forecasting accuracy. Our results suggest that reconciliation is a useful tool for forecasting both the stock market index and its underlying constituents, even without clustering. But with the clustering of stocks included in the reconciliation procedure, even better forecasts are obtained. We also evaluate the usefulness of forecast reconciliation and clustering from a financial viewpoint, considering investment strategies built on alternative forecasts. We show that reconciliation and clustering can be successfully used for constructing profitable trading strategies based on forecasts.

By comparing the out-of-sample forecast accuracy and the Sharpe ratios associated with alternative reconciliation procedures at different forecasting horizons, we find that MinT: ALL (i.e., MinT reconciliation combining different clustering structures) provides the best performance compared with other approaches for the top-level series. Indeed, MinT: ALL is often ranked as the best model while forecasting the top-level time series in out-of-sample and, when we consider the economic significance of the forecasts through the investment exercise, we find that it provides the best performances in all the considered experiments. We notice that the MinT: ALL approach ensembles different hierarchies. Therefore, we suggest using it to avoid relying on a single cluster structure. We find this choice to be economically worthy, at least for short investment horizons.

However, the proposed methodology presents some limitations. Changes in market conditions (e.g. structural breaks) that are not well-represented in the historical data might reduce the out-of-sample forecast accuracy if the cluster structure is not updated to account for such changes. This could potentially be addressed by future research developing more sophisticated approaches for computing hierarchies. For instance, an online reconciliation approach (e.g. Brégère and Huard 2022) which also updates the hierarchical structure in real-time could be a promising solution to this issue. Another solution may be found in the market state forecasting techniques (e.g. Marsili 2002). These are beyond the scope of this paper and we leave these to future studies.

Disclosure of interest

There are no interests to declare.

Declaration of funding

No funding was received.

References

- Alexander, C. (2008). *Market risk analysis, pricing, hedging and trading financial instruments*, volume 3. John Wiley & Sons.
- Allen, D., Lizieri, C., and Satchell, S. (2019). In defense of portfolio optimization: What if we can forecast? *Financial Analysts Journal*, 75(3):20–38.
- Anatolyev, S. and Gerko, A. (2005). A trading approach to testing for predictability. *Journal of Business & Economic Statistics*, 23(4):455–461.
- Arbelaitz, O., Gurrutxaga, I., Muguerza, J., Pérez, J. M., and Perona, I. (2013). An extensive comparative study of cluster validity indices. *Pattern Recognition*, 46(1):243–256.
- Athanasopoulos, G., Ahmed, R. A., and Hyndman, R. J. (2009). Hierarchical forecasts for Australian domestic tourism. *International Journal of Forecasting*, 25(1):146–166.
- Athanasopoulos, G., Gamakumara, P., Panagiotelis, A., Hyndman, R. J., and Affan, M. (2020). Hierarchical forecasting. In Fuleky, P., editor, *Macroeconomic Forecasting in the Era of Big Data. Advanced Studies in Theoretical and Applied Econometrics.*, chapter 21, pages 689–719. Springer, Cham, vol 52 edition.
- Athanasopoulos, G., Hyndman, R. J., Kourentzes, N., and Panagiotelis, A. (2023). Forecast reconciliation: A review. *International Journal of Forecasting*, forthcoming.
- Batool, F. and Hennig, C. (2021). Clustering with the average silhouette width. *Computational Statistics & Data Analysis*, 158:107190.
- Bisaglia, L., Di Fonzo, T., and Girolimetto, D. (2020). Fully reconciled GDP forecasts from income and expenditure sides. In Pollice, A., Salvati, N., and Schirripa Spagnolo, F., editors, *Book of Short Papers SIS 2020*, pages 951–956. Pearson.
- Blaskowitz, O. and Herwartz, H. (2011). On economic evaluation of directional forecasts. *International Journal of Forecasting*, 27(4):1058–1065.
- Box, G. E. P., Jenkins, G. M., Reinsel, G. C., and Ljung, G. M. (2016). *Time Series Analysis: Forecasting and Control*. John Wiley and Sons, 5th edition.
- Brégère, M. and Huard, M. (2022). Online hierarchical forecasting for power consumption data. *International Journal of Forecasting*, 38(1):339–351.
- Brown, S. J. and Goetzmann, W. N. (1997). Mutual fund styles. *Journal of Financial Economics*, 43(3):373–399.
- Brown, S. J., Goetzmann, W. N., and Kumar, A. (1998). The Dow theory: William peter hamilton’s track record reconsidered. *The Journal of Finance*, 53(4):1311–1333.
- Brown, S. J., Lajbcygier, P., and Wong, W. W. (2012). Estimating the cost of capital with basis assets. *Journal of Banking & Finance*, 36(11):3071–3079.

- Caporin, M., Di Fonzo, T., and Girolimetto, D. (2024). Exploiting intraday decompositions in realized volatility forecasting: A forecast reconciliation approach. *Journal of Financial Econometrics*, pages 1–26.
- Darrrough, M. N. and Russell, T. (2002). A positive model of earnings forecasts: Top down versus bottom up. *Journal of Business*, 75(1):127–152.
- de Prado, M. L. (2016). Building diversified portfolios that outperform out of sample. *The Journal of Portfolio Management*, 42(4):59–69.
- de Prado, M. M. L. (2020). *Machine learning for asset managers*. Cambridge University Press.
- DeMiguel, V., Garlappi, L., and Uppal, R. (2009). Optimal versus naive diversification: How inefficient is the $1/n$ portfolio strategy? *The Review of Financial Studies*, 22(5):1915–1953.
- Di Fonzo, T. and Girolimetto, D. (2023). Spatio-temporal reconciliation of solar forecasts. *Solar Energy*, 251:13–29.
- Diebold, F. X. and Mariano, R. S. (2002). Comparing predictive accuracy. *Journal of Business & Economic Statistics*, 20(1):134–144.
- Dunn, D. M., Williams, W. H., and DeChaine, T. (1976). Aggregate versus subaggregate models in local area forecasting. *Journal of the American Statistical Association*, 71(353):68–71.
- Eckert, F., Hyndman, R. J., and Panagiotelis, A. (2021). Forecasting Swiss exports using Bayesian forecast reconciliation. *European Journal of Operational Research*, 291(2):693–710.
- Fama, E. F. (1995). Random walks in stock market prices. *Financial Analysts Journal*, 51(1):75–80.
- Giada, L. and Marsili, M. (2001). Data clustering and noise undressing of correlation matrices. *Physical Review E*, 63(6):061101.
- Giamattei, M., Huber, J., Lambsdorff, J. G., Nicklisch, A., and Palan, S. (2020). Who inflates the bubble? forecasters and traders in experimental asset markets. *Journal of Economic Dynamics and Control*, 110:103718.
- Green, J., Hand, J. R., and Zhang, X. F. (2013). The supraview of return predictive signals. *Review of Accounting Studies*, 18(3):692–730.
- Gross, C. W. and Sohl, J. E. (1990). Disaggregation methods to expedite product line forecasting. *Journal of Forecasting*, 9(3):233–254.
- Hollyman, R., Petropoulos, F., and Tipping, M. E. (2021). Understanding forecast reconciliation. *European Journal of Operational Research*, 294(1):149–160.
- Hyndman, R. J., Ahmed, R. A., Athanasopoulos, G., and Shang, H. L. (2011). Optimal combination forecasts for hierarchical time series. *Computational Statistics & Data Analysis*, 55(9):2579–2589.
- Hyndman, R. J. and Athanasopoulos, G. (2021). *Forecasting: Principles and Practice*. OTexts, Melbourne, Australia, 3rd edition.
- Hyndman, R. J. and Khandakar, Y. (2008). Automatic time series forecasting: the forecast package for R. *Journal of Statistical Software*, 27:1–22.
- Hyndman, R. J., Lee, A. J., and Wang, E. (2016). Fast computation of reconciled forecasts for hierarchical and grouped time series. *Computational Statistics & Data Analysis*, 97:16–32.
- Jeon, J., Panagiotelis, A., and Petropoulos, F. (2019). Probabilistic forecast reconciliation with applications to wind power and electric load. *European Journal of Operational Research*, 279(2):364–379.
- Jobson, J. D. and Korkie, B. M. (1981). Performance hypothesis testing with the sharpe and treynor measures. *Journal of Finance*, pages 889–908.
- Kaufman, L. and Rousseeuw, P. J. (1990). *Finding groups in data: an introduction to cluster analysis*. John Wiley & Sons.
- Kim, J. H., Shamsuddin, A., and Lim, K.-P. (2011). Stock return predictability and the adaptive markets hypothesis: Evidence from century-long US data. *Journal of Empirical Finance*, 18(5):868–879.
- King, B. F. (1966). Market and industry factors in stock price behavior. *Journal of Business*, 39(1):139–190.
- Kong, A., Rapach, D. E., Strauss, J. K., and Zhou, G. (2011). Predicting market components out of sample: asset allocation implications. *Journal of Portfolio Management*, 37(4):29–41.
- Kumbure, M. M., Lohrmann, C., Luukka, P., and Porras, J. (2022). Machine learning techniques and data for stock market forecasting: A literature review. *Expert Systems with Applications*, 197:116659.
- Lee, C. M., Myers, J., and Swaminathan, B. (1999). What is the intrinsic value of the Dow? *The Journal of Finance*, 54(5):1693–1741.
- Lee, C. M. C. and Swaminathan, B. (1999). Valuing the Dow: A bottom-up approach. *Financial Analysts Journal*, 55(5):4–23.
- Lee, Y., Thompson, J. R., Kim, J. H., Kim, W. C., Fabozzi, F. A., Fabozzi, F. J., Musumeci, J., Feibel,

- B., Cornell, B., Nagy, Z., et al. (2023). An overview of machine learning for asset management. *The Journal of Portfolio Management*, 49(9):31–63.
- Li, H. and Tang, Q. (2019). Analyzing mortality bond indexes via hierarchical forecast reconciliation. *ASTIN Bulletin*, 49(3):823–846.
- Lila, M. F., Meira, E., and Oliveira, F. L. C. (2022). Forecasting unemployment in Brazil: A robust reconciliation approach using hierarchical data. *Socio-Economic Planning Sciences*, 82:101298.
- Livingston, M. (1977). Industry movements of common stocks. *Journal of Finance*, 32(3):861–874.
- Lohre, H., Rother, C., and Schäfer, K. A. (2020). Hierarchical risk parity: accounting for tail dependencies in multi-asset multi-factor allocations. *Machine learning for asset management: new developments and financial applications*, pages 329–368.
- Maharaj, E. A., D’Urso, P., and Caiado, J. (2019). *Time series clustering and classification*. Chapman and Hall/CRC.
- Makridakis, S., Spiliotis, E., and Assimakopoulos, V. (2022). M5 accuracy competition: Results, findings, and conclusions. *International Journal of Forecasting*, 38(4):1346–1364.
- Mantegna, R. N. (1999). Hierarchical structure in financial markets. *European Physical Journal B-Condensed Matter and Complex Systems*, 11(1):193–197.
- Marquering, W. and Verbeek, M. (2004). The economic value of predicting stock index returns and volatility. *Journal of Financial and Quantitative Analysis*, 39(2):407–429.
- Marsili, M. (2002). Dissecting financial markets: sectors and states. *Quantitative Finance*, 2(4):297.
- Memmel, C. (2003). Performance hypothesis testing with the sharpe ratio. *Finance Letters*, 1:21–23.
- Panagiotelis, A., Athanasopoulos, G., Gamakumara, P., and Hyndman, R. J. (2021). Forecast reconciliation: A geometric view with new insights on bias correction. *International Journal of Forecasting*, 37(1):343–359.
- Panagiotelis, A., Gamakumara, P., Athanasopoulos, G., and Hyndman, R. J. (2023). Probabilistic forecast reconciliation: Properties, evaluation and score optimisation. *European Journal of Operational Research*, 306(2):693–706.
- Piccolo, D. (1990). A distance measure for classifying ARIMA models. *Journal of Time Series Analysis*, 11(2):153–164.
- Raffinot, T. (2017). Hierarchical clustering-based asset allocation. *Journal of Portfolio Management*, 44(2):89–99.
- Rapach, D. E., Strauss, J. K., and Zhou, G. (2010). Out-of-sample equity premium prediction: Combination forecasts and links to the real economy. *The Review of Financial Studies*, 23(2):821–862.
- Sáenz, J. V., Quiroga, F. M., and Bariviera, A. F. (2023). Data vs. information: Using clustering techniques to enhance stock returns forecasting. *International Review of Financial Analysis*, 88:102657.
- Sharpe, W. F. (1963). A simplified model for portfolio analysis. *Management Science*, 9(2):277–293.
- Sutcliffe, C. M. (2018). *Stock index futures*. Routledge.
- Tola, V., Lillo, F., Gallegati, M., and Mantegna, R. N. (2008). Cluster analysis for portfolio optimization. *Journal of Economic Dynamics and Control*, 32(1):235–258.
- Trippi, R. R. and DeSieno, D. (1992). Trading equity index futures with a neural network. *Journal of Portfolio management*, 19:27–27.
- Tumminello, M., Lillo, F., and Mantegna, R. N. (2010). Correlation, hierarchies, and networks in financial markets. *Journal of Economic Behavior & Organization*, 75(1):40–58.
- van Dijk, D. and Franses, P. H. (2003). Selecting a nonlinear time series model using weighted tests of equal forecast accuracy. *Oxford Bulletin of Economics and Statistics*, 65:727–744.
- Wickramasuriya, S. L., Athanasopoulos, G., and Hyndman, R. J. (2019). Optimal forecast reconciliation for hierarchical and grouped time series through trace minimization. *Journal of the American Statistical Association*, 114(526):804–819.
- Wickramasuriya, S. L., Turlach, B. A., and Hyndman, R. J. (2020). Optimal non-negative forecast reconciliation. *Statistics and Computing*, 30(5):1167–1182.
- Yang, Y., Shang, H. L., and Raymer, J. (2022). Forecasting Australian fertility by age, region, and birthplace. *International Journal of Forecasting*.
- Yuan, M. and Zhou, G. (2023). Why naive diversification is not so naive, and how to beat it? *Journal of Financial and Quantitative Analysis*, pages 1–32.

Appendix A: Additional tables

Table A1. DJIA composition after last revision occurred the 31/08/2020

Company	Exchange	Symbol	Industry
Procter & Gamble	NYSE	PG	Fast-moving consumer goods
3M	NYSE	MMM	Conglomerate
IBM	NYSE	IBM	Information technology
Merck	NYSE	MRK	Pharmaceutical industry
American Express	NYSE	AXP	Financial services
McDonald's	NYSE	MCD	Food industry
Boeing	NYSE	BA	Aerospace and defense
Coca-Cola	NYSE	KO	Drink industry
Caterpillar	NYSE	CAT	Construction and Mining
Disney	NYSE	DIS	Broadcasting and entertainment
JPMorgan Chase	NYSE	JPM	Financial services
Johnson & Johnson	NYSE	JNJ	Pharmaceutical industry
Walmart	NYSE	WMT	Retailing
Home Depot	NYSE	HD	Home Improvement
Intel	NASDAQ	INTC	Semiconductor industry
Microsoft	NASDAQ	MSFT	Information technology
Verizon	NYSE	VZ	Telecommunications industry
Chevron	NYSE	CVX	Petroleum industry
Cisco	NASDAQ	CSCO	Information technology
Travelers	NYSE	TRV	Insurance
UnitedHealth	NYSE	UNH	Managed health care
Goldman Sachs	NYSE	GS	Financial services
Nike	NYSE	NKE	Clothing industry
Visa	NYSE	V	Financial services
Apple	NASDAQ	AAPL	Information technology
Walgreens Boots Alliance	NASDAQ	WBA	Retailing
Dow	NYSE	DOW	Chemical industry
Amgen	NASDAQ	AMGN	Biopharmaceutical
Honeywell	NASDAQ	HON	Conglomerate
Salesforce	NYSE	CRM	Information technology

Table A2. Best ARIMA models using the Hyndman-Khandakar algorithm with the AICc criterion.

Stock	Model	Stock	Model	Stock	Model
AAPL	ARIMA(3,1,1)	GS	ARIMA(0,1,1)	MRK	ARIMA(1,1,1)
AMGN	ARIMA(1,1,1)	HD	ARIMA(2,1,3)	MSFT	ARIMA(2,1,2)
AXP	ARIMA(0,1,0)	HON	ARIMA(2,1,0)	NKE	ARIMA(0,1,0)
BA	ARIMA(0,1,0)	IBM	ARIMA(2,1,0)	PG	ARIMA(1,1,0)
CAT	ARIMA(0,1,0)	INTC	ARIMA(0,1,0)	TRV	ARIMA(2,1,1)
CRM	ARIMA(2,1,1)	JNJ	ARIMA(2,1,0)	UNH	ARIMA(0,1,0)
CSCO	ARIMA(2,1,0)	JPM	ARIMA(0,1,0)	V	ARIMA(0,1,0)
CVX	ARIMA(0,1,0)	KO	ARIMA(1,1,1)	VZ	ARIMA(2,1,1)
DIS	ARIMA(0,1,0)	MCD	ARIMA(1,1,3)	WBA	ARIMA(0,1,0)
DOW	ARIMA(1,1,0)	MMM	ARIMA(1,1,2)	WMT	ARIMA(2,1,3)

Absolute Errors	Base	RW	BU	MinT	Squared Errors	Base	RW	BU	MinT
$h = 1$ step ahead									
MinT: IND	0.01	0.89	0.04	0.47	MinT: IND	0.01	0.98	0.62	0.74
MinT: EXCH	0.02	0.38	0.02	0.12	MinT: EXCH	0.01	0.34	0.07	0.11
MinT: EUCL	0.02	0.97	0.41	0.82	MinT: EUCL	0.02	1.00	0.70	0.89
MinT: COR	0.02	0.99	0.21	0.89	MinT: COR	0.01	1.00	0.88	0.99
MinT: ARMA	0.01	0.03	0.00	0.01	MinT: ARMA	0.01	0.55	0.01	0.04
MinT: ALL	0.02	0.28	0.12	0.21	MinT: ALL	0.02	0.53	0.26	0.31
$h = 3$ steps ahead									
MinT: IND	0.02	0.57	0.66	0.72	MinT: IND	0.02	0.76	0.94	0.94
MinT: EXCH	0.02	0.43	0.29	0.49	MinT: EXCH	0.02	0.14	0.13	0.16
MinT: EUCL	0.02	0.44	0.42	0.59	MinT: EUCL	0.03	0.40	0.97	0.98
MinT: COR	0.02	0.55	0.75	0.85	MinT: COR	0.02	0.39	0.98	0.98
MinT: ARMA	0.02	0.86	0.74	0.79	MinT: ARMA	0.02	0.13	0.65	0.70
MinT: ALL	0.02	0.81	0.79	0.83	MinT: ALL	0.04	0.80	0.97	0.97
$h = 6$ steps ahead									
MinT: IND	0.02	0.03	0.16	0.24	MinT: IND	0.02	0.05	0.81	0.88
MinT: EXCH	0.02	0.05	0.16	0.17	MinT: EXCH	0.02	0.03	0.04	0.11
MinT: EUCL	0.02	0.05	0.64	0.83	MinT: EUCL	0.02	0.04	0.72	0.98
MinT: COR	0.02	0.04	0.18	0.20	MinT: COR	0.02	0.04	0.31	0.88
MinT: ARMA	0.02	0.03	0.48	0.59	MinT: ARMA	0.02	0.02	0.12	0.36
MinT: ALL	0.02	0.02	0.16	0.17	MinT: ALL	0.02	0.07	0.71	0.83
$h = 12$ steps ahead									
MinT: IND	0.01	0.04	0.05	0.09	MinT: IND	0.01	0.01	0.16	0.43
MinT: EXCH	0.01	0.12	0.31	0.87	MinT: EXCH	0.01	0.03	0.07	0.48
MinT: EUCL	0.01	0.09	0.00	0.00	MinT: EUCL	0.01	0.03	0.01	0.00
MinT: COR	0.01	0.08	0.02	0.02	MinT: COR	0.01	0.02	0.01	0.00
MinT: ARMA	0.01	0.09	0.36	0.49	MinT: ARMA	0.01	0.02	0.10	0.25
MinT: ALL	0.01	0.02	0.02	0.02	MinT: ALL	0.01	0.01	0.10	0.20

Table A3. P-values from the standard Diebold and Mariano (2002) test for the top-level series (the DJIA Index). Under the null, the difference in the out-of-sample forecast errors between the model in the row and the model in the column is equal to zero. Entries in bold indicate a rejection of the null at a 5% level of significance.

Absolute errors	Base	RW	MinT	Squared errors	Base	RW	MinT
$h = 1$ step ahead							
MinT: IND	0.15	1.00	0.07	MinT: IND	0.47	1.00	0.53
MinT: EXCH	0.08	1.00	0.09	MinT: EXCH	0.08	0.99	0.09
MinT: EUCL	0.18	1.00	0.04	MinT: EUCL	0.10	1.00	0.10
MinT: COR	0.73	1.00	0.64	MinT: COR	0.86	1.00	0.91
MinT: ARMA	0.00	1.00	0.00	MinT: ARMA	0.00	1.00	0.00
MinT: ALL	0.02	0.88	0.02	MinT: ALL	0.04	1.00	0.06
$h = 3$ steps ahead							
MinT: IND	0.08	0.99	0.19	MinT: IND	0.35	1.00	0.57
MinT: EXCH	0.04	0.93	0.12	MinT: EXCH	0.04	0.89	0.09
MinT: EUCL	0.04	0.95	0.10	MinT: EUCL	0.04	0.97	0.12
MinT: COR	0.19	0.99	0.46	MinT: COR	0.21	1.00	0.49
MinT: ARMA	0.03	0.97	0.07	MinT: ARMA	0.00	0.93	0.01
MinT: ALL	0.08	0.72	0.11	MinT: ALL	0.17	0.92	0.25
$h = 6$ steps ahead							
MinT: IND	0.01	0.73	0.01	MinT: IND	0.01	0.67	0.02
MinT: EXCH	0.02	0.77	0.05	MinT: EXCH	0.03	0.69	0.05
MinT: EUCL	0.07	0.92	0.25	MinT: EUCL	0.00	0.82	0.02
MinT: COR	0.05	0.86	0.10	MinT: COR	0.00	0.73	0.01
MinT: ARMA	0.01	0.83	0.01	MinT: ARMA	0.00	0.63	0.00
MinT: ALL	0.02	0.44	0.02	MinT: ALL	0.01	0.46	0.02
$h = 12$ steps ahead							
MinT: IND	0.07	0.36	0.21	MinT: IND	0.01	0.35	0.04
MinT: EXCH	0.01	0.39	0.04	MinT: EXCH	0.01	0.41	0.04
MinT: EUCL	0.00	0.37	0.01	MinT: EUCL	0.00	0.38	0.00
MinT: COR	0.04	0.35	0.08	MinT: COR	0.00	0.34	0.00
MinT: ARMA	0.00	0.28	0.00	MinT: ARMA	0.00	0.29	0.00
MinT: ALL	0.02	0.17	0.04	MinT: ALL	0.00	0.24	0.01

Table A4. P-values from the modified Diebold and Mariano (2002) test of van Dijk and Franses (2003) for the bottom-level series (the DJIA Index constituents). Under the null, the difference in the (average) out-of-sample forecast errors between the model in the row and the model in the column is equal to zero. Entries in bold indicate a rejection of the null at a 5% level of significance.

Absolute errors	Base	RW	BU	MinT	Squared errors	Base	RW	BU	MinT
$h = 1$ step ahead									
MinT: IND	0.00	0.10	0.05	0.20	MinT: IND	0.21	0.11	0.02	0.14
MinT: EUCL	0.00	0.05	0.31	0.14	MinT: EUCL	0.12	0.34	0.22	0.39
MinT: COR	0.00	0.12	0.05	0.24	MinT: COR	0.17	0.17	0.03	0.05
MinT: ARMA	0.00	0.05	0.05	0.21	MinT: ARMA	0.43	0.05	0.04	0.12
MinT: ALL	0.02	0.18	0.08	0.18	MinT: ALL	0.50	0.03	0.05	0.10
$h = 3$ steps ahead									
MinT: IND	0.09	0.22	0.30	0.24	MinT: IND	0.28	0.10	0.28	0.27
MinT: EUCL	0.07	0.18	0.33	0.38	MinT: EUCL	0.40	0.25	0.37	0.37
MinT: COR	0.08	0.20	0.37	0.42	MinT: COR	0.41	0.29	0.22	0.16
MinT: ARMA	0.05	0.18	0.49	0.42	MinT: ARMA	0.39	0.23	0.37	0.37
MinT: ALL	0.18	0.28	0.21	0.17	MinT: ALL	0.15	0.01	0.14	0.13
$h = 6$ steps ahead									
MinT: IND	0.04	0.07	0.28	0.24	MinT: IND	0.35	0.33	0.43	0.49
MinT: EUCL	0.03	0.07	0.16	0.10	MinT: EUCL	0.37	0.37	0.42	0.49
MinT: COR	0.03	0.06	0.07	0.03	MinT: COR	0.32	0.42	0.06	0.07
MinT: ARMA	0.02	0.05	0.15	0.12	MinT: ARMA	0.33	0.40	0.29	0.34
MinT: ALL	0.03	0.03	0.48	0.50	MinT: ALL	0.32	0.44	0.33	0.35
$h = 12$ steps ahead									
MinT: IND	0.00	0.38	0.00	0.00	MinT: IND	0.00	0.48	0.00	0.00
MinT: EUCL	0.04	0.09	0.08	0.26	MinT: EUCL	0.03	0.13	0.10	0.40
MinT: COR	0.02	0.15	0.00	0.00	MinT: COR	0.02	0.19	0.00	0.00
MinT: ARMA	0.01	0.19	0.05	0.09	MinT: ARMA	0.00	0.32	0.05	0.10
MinT: ALL	0.00	0.29	0.00	0.00	MinT: ALL	0.00	0.24	0.00	0.00

Table A5. P-values from the modified Diebold and Mariano (2002) test of van Dijk and Franses (2003) for the top-level series (the DJIA Index) in the period 2015–2017. Under the null, the difference in the out-of-sample forecast errors between the model in the row and the model in the column is equal to zero. Entries in bold indicate a rejection of the null at a 5% level of significance.

Absolute errors	Base	RW	BU	MinT	Squared errors	Base	RW	BU	MinT
$h = 1$ step ahead									
MinT: IND	0.02	0.00	0.02	0.03	MinT: IND	0.01	0.01	0.01	0.01
MinT: EUCL	0.01	0.00	0.00	0.00	MinT: EUCL	0.00	0.00	0.00	0.00
MinT: COR	0.24	0.20	0.00	0.00	MinT: COR	0.04	0.02	0.00	0.00
MinT: ARMA	0.00	0.00	0.00	0.00	MinT: ARMA	0.00	0.00	0.00	0.00
MinT: ALL	0.28	0.24	0.39	0.26	MinT: ALL	0.39	0.48	0.36	0.24
$h = 3$ steps ahead									
MinT: IND	0.00	0.00	0.00	0.00	MinT: IND	0.00	0.01	0.01	0.02
MinT: EUCL	0.00	0.06	0.08	0.28	MinT: EUCL	0.00	0.01	0.00	0.01
MinT: COR	0.01	0.00	0.00	0.00	MinT: COR	0.13	0.00	0.00	0.00
MinT: ARMA	0.00	0.00	0.00	0.00	MinT: ARMA	0.00	0.00	0.00	0.00
MinT: ALL	0.17	0.21	0.11	0.06	MinT: ALL	0.14	0.31	0.20	0.08
$h = 6$ steps ahead									
MinT: IND	0.00	0.06	0.02	0.05	MinT: IND	0.00	0.05	0.02	0.05
MinT: EUCL	0.00	0.03	0.00	0.00	MinT: EUCL	0.00	0.01	0.00	0.00
MinT: COR	0.09	0.00	0.00	0.00	MinT: COR	0.08	0.00	0.00	0.00
MinT: ARMA	0.00	0.00	0.00	0.00	MinT: ARMA	0.00	0.00	0.00	0.00
MinT: ALL	0.16	0.12	0.15	0.07	MinT: ALL	0.14	0.13	0.16	0.07
$h = 12$ steps ahead									
MinT: IND	0.00	0.04	0.02	0.06	MinT: IND	0.00	0.08	0.04	0.12
MinT: EUCL	0.00	0.00	0.00	0.00	MinT: EUCL	0.00	0.00	0.00	0.00
MinT: COR	0.32	0.00	0.00	0.00	MinT: COR	0.44	0.00	0.00	0.00
MinT: ARMA	0.00	0.00	0.00	0.00	MinT: ARMA	0.00	0.00	0.00	0.00
MinT: ALL	0.03	0.47	0.42	0.18	MinT: ALL	0.07	0.26	0.20	0.08

Table A6. P-values from the modified Diebold and Mariano (2002) test of van Dijk and Franses (2003) for the top-level series (the S&P500 Index). Under the null, the difference in the out-of-sample forecast errors between the model in the row and the model in the column is equal to zero. Entries in bold indicate a rejection of the null at a 5% level of significance.

Appendix B: Additional figures

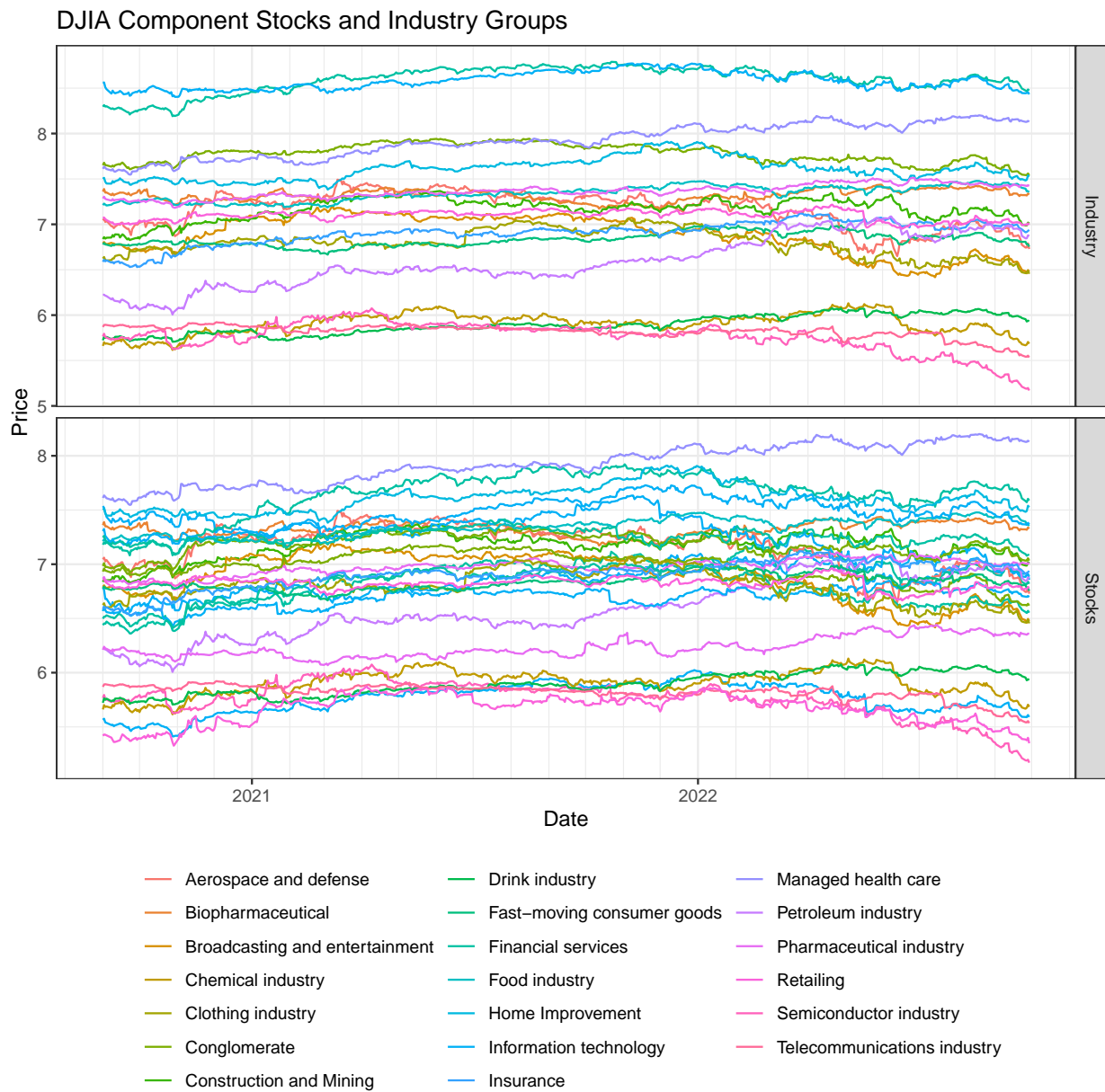


Figure B1. Aggregated series from the hierarchical structure implied by the Industry-based clustering. The upper plot shows the time series for the industry aggregates and the bottom plot shows the DJIA Index constituents. The colors in both plots correspond to different industries. Prices are shown on a logarithmic scale.

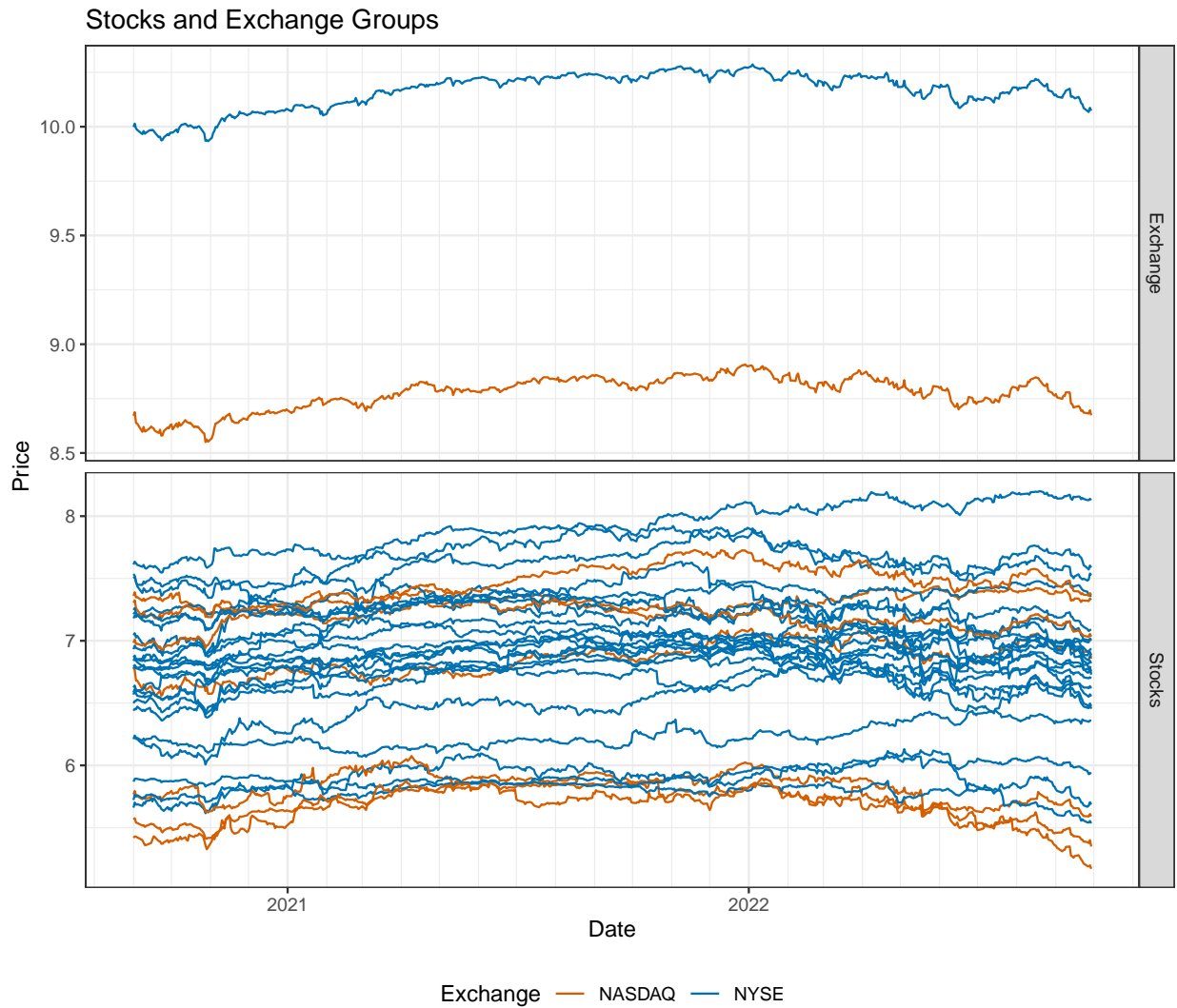


Figure B2. Aggregated series from the hierarchical structure implied by the Exchange-based clustering. The upper plot shows the time series for the exchange market aggregates and the bottom plot shows the DJIA Index constituents. Most of the stocks included in the DJIA Index during the period 2020-2022 are traded on the NYSE (blue color). The colors in both plots correspond to different markets. Prices are shown on a logarithmic scale.

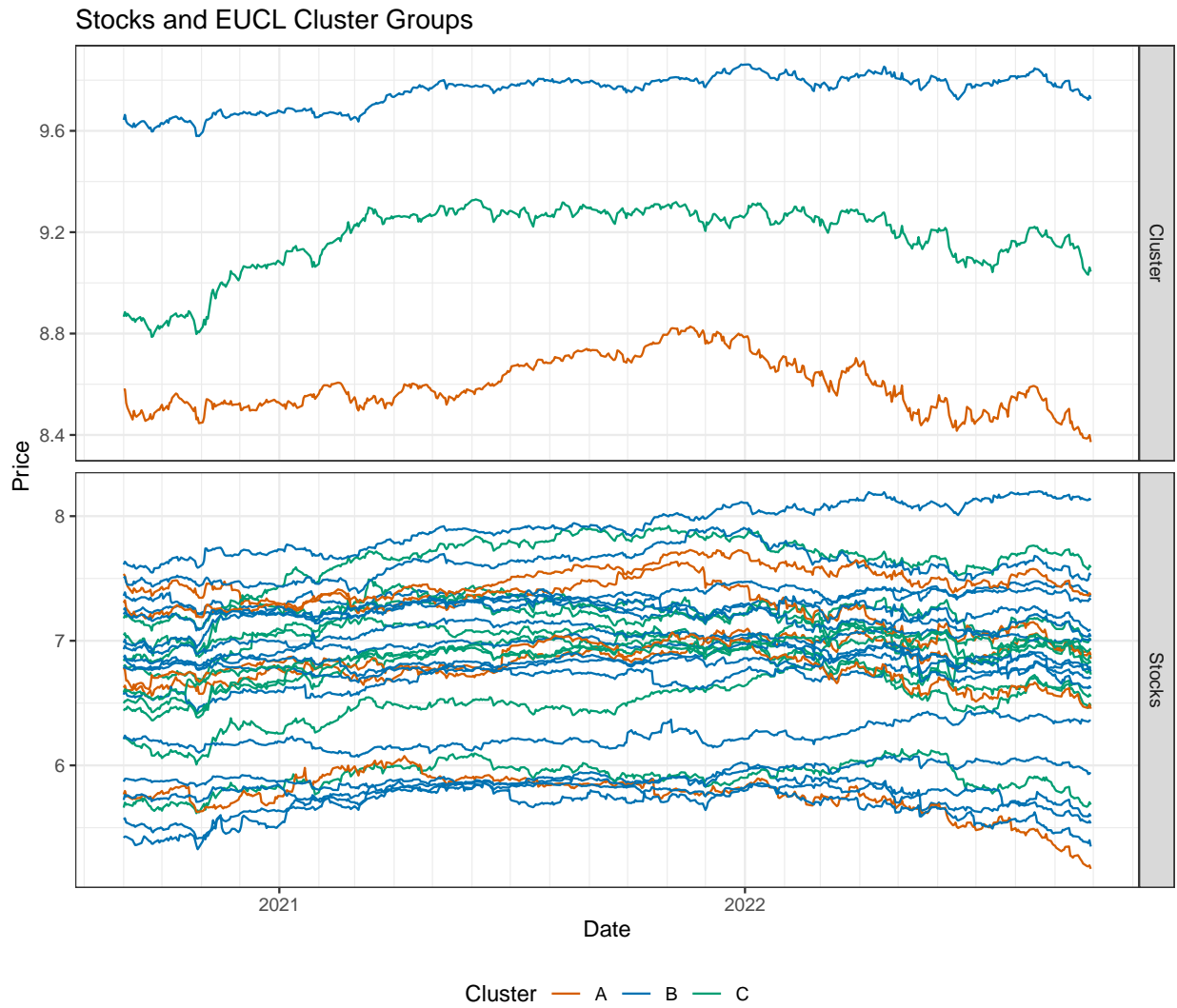


Figure B3. Aggregated series from the hierarchical structure implied by the EUCL-based PAM clustering algorithm. The upper plot shows the time series for the three returns-based clusters and the bottom plot shows the DJIA Index constituents. The colors in both plots correspond to different clusters. Prices are shown on a logarithmic scale.

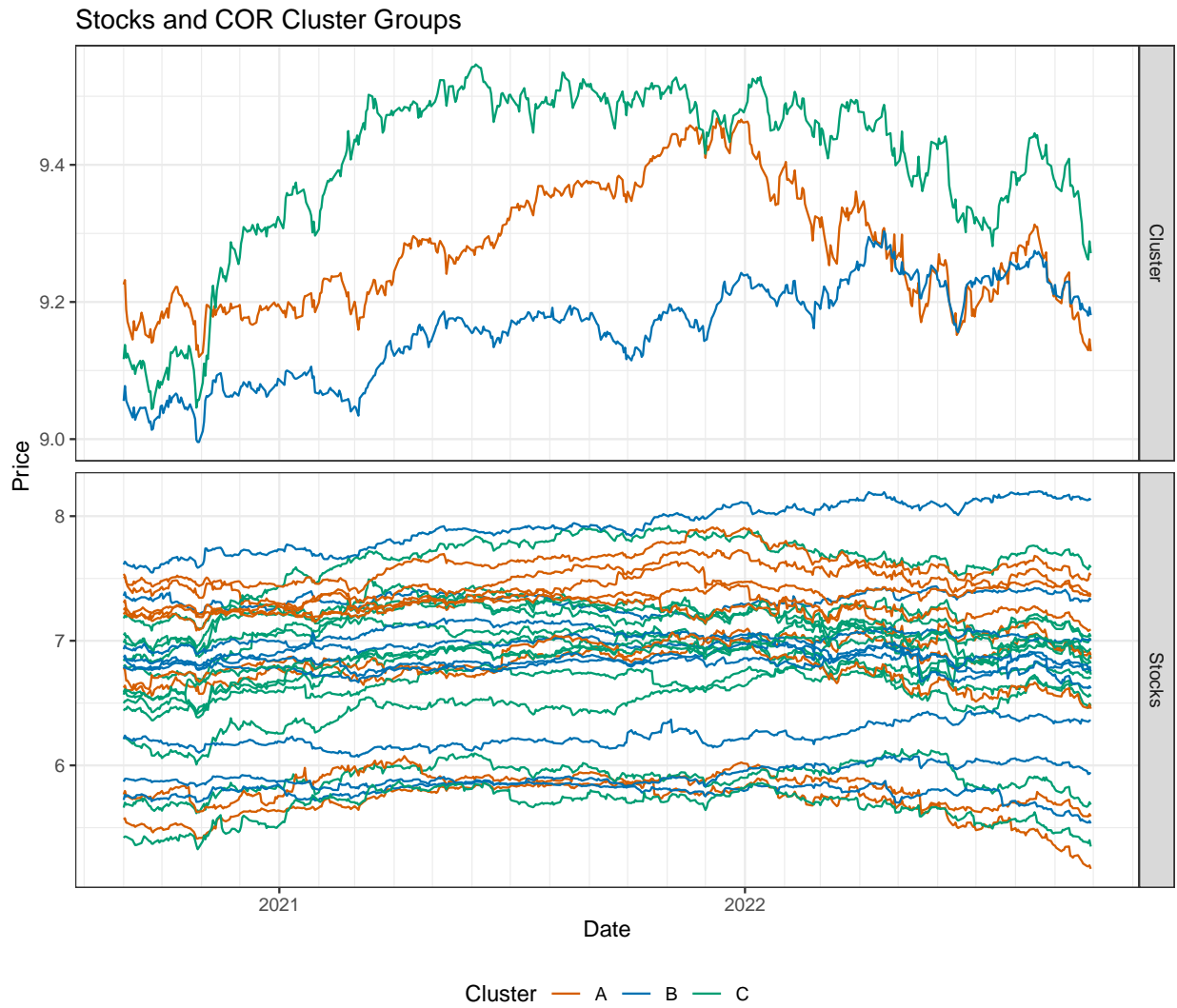


Figure B4. Aggregated series from the hierarchical structure implied by the COR-based PAM clustering algorithm. The upper plot shows the time series for the three correlation-based clusters and the bottom plot shows the DJIA Index constituents. The colors in both plots correspond to different clusters. Prices are shown on a logarithmic scale.

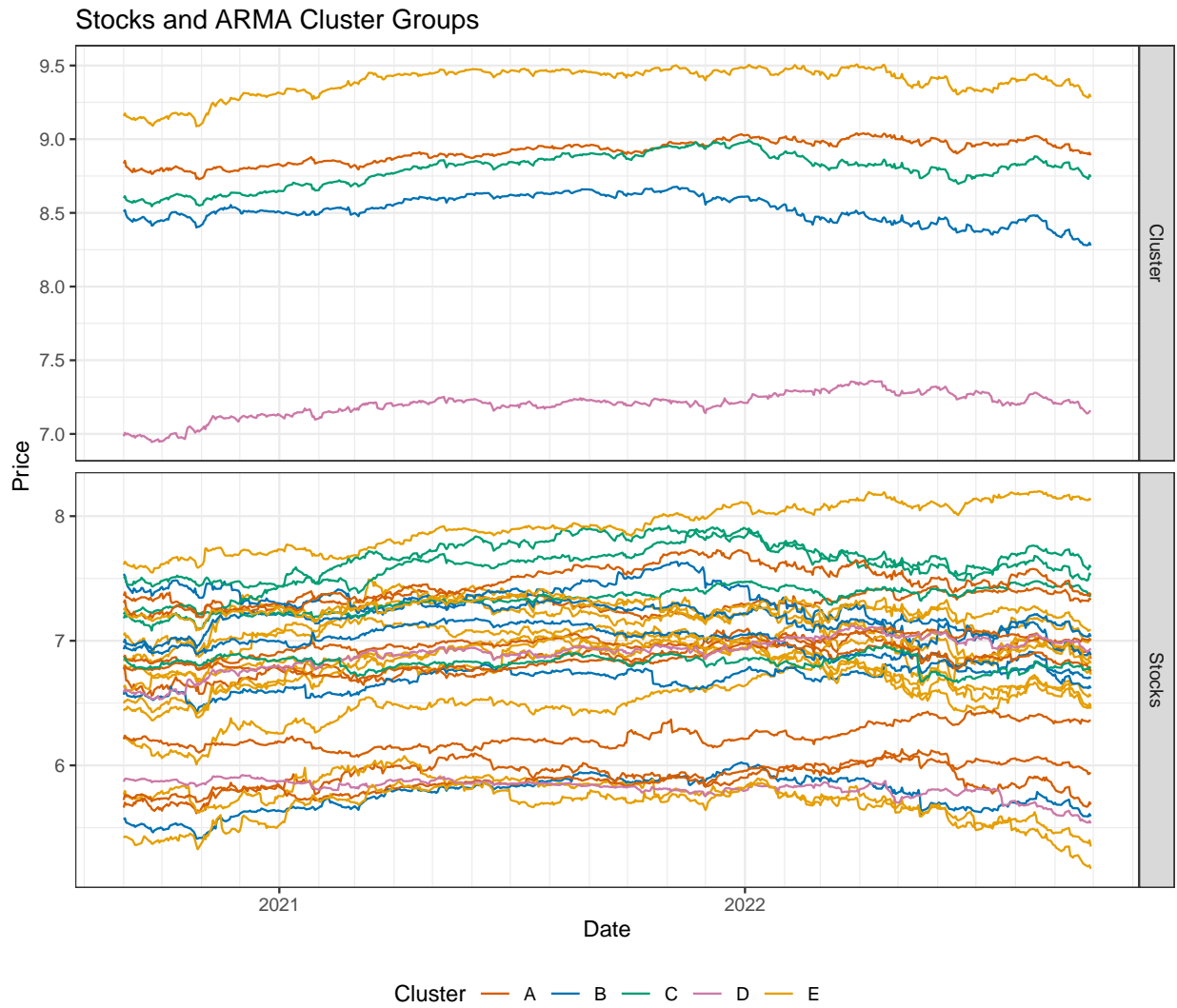


Figure B5. Aggregated series from the hierarchical structure implied by the ARIMA-based PAM clustering algorithm. The upper plot shows the time series for the five ARIMA-based clusters and the bottom plot shows the DJIA Index constituents. The colors in both plots correspond to different clusters. Prices are shown on a logarithmic scale.



Figure B6. Average Silhouette Width with EUCL-based PAM partition

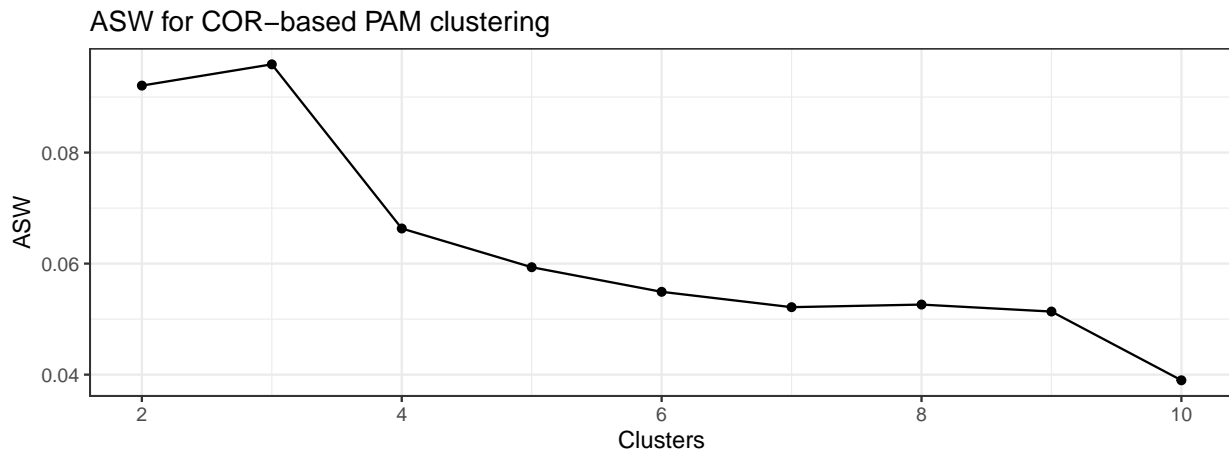


Figure B7. Average Silhouette Width with COR-based PAM partition

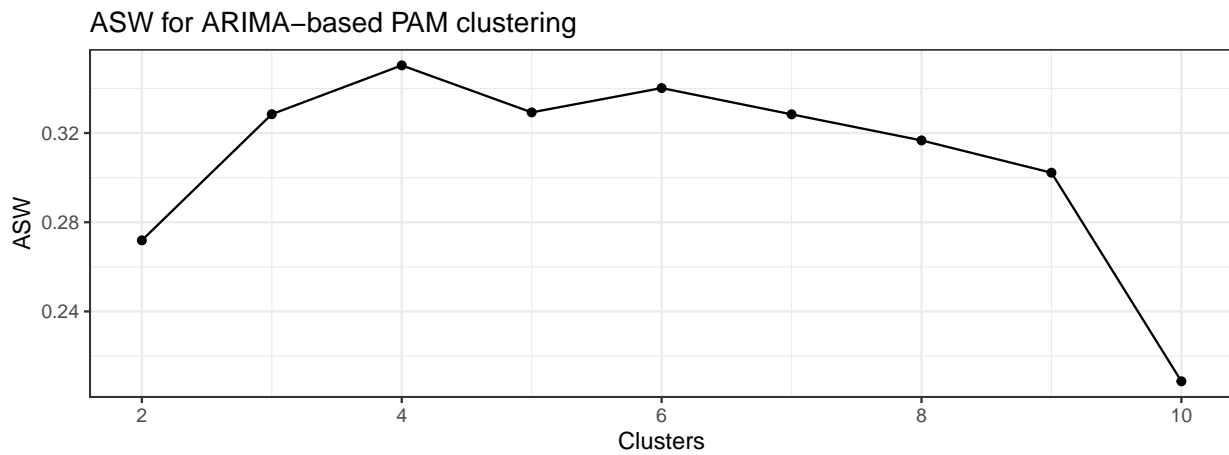


Figure B8. Average Silhouette Width with ARIMA-based PAM partition

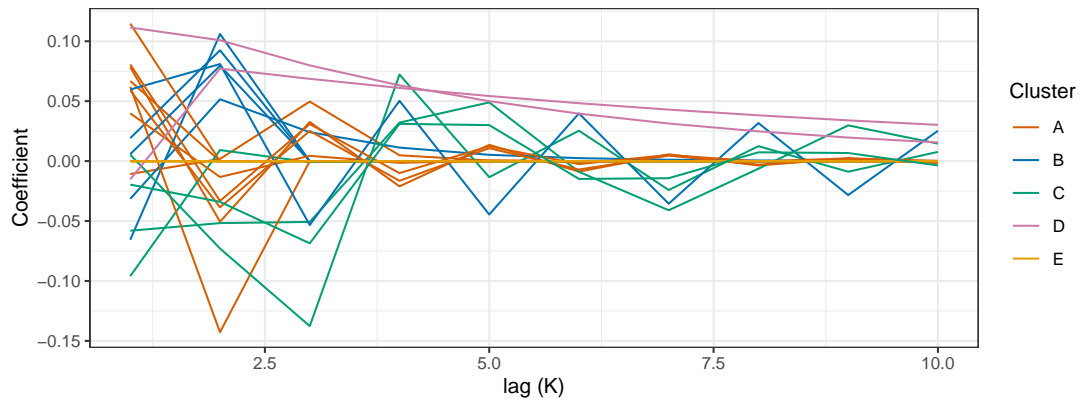


Figure B9. $AR(\infty)$ weights of the clustered time series