

High Dimensional Forecast Combinations Under Latent Structures

Zhentao Shi, Liangjun Su and Tian Xie

October 20, 2020

Abstract

This paper proposes a novel high dimensional forecast combination estimator in the presence of many forecasts and potential latent group structures. The new algorithm, which we call ℓ_2 -relaxation, minimizes the squared ℓ_2 -norm of the weight vector subject to a relaxed version of the first-order conditions, instead of minimizing the mean squared forecast error as those standard optimal forecast combination procedures. A proper choice of the tuning parameter achieves bias and variance trade-off, and incorporates as special cases the simple average (equal-weight) strategy and the conventional optimal weighting scheme. When the variance-covariance (VC) matrix of the individual forecast errors exhibits latent group structures — a block equicorrelation matrix plus a VC for idiosyncratic noises, ℓ_2 -relaxation delivers combined forecasts with roughly equal within-group weights. Asymptotic optimality of the new method is established by exploiting the duality between the sup-norm restriction and the high-dimensional sparse ℓ_1 -norm penalization. Excellent finite sample performance of our method is demonstrated in Monte Carlo simulations. Its wide applicability is highlighted in three real data examples concerning empirical applications of microeconomics, macroeconomics, and finance.

Key Words: Factor models; forecast combination puzzle; high dimension; Lasso; latent group; machine learning; optimality.

JEL Classification: C22, C53, C55

Shi acknowledges financial support from the Research Grants Council (RGC) No.14500118. Su acknowledges financial support from Tsinghua University. Xie's research is supported by the Natural Science Foundation of China (71701175), the Chinese Ministry of Education Project of Humanities and Social Sciences (17YJC790174), and the Fundamental Research Funds for the Central Universities.

Address correspondence: Zhentao Shi: zhentaoshi@cuhk.edu.hk, Department of Economics, 928 Esther Lee Building, the Chinese University of Hong Kong, Shatin, New Territories, Hong Kong SAR, China; Tel: (852) 3943-1432; Fax (852) 2603-5805. Liangjun Su: sulj@sem.tsinghua.edu.cn, School of Economics and Management, Tsinghua University, Beijing, China; Tel: (86 10) 6278-9506. Tian Xie: xietian001@hotmail.com, College of Business, Shanghai University of Finance and Economics, China.

1 Introduction

Forecast has been one of the fundamental tasks since the dawn of econometrics. In his paper presented at the inaugural meeting of the Econometric Society, Cowles (1933) found little evidence of forecast ability to the stock market among dozens of financial services and publications. As part of the postwar development of national accounting and statistics, Theil (1992, originally 1955) examined accuracy of forecasting models of macroeconomic indicators compiled by the Economic Commission for Europe. Besides quantitative economic data analysis, forecast is the backbone of macroeconomic theory, under the names of *adaptive expectation*, *rational expectation*, and so on. Individuals and firms rely on their perception of the future to make decisions about consumption, investment and production. Monetary and fiscal authorities communicate their future targets to stabilize price levels and unemployment rates. Well-functioning financial markets and economies count on prediction, confidence, and interaction of market participants and regulators.

Thanks to information technology advancements in recent decades, the costs of collecting market forecaster data and running predictive algorithms drop tremendously. We have fastforwarded into the era of big data, when machine learning and artificial intelligence profoundly change the routines of our daily life and the ways we conduct economic research. While we indulge in the rich-data environment and the awesome computing power which those pioneer econometricians could only envy, we must search for effective approaches of ensembling disaggregate information scattered in many individual forecasters or forecasting models in order to distill signals out of noises. In this paper, we seek enlightenment from Bates and Granger (1969)'s seminal work of forecast combination, which has become part of the core technical toolkits of modern forecasting practices.

We consider high dimensional forecast combinations in the presence of many forecasts. Suppose that y_{t+1} is an outcome variable and there are N forecasts, $\{f_{it}\}_{i \in [N]}$, available at time t for y_{t+1} , where $t \in [T] := \{1, 2, \dots, T\}$, $[N] := \{1, 2, \dots, N\}$ and “:=” signifies definition. Let $\mathbf{f}_t = (f_{1t}, \dots, f_{Nt})'$. We are interested in finding an $N \times 1$ weight vector $\mathbf{w} = (w_i)_{i \in [N]}$ to form a linear combination $\mathbf{w}'\mathbf{f}_t$ whose mean squared forecast error (MSFE) is minimized. One way to estimate \mathbf{w} is to run the restricted least squares (RLS):

$$\min_{\mathbf{w} \in \mathbb{R}^N} \frac{1}{2T} \sum_{t=1}^T (y_{t+1} - \mathbf{w}'\mathbf{f}_t)^2 \quad \text{subject to } \mathbf{w}'\mathbf{1}_N = 1, \quad (1.1)$$

where $\mathbf{1}_N$ is an $N \times 1$ vector of ones. Alternatively, we collect the forecast error $e_{it} = y_{t+1} - f_{it}$, compute its sample variance-covariance (VC) $\widehat{\Sigma} := T^{-1} \sum_{t=1}^T \mathbf{e}_t \mathbf{e}_t'$ with $\mathbf{e}_t = (e_{it})_{i \in [N]}$, and follow Bates and Granger (1969) to work with the following minimization problem:

$$\min_{\mathbf{w} \in \mathbb{R}^N} \frac{1}{2} \mathbf{w}' \widehat{\Sigma} \mathbf{w} \quad \text{subject to } \mathbf{w}'\mathbf{1}_N = 1. \quad (1.2)$$

Denote the solution to the above constrained optimization problem as $\widehat{\mathbf{w}}^{\text{BG}}$. When $\widehat{\Sigma}$ is invertible, the explicit solution is

$$\widehat{\mathbf{w}}^{\text{BG}} = (\mathbf{1}'_N \widehat{\Sigma}^{-1} \mathbf{1}_N)^{-1} \widehat{\Sigma}^{-1} \mathbf{1}_N. \quad (1.3)$$

The two formulations in (1.2) and (1.1) are numerically equivalent, as demonstrated by Granger and Ramanathan (1984). Apparently, the requirement of the invertibility of $\widehat{\Sigma}$ is not innocuous in high dimensional settings, and in fact $\widehat{\Sigma}$ must be singular when $N > T$.

Despite the MSFE-optimality of the above Bates and Granger's combined forecasts, the optimal weights often do not yield the best forecast in empirical work when they are replaced with their sample

estimates; rather, the simple average (namely, *equal-weight* forecast combination) often performs better. This empirical fact has been best known as the “*forecast combination puzzle*,” which was first noted by Clemen (1989) and formally named by Stock and Watson (2004). In particular, Clemen (1989) reviews the literature up to the late 1980s and notes that over a large number of papers averaging forecasts appear to be a more robust procedure than optimal combination. A reasonable explanation suggests that the errors on the estimation of the weights can be large and thus dominate the gains from the use of optimal combination; see, e.g., Smith and Wallis (2009) and Claeskens et al. (2016). See Timmermann (2006) for an excellent review of forecast combination.

A lesson learned from this literature is that it is unwise to include all possible variables. Limiting the number of unknown parameters can help reduce estimation error. This has led to the adoption of various shrinkage and variable selection techniques. For example, Elliott et al. (2013) propose a complete subset regression (CSR) approach to forecast combinations by using equal weights to combine forecasts based on the same number of predictive variables. They show that in many cases subset regression combinations amount to a form of shrinkage that is more general than the conventional variable-by-variable shrinkage implied by ridge regression. In contrast, Diebold and Shin (2019) bring forth the partially egalitarian Lasso (peLASSO) procedures that discard some forecasts and then select and shrink the remaining forecasts toward the equal weights. In essence, both Elliott et al. (2013) and Diebold and Shin (2019) assign two distinct weights to two subsets of models: zero weight to a large subset of the models and equal or roughly equal weights to the remaining subset of the models.

In this paper, we extend the ideas of Elliott et al. (2013) and Diebold and Shin (2019) and propose a new shrinkage technique for forecast combinations. Specifically, we propose to minimize the squared ℓ_2 -norm of the weight vector subject to a relaxed version of the first order conditions from the minimization problem in (1.2), yielding the ℓ_2 -relaxation problem. The strategy is similar in spirit to the ℓ_1 -relaxation in Dantzig selector *a la* Candes and Tao (2007). Interestingly, the ℓ_2 -relaxed optimal forecasts incorporate both the simple average (equal-weight) strategy and the conventional optimal weighting scheme as special cases by setting the tuning parameter to be sufficiently large or zero, respectively. A proper choice of the tuning parameter can achieve the usual bias and variance trade-off and deliver combined forecasts with roughly equal groupwise weights when the variance-covariance (VC) matrix of the individual forecasts exhibit a certain latent group structure. This is consistent with the intuition that one should assign the same weights to all individuals within the same group and potentially distinct weights to individuals in different groups when the forecast error VC matrix exactly exhibits a block equicorrelation structure. When the VC matrix is contaminated with a noisy component, we show that the resultant ℓ_2 -relaxed weights are close to the infeasible groupwise equal weights.

The approximate latent group structure is not a man-made artifact. It is inherent in many forecast combination problems. For example, it emerges when the forecasting errors exhibit a factor structure as in Hsiao and Wan (2014) and Chan and Pauwels (2018) and the factor loadings are directly governed by certain latent group structures or are approximable by a few values. It is also present if one considers forecast combinations based on a large number of forecast models (say, $2^{10} = 1024$) with a fixed number of predictive regressors (say, $p = 10$). In the latter case, we argue that the p regressors serve as a part of the “latent” factors.

We develop original asymptotic results to support this new ℓ_2 -relaxation estimation method. Noticing the duality between the sup-norm constraint and the ℓ_1 -penalization of Lasso (Tibshirani, 1996), we first develop the asymptotic convergence in its dual problem (see (2.6) in the next section) which resembles Lasso in view of its ℓ_1 -penalty, instead of directly working with the primal problem of ℓ_2 -relaxation (see (2.5)) which is a linearly constrained quadratic optimization. Studies of high

dimensional regression have prepared a set of inequalities for Lasso to handle sparse regressions. We sharpen these techniques in our context to cope with groupwise sparsity. Once the convergence of the high dimensional parameters in the dual problem is established, the convergence of the combination weights proceeds immediately, and then follows the asymptotic optimality of ℓ_2 -relaxation.

We assess the finite sample behavior of ℓ_2 -relaxation in Monte Carlo simulations. Compared with the oracle estimator and popular off-the-shelf machine learning estimators, ℓ_2 -relaxation performs well under various data generating processes (DGPs). We further evaluate its empirical accuracy in three real data examples covering box office prediction, inflation forecast by surveyed experts, and stock market volatility forecast based on regressions. These examples demonstrate wide applicability of ℓ_2 -relaxation.

Literature Review. It is worth mentioning that various shrinkage and machine learning techniques have been applied to the forecasting literature since the pioneer work of Tibshirani (1996). See, e.g., Li and Chen (2014), Conflitti et al. (2015), Konzen and Ziegelmann (2016), Stasinakis et al. (2016), Bayer (2018), Wilms et al. (2018), Kotchoni et al. (2019), Coulombe et al. (2020), and Roccazzella et al. (2020). We complement these works by considering an ℓ_2 -relaxation of the regularized weights estimation problem, which exhibits certain optimality properties when the forecast error VC matrix can be decomposed into the sum of a low-rank matrix and a VC of idiosyncratic shocks.

Our paper is also related to the recent literature on latent group structures in panel data analysis; see, e.g., Bonhomme and Manresa (2015), Su et al. (2016), Bonhomme et al. (2017), Su and Ju (2018), Su et al. (2019), Vogt and Linton (2017), and Vogt and Linton (2020). Except Bonhomme and Manresa (2015) and Bonhomme et al. (2017), all these previous studies focus on the recovery of the latent group structures. In this paper, although we assume that a community or latent group structure lies in the dominant term in the forecast error VC matrix, we do not attempt to recover the membership for each individual forecast.

Lastly, our paper adds to the vast literature on portfolio optimization as well; see Ledoit and Wolf (2004), Disatnik and Katz (2012), Fan et al. (2012), Fan et al. (2016), Ledoit and Wolf (2017), and Ao et al. (2019), among many others. In particular, Ledoit and Wolf (2004) use Bayesian methods for shrinking the sample correlation matrix to an equicorrelated target and show that this helps select portfolios with low volatility compared to those based on the sample correlation; and Ledoit and Wolf (2017) promote a nonlinear shrinkage estimator that is more flexible than previous linear shrinkage estimators and has just the right number of free parameters. Apparently, our method can also be used to estimate the optimal portfolios.

Organization. The rest of the paper is organized as follows. In Section 2, we introduce the ℓ_2 -relaxation primal problem and its dual problem, and discuss their attributes without imposing structures on the VC matrix of the forecast errors. Section 3 considers the latent group structures in the VC matrix and develops the finite sample numerical properties of the optimization problems. In Section 4, we study the asymptotic behaviors of the estimators in both the dual and primal problems and establish the asymptotic optimality of the ℓ_2 -relaxed estimator of the combination weights. Section 5 reports Monte Carlo simulation results. The new method is applied to three datasets in Section 6. We provide the proofs of all theoretical results Appendix A. Additional numerical results are contained in Appendix B.

Notation. For a random variable x , we write its population mean as $E[x]$; for a sample (x_1, \dots, x_T) , we write its sample mean as $\mathbb{E}_T[x_t] := T^{-1} \sum_{t=1}^T x_t$. A plain b denotes a scalar, a boldface lowercase \mathbf{b} denotes a vector, and a boldface uppercase \mathbf{B} denotes a matrix. The ℓ_1 -norm and ℓ_2 -norm of $\mathbf{b} = (b_1, \dots, b_n)'$ are denoted as $\|\mathbf{b}\|_1 := \sum_{i=1}^n |b_i|$ and $\|\mathbf{b}\|_2 := (\sum_{i=1}^n b_i^2)^{1/2}$, respectively. For a generic $n \times m$ matrix \mathbf{B} , we denote \mathbf{B}_i as the i -th row ($1 \times m$ vector), and \mathbf{B}_j as

the j -th column ($n \times 1$ vector). $\phi_{\max}(\cdot)$ and $\phi_{\min}(\cdot)$ denote the largest and smallest eigenvalues of a real symmetric matrix, respectively. Define the spectral norm, sup-norm, and maximum column ℓ_2 matrix norm as $\|\mathbf{B}\|_{\text{sp}} := \phi_{\max}^{1/2}(\mathbf{B}'\mathbf{B})$, $\|\mathbf{B}\|_{\infty} := \max_{i \leq n, j \leq m} |b_{ij}|$, and $\|\mathbf{B}\|_{c2} := \max_{j \leq m} \|\mathbf{B}_{\cdot j}\|_2$, respectively. $\mathbf{0}_n$ and $\mathbf{1}_n$ are $n \times 1$ vectors of zeros and ones, respectively, and \mathbf{I}_n is the $n \times n$ identity matrix. “w.p.a.1” is short for “with probability approaching one” in asymptotic statements. We write $a \asymp b$ when both a/b and b/a are stochastically bounded. Let $a \wedge b = \min\{a, b\}$.

The cross-sectional units are indexed by $i \in [N] := \{1, \dots, N\}$. For a generic index set $\mathcal{G} \subset [N]$, we denote $|\mathcal{G}|$ as the cardinality of \mathcal{G} , and $\mathbf{b}_{\mathcal{G}} = (b_i)_{i \in \mathcal{G}}$ as the $|\mathcal{G}|$ -dimensional subvector. We call a vector of *similar sign* if no pair of its elements takes opposite signs. Formally, an n -vector \mathbf{b} is of similar sign if

$$\mathbf{b} \in \mathbb{S}^n := \{\mathbf{b} \in \mathbb{R}^n : b_i b_j \geq 0 \text{ for all } i, j \in [N]\}.$$

Let $\mathcal{G}_1, \dots, \mathcal{G}_K$ be a partition of $[N]$, and denote $N_k = |\mathcal{G}_k|$. We call a generic N -vector \mathbf{b} of *similar sign within all groups* if $\mathbf{b} \in \mathbb{S}^{\text{all}} := \{\mathbf{b} \in \mathbb{R}^N : \mathbf{b}_{\mathcal{G}_k} \in \mathbb{S}^{N_k} \text{ for all } k \in [K]\}$, and define $\tilde{\mathbb{S}}^{\text{all}} := \{\mathbf{b} \in \mathbb{S}^{\text{all}} : \mathbf{1}'_N \mathbf{b} = 0\}$ with a further restriction that the elements in \mathbf{b} add up to 0.

2 ℓ_2 -Relaxation for the Optimal Forecast Combination

In this section, we introduce the idea of ℓ_2 -relaxation to the classical forecast combination problem. We first discuss the low-dimensional case, and then proceed to the high-dimensional case.

2.1 ℓ_2 -Relaxation in the Low Dimensional Case

To solve the constrained optimization problem in (1.2), we rewrite it as the unconstrained Lagrangian problem:

$$\frac{1}{2} \mathbf{w}' \widehat{\Sigma} \mathbf{w} + \gamma (\mathbf{w}' \mathbf{1}_N - 1), \quad (2.1)$$

where γ is the Lagrangian multiplier. If $\widehat{\Sigma}$ is invertible, the explicit solution is given in (1.3). Nevertheless, $\widehat{\Sigma}$ is frequently non-invertible in high-dimensional settings. In this case, (1.2) has multiple solutions, all of which satisfy the Kuhn-Karush-Tucker (KKT) conditions:

$$\begin{aligned} \widehat{\Sigma} \mathbf{w} + \gamma \mathbf{1}_N &= \mathbf{0}_N \\ \mathbf{w}' \mathbf{1}_N - 1 &= 0. \end{aligned} \quad (2.2)$$

The non-unique solution due to the singularity of $\widehat{\Sigma}$ is parallel to the familiar multi-collinearity issue in ordinary least squares (OLS) regression. To motivate the ℓ_2 -relaxation, we take a digression and consider the generic OLS regression in the following example.

Example 1 *The OLS seeks a parameter $\boldsymbol{\theta} \in \mathbb{R}^p$ that minimizes the sum of squared residuals (SSR) $\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2$, where \mathbf{y} is a T -vector of dependent variable and \mathbf{X} is a $T \times p$ matrix of independent variables. When the columns in \mathbf{X} are collinear, any solution that satisfies the first order condition $\mathbf{X}'\mathbf{X}\boldsymbol{\theta} = \mathbf{X}'\mathbf{y}$ minimizes the SSR. One solution to the collinearity problem is provided by the well-known ridge regression (Tikhonov, 1977): $\widehat{\boldsymbol{\theta}}^{\text{ridge}} = \min_{\boldsymbol{\theta} \in \mathbb{R}^p} \left\{ \frac{1}{2T} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_2^2 \right\}$, where λ is a tuning parameter. Alternatively, one can set up a criterion to pick one out of these multiple solutions. For example, under the ℓ_2 -norm, one of the most popular criteria, the problem can be written as*

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^p} \frac{1}{2} \|\boldsymbol{\theta}\|_2^2 \quad \text{subject to} \quad \mathbf{X}'\mathbf{X}\boldsymbol{\theta} = \mathbf{X}'\mathbf{y}. \quad (2.3)$$

Unlike the ridge regression, no tuning parameter is involved in (2.3). Apparently, when $\mathbf{X}'\mathbf{X}$ is non-singular, the solution to the last minimization problem is given by the OLS solution $\widehat{\boldsymbol{\theta}}^{\text{OLS}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ which is the only feasible point satisfying the constraint in (2.3).

Following the spirit of (2.3), a unique solution to (1.2) can be found by

$$\min_{(\mathbf{w}, \gamma) \in \mathbb{R}^{N+1}} \frac{1}{2} \|\mathbf{w}\|_2^2 \quad \text{subject to} \quad \mathbf{w}'\mathbf{1}_N = 1 \text{ and } \widehat{\boldsymbol{\Sigma}}\mathbf{w} + \gamma\mathbf{1}_N = \mathbf{0}_N. \quad (2.4)$$

The above minimization problem can be solved regardless of the invertibility of $\widehat{\boldsymbol{\Sigma}}$, and (1.3) is its solution when $\widehat{\boldsymbol{\Sigma}}$ is indeed invertible.

2.2 ℓ_2 -Relaxation in the High Dimensional Case

Potential problems arise when N is large relatively to the time dimension T . “High dimensional” here means that the number of unknown parameters, in our context N , is comparable to the sample size T . For example, we allow $N/T \rightarrow c \in (0, \infty]$ as $(N, T) \rightarrow \infty$.

Consider the case where N is of similar magnitude to T but $N < T$, say $N = 80$ and $T = 100$. Even if $\widehat{\boldsymbol{\Sigma}}$ is non-singular, a few small sample eigenvalues of $\widehat{\boldsymbol{\Sigma}}$ are likely to be close to zero, leading to a numerically unstable solution from (1.3). The numerical instability is due to the fact that the weights are solely decided by $\widehat{\boldsymbol{\Sigma}}$ via the $(N + 1)$ -equation linear system in (2.4).

An idea to stabilize the numerical solution is expanding the feasible set. Inspired by the Dantzig selector of Candès and Tao (2007) and the relaxed empirical likelihood of Shi (2016), we consider relaxing the sup-norm of the KKT condition as follows:

$$\min_{(\mathbf{w}, \gamma) \in \mathbb{R}^{N+1}} \frac{1}{2} \|\mathbf{w}\|_2^2 \quad \text{subject to} \quad \mathbf{w}'\mathbf{1}_N = 1 \text{ and } \|\widehat{\boldsymbol{\Sigma}}\mathbf{w} + \gamma\mathbf{1}_N\|_\infty \leq \tau, \quad (2.5)$$

where τ is a tuning parameter to be specified by the user. We call the problem in (2.5) the ℓ_2 -relaxation primal problem, and denote the solution to (2.5) as $\widehat{\mathbf{w}} = \widehat{\mathbf{w}}_\tau$, where the dependence of $\widehat{\mathbf{w}}$ on τ is often suppressed for notational conciseness.

Constraints in (2.5) are feasible for any $\tau \geq 0$. The solution to (2.5) is unique because the objective is a strictly convex function and the feasible set is a closed convex set. Therefore, while (1.2) may have multiple solutions, the objective function in (2.5) selects in the feasible set the solution with the smallest ℓ_2 -norm. With modern convex optimization modeling languages and open-source convex solvers, the quadratic optimization with constraints in (2.5) can be handled with ease even when N is in hundreds or thousands. Proprietary convex solvers can also be called upon for further speed gain in numerical operations; see Gao and Shi (2020).

Remark 2.1. The tuning parameter τ plays a crucial role for the ℓ_2 -relaxation problem. Interestingly, the penalized scheme in (2.5) incorporates the two extremes, simple average and optimal weighting, as special cases.

- (a) If $\tau = 0$, the solution $\widehat{\mathbf{w}}$ is characterized by the KKT conditions in (2.2); if, in addition, $\widehat{\boldsymbol{\Sigma}}$ is invertible, the unique solution is then given by the optimal weighting $\widehat{\mathbf{w}}^{\text{BG}}$ defined in (1.3).
- (b) When $\tau > 0$, it relaxes the high dimensional component of the KKT condition in (2.2). If τ is sufficiently large, say $\tau \geq \max_{i \in [N]} |\widehat{\boldsymbol{\Sigma}}_i \mathbf{1}_N|/N$, the second constraint in (2.5) will not play a role in minimizing the objective function and it is easy to see that $\widehat{\mathbf{w}} = N^{-1}\mathbf{1}_N$ is the unique solution to \mathbf{w} in (2.5).¹ That is, the simple average is optimal for the primal problem in (2.5)

¹If $\tau = \max_{i \in [N]} |\widehat{\boldsymbol{\Sigma}}_i \mathbf{1}_N|/N$, then $\widehat{\gamma} = 0$ is also the unique solution to γ .

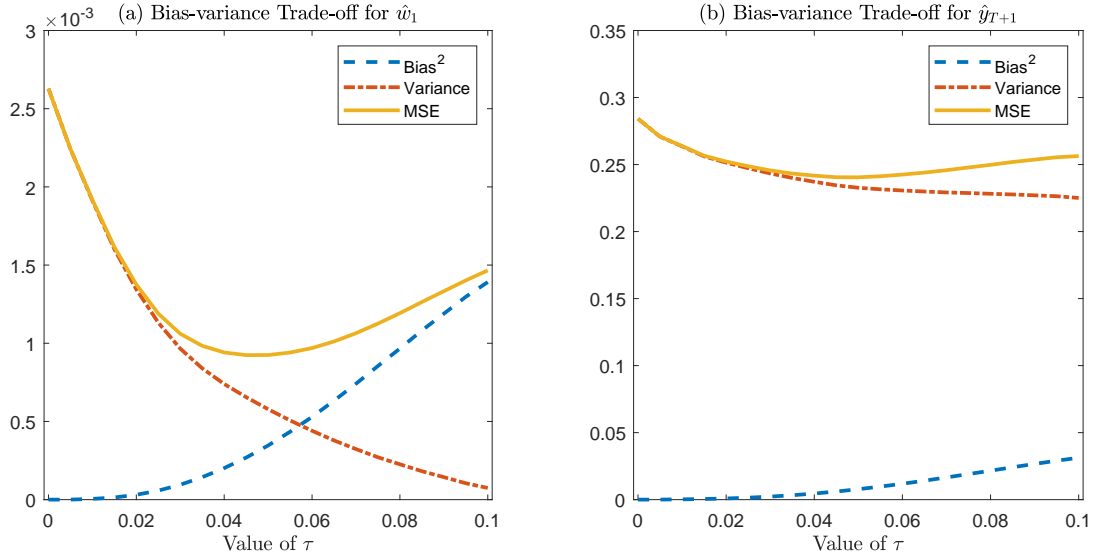


Figure 1: The Bias and Variance Trade-off under Different τ Values

provided τ is sufficiently large.

- (c) The optimal weighting strategy in (1.2) often performs unsatisfactorily in practice because the estimation of the optimal weights yield biases and large estimation errors. If the tuning parameter τ is chosen properly, the ℓ_2 -relaxation can achieve the right balance between the optimal weighting and the simple average by exploring their bias-variance trade-off. Figure 1 demonstrates the bias and variance trade-off under a range of τ , where the DGP is described in Section B.1 in the appendix. This graphs plots the squared bias, variance and MSE of the first element ($\hat{w}_1 = \hat{w}_{1\tau}$) of $\hat{\mathbf{w}}_\tau$ and those of the one-step-ahead forecast \hat{y}_{T+1} with $T = 100$ as a function of τ . As can be seen clearly from the figure, both the ℓ_2 -relaxation weight estimator and the one-step-ahead forecast associated with a small value of τ tend to have small biases but large variances whereas those associated with a large value of τ tend to have large biases but small variances. In the middle, there is a large range of values for τ where the combined forecast yields MSFE that is smaller than either that of the simple average estimator (attainable for sufficiently large τ) or that of the Bates-Granger optimally combined estimator (attainable for $\tau = 0$).

The primal problem in (2.5) is accompanied by the dual problem as stated in the following lemma.

Lemma 1 *The dual problem of (2.5) is*

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^N} \left\{ \frac{1}{2} \boldsymbol{\alpha}' \hat{\mathbf{A}}' \hat{\mathbf{A}} \boldsymbol{\alpha} + \frac{1}{N} \mathbf{1}'_N \hat{\boldsymbol{\Sigma}} \boldsymbol{\alpha} + \tau \|\boldsymbol{\alpha}\|_1 - \frac{1}{2N} \right\} \quad \text{subject to } \mathbf{1}'_N \boldsymbol{\alpha} = 0, \quad (2.6)$$

where $\hat{\mathbf{A}} = (\mathbf{I}_N - N^{-1} \mathbf{1}_N \mathbf{1}'_N) \hat{\boldsymbol{\Sigma}}$ is the demeaned version of $\hat{\boldsymbol{\Sigma}}$. Denote $\hat{\boldsymbol{\alpha}} = \hat{\boldsymbol{\alpha}}_\tau$ as a solution to the dual problem in (2.6), and it is connected with the solution to the primal problem in (2.5) via

$$\hat{\mathbf{w}} = \hat{\mathbf{A}} \hat{\boldsymbol{\alpha}} + \frac{\mathbf{1}_N}{N}. \quad (2.7)$$

Remark 2.2. The dual problem of (2.4) can be produced by setting $\tau = 0$ in (2.6). When $\tau > 0$, (2.6) is a constrained ℓ_1 -penalized optimization where the criterion function is the summation of a quadratic form of $\boldsymbol{\alpha}$, a linear combination of $\boldsymbol{\alpha}$, and the ℓ_1 -norm of $\boldsymbol{\alpha}$, while the constraint is linear in $\boldsymbol{\alpha}$. The dual problem is instrumental in our theoretical analyses due to its similarity to the familiar Lasso, the ℓ_1 -penalized sparse regression problem (Tibshirani, 1996).

Remark 2.3. Since $\text{rank}(\widehat{\mathbf{A}}) \leq \text{rank}(\mathbf{I}_N - N^{-1}\mathbf{1}_N\mathbf{1}'_N) = N - 1$, the singularity of $\widehat{\mathbf{A}}$ may induce multiple solutions to the dual problem in (2.6). Despite this, the uniqueness of $\widehat{\mathbf{w}}$ as a solution to the primal problem in (2.5) implies $\widehat{\mathbf{A}}\widehat{\boldsymbol{\alpha}}^{(1)} = \widehat{\mathbf{A}}\widehat{\boldsymbol{\alpha}}^{(2)}$ for any $\widehat{\boldsymbol{\alpha}}^{(1)}$ and $\widehat{\boldsymbol{\alpha}}^{(2)}$ that solve (2.6). It is sufficient to find any solution to the dual problem to recover the same $\widehat{\mathbf{w}}$ in the primal.

In the next section, we assume that $\widehat{\boldsymbol{\Sigma}}$ has a certain structure and then examine its implications on the optimal forecast combination based on the ℓ_2 -relaxation.

3 Latent Group Structure and Its Implications

In this section, we impose latent group structures on $\widehat{\boldsymbol{\Sigma}}$ (or its population version) and then study its implications on the ℓ_2 -relaxed estimates of the weights.

3.1 Decomposition of $\widehat{\boldsymbol{\Sigma}}$ and Latent Group Structure

Statistical analysis of high dimensional problems typically postulates structures on the data generating process for dimension reduction. For example, variable selection methods such as Lasso (Tibshirani, 1996) and SCAD (Fan and Li, 2001) are motivated from regressions with sparsity, meaning most of the regression coefficients are either exactly zero or approximately zero. Similarly, in large VC estimation various structures have been considered in the literature. Bickel and Levina (2008) imposes many off-diagonal elements to be zero, Engle and Kelly (2012) assume a block equicorrelation structure, and Ledoit and Wolf (2004) use Bayesian methods for shrinking the sample correlation matrix to an equicorrelated target, to name a few.

Latent group structures in panel data is an alternative way to reduce dimensions, which now has grown into a burgeoning literature. While the optimization problem (2.5) is formulated for a generic covariance matrix $\widehat{\boldsymbol{\Sigma}}$, to analyze it in depth in the high dimensional framework we assume $\widehat{\boldsymbol{\Sigma}} = \{\widehat{\Sigma}_{ij}\}_{i,j \in [N]}$ can be approximated by a block equicorrelation matrix:

$$\widehat{\boldsymbol{\Sigma}} = \widehat{\boldsymbol{\Sigma}}^* + \widehat{\boldsymbol{\Sigma}}^e, \quad (3.1)$$

where $\widehat{\boldsymbol{\Sigma}}^* = \{\widehat{\Sigma}_{ij}^*\}_{i,j \in [N]}$ is a block equicorrelation matrix and $\widehat{\boldsymbol{\Sigma}}^e = \{\widehat{\Sigma}_{ij}^e\}_{i,j \in [N]}$ denotes the deviation of $\widehat{\boldsymbol{\Sigma}}$ from the block equicorrelation matrix. We will treat $\widehat{\boldsymbol{\Sigma}}^*$ as the oracle object in Section 3.2, instead of its population counterpart, so that we can study the finite-sample numerical properties in Section 3.3 without resorting to asymptotics. We write

$$\widehat{\boldsymbol{\Sigma}}^*_{(N \times N)} = \mathbf{Z}_{(N \times K)} \widehat{\boldsymbol{\Sigma}}^{\text{co}}_{(K \times K)} \mathbf{Z}' \quad (3.2)$$

where $\mathbf{Z} = \{Z_{ik}\}$ denotes an $N \times K$ binary matrix providing the cluster membership of each individual forecast, i.e., $Z_{ik} = 1$ if forecast i belongs to group $\mathcal{G}_k \subset [N]$ and $Z_{ik} = 0$ otherwise, and $\widehat{\boldsymbol{\Sigma}}^{\text{co}} = \{\widehat{\Sigma}_{kl}^{\text{co}}\}_{k,l \in [K]}$ is a $K \times K$ symmetric positive definite matrix. Here, the superscript ‘‘co’’ stands for ‘‘core’’. Note that

$$\widehat{\Sigma}_{ij}^* = \widehat{\Sigma}_{kl}^{\text{co}} \text{ if } i \in \mathcal{G}_k \text{ and } j \in \mathcal{G}_l.$$

An econometrician observes $\widehat{\Sigma}$ from the data but not $\widehat{\Sigma}^*$. We will be precise about the definition of “approximation” for $\widehat{\Sigma}^e$ in Assumption 1 later. Recall that $N_k := |\mathcal{G}_k|$ denotes the number of individuals in the k -th group, and $\sum_{k=1}^K N_k = N$.

For ease of notation and after necessary re-ordering the N forecast units, we write

$$\widehat{\Sigma}^* = (\widehat{\Sigma}_{kl}^{\text{co}} \cdot \mathbf{1}_{N_k} \mathbf{1}'_{N_l})_{k,l \in [K]}, \quad (3.3)$$

in which the units in the same group cluster together in a block. The re-ordering is for the convenience of notation only. The theory to be developed is irrelevant to the ordering of the forecast units, and does not require the knowledge on the membership matrix \mathbf{Z} .

We now motivate the above decomposition via two examples.

Example 2 *Chan and Pauwels (2018) assume the existence of a “best” unbiased forecast f_{0t} of variable y_{t+1} with an associated forecast error e_{0t} , and the forecast error e_{it} of model i can be decomposed as*

$$e_{it} = e_{0t} + u_{it},$$

where e_{0t} represents the forecast error from the best forecasting model, and u_{it} is the deviation of e_{it} from the best forecasting model. Assuming that $E[u_{it}] = 0$ and $E[e_{0t}u_{it}] = 0$ for each i , the VC of $\mathbf{e}_t = (e_{1t}, \dots, e_{Nt})'$ can be written as $\Sigma_0 = E[\mathbf{e}_t \mathbf{e}'_t] = E[e_{0t}^2] \mathbf{1}_N \mathbf{1}'_N + E[\mathbf{u}_t \mathbf{u}'_t]$, where $\mathbf{u}_t = (u_{1t}, \dots, u_{Nt})'$. At the sample level, we have $\widehat{\Sigma} = \widehat{\Sigma}^* + \widehat{\Sigma}^e$, where

$$\widehat{\Sigma} = \mathbb{E}_T[\mathbf{e}_t \mathbf{e}'_t], \quad \widehat{\Sigma}^* = \mathbb{E}_T[e_{0t}^2] \mathbf{1}_N \mathbf{1}'_N, \quad \text{and} \quad \widehat{\Sigma}^e = \mathbb{E}_T[\mathbf{u}_t \mathbf{u}'_t] + \mathbf{1}_N \mathbb{E}_T[e_{0t} \mathbf{u}'_t] + \mathbb{E}_T[e_{0t} \mathbf{u}_t] \mathbf{1}'_N.$$

In this case, all the N forecast units belong to the same group \mathcal{G}_1 as $\text{rank}(\widehat{\Sigma}^*) = 1$.

Example 3 *Consider that each individual forecast f_{it} is generated from a factor model*

$$f_{it} = \boldsymbol{\lambda}'_{g_i} \boldsymbol{\eta}_t + u_{it}, \quad (3.4)$$

where $\boldsymbol{\lambda}_{g_i}$ is a $q \times 1$ vector of factor loadings, $\boldsymbol{\eta}_t$ is a $q \times 1$ vector of latent factors, and u_{it} is an idiosyncratic shock. Here g_i denotes individual i 's membership, i.e., it takes value k if individual i belongs to group \mathcal{G}_k for $k \in [K]$ and $i \in [N]$. Similarly, assume $y_{t+1} = \boldsymbol{\lambda}'_y \boldsymbol{\eta}_t + u_{y,t+1}$. Assume that $E[u_{it}|\boldsymbol{\eta}_t] = 0$ and $E[u_{y,t+1}|\boldsymbol{\eta}_t] = 0$ and $E[u_{it}u_{y,t+1}|\boldsymbol{\eta}_t] = 0$.² For simplicity, we also assume conditional homoskedasticity $\text{var}(\mathbf{u}_t|\boldsymbol{\eta}_t) = \boldsymbol{\Omega}_u$ and the factor loadings are nonstochastic. Then individual i 's forecast error is

$$e_{it} := y_{t+1} - f_{it} = [(\boldsymbol{\lambda}_y - \boldsymbol{\lambda}_{g_i})' \boldsymbol{\eta}_t + u_{y,t+1}] - u_{it} = \boldsymbol{\lambda}'_{g_i} \boldsymbol{\eta}_t^\dagger - u_{it},$$

where $\boldsymbol{\eta}_t^\dagger = (\boldsymbol{\eta}_t', u_{y,t+1})'$ and $\boldsymbol{\lambda}_{g_i}^\dagger = ((\boldsymbol{\lambda}_y - \boldsymbol{\lambda}_{g_i})', 1)'$, or equivalently $\mathbf{e}_t = \boldsymbol{\Lambda}^\dagger \boldsymbol{\eta}_t^\dagger - \mathbf{u}_t$ in a vector form, where $\boldsymbol{\Lambda}^\dagger := (\boldsymbol{\lambda}_{g_1}^\dagger, \dots, \boldsymbol{\lambda}_{g_N}^\dagger)'$. The population VC of \mathbf{e}_t is given by

$$\Sigma_0 = E[\mathbf{e}_t \mathbf{e}'_t] = \boldsymbol{\Lambda}^\dagger E[\boldsymbol{\eta}_t^\dagger \boldsymbol{\eta}_t^{\dagger'}] \boldsymbol{\Lambda}^{\dagger'} + \boldsymbol{\Omega}_u.$$

²Other than those q factors in $\boldsymbol{\eta}_t$, the additional latent factor $u_{y,t+1}$ in y_{t+1} is unforeseeable at time t . In other words, given the information set \mathcal{I}_t that contains $(\{f_{it}\}_{i \in [N]}, \boldsymbol{\eta}_t)$ and \mathbf{u}_t at time t , the error $u_{y,t+1} = y_{t+1} - \boldsymbol{\lambda}'_y \boldsymbol{\eta}_t = y_{t+1} - E(y_{t+1}|\mathcal{I}_t)$ must be orthogonal to \mathcal{I}_t . Then $E[u_{it}u_{y,t+1}|\mathcal{I}_t] = 0$ implies $E[u_{it}u_{y,t+1}|\boldsymbol{\eta}_t] = 0$ by the law of iterated expectations.

Decompose the sample VC as $\widehat{\Sigma} = \widehat{\Sigma}^* + \widehat{\Sigma}^e$, where

$$\widehat{\Sigma} = \mathbb{E}_T [\mathbf{e}_t \mathbf{e}_t'], \quad \widehat{\Sigma}^* = \mathbf{\Lambda}^\dagger \mathbb{E}_T [\boldsymbol{\eta}_t^\dagger \boldsymbol{\eta}_t^{\dagger'}] \mathbf{\Lambda}^{\dagger'}, \quad \text{and} \quad \widehat{\Sigma}^e = \mathbb{E}_T [\mathbf{u}_t \mathbf{u}_t'] - \mathbf{\Lambda}^\dagger \mathbb{E}_T [\boldsymbol{\eta}_t^\dagger \mathbf{u}_t'] - \mathbb{E}_T [\mathbf{u}_t \boldsymbol{\eta}_t^{\dagger'}] \mathbf{\Lambda}^{\dagger'}.$$

By construction, the core matrix has element $\widehat{\Sigma}_{kl}^{\text{co}} = \lambda_k^\dagger \mathbb{E}_T [\boldsymbol{\eta}_t^\dagger \boldsymbol{\eta}_t^{\dagger'}] \lambda_l^\dagger$ for $k, l \in [K]$, the equicorrelation matrix has element $\widehat{\Sigma}_{ij}^* = \widehat{\Sigma}_{kl}^{\text{co}}$ if $i \in \mathcal{G}_k$ and $j \in \mathcal{G}_l$, and $\text{rank}(\widehat{\Sigma}^*) \leq (q+1) \wedge K$.

Remark 3.1. We emphasize that our theory below does not require the knowledge on the group membership for individual forecasts. If one believes in the multi-factor structure in (3.4), he can always conduct the principal component analysis (PCA) to estimate the factor loadings and factors first in the large N and large T framework. Then he can apply either the K-means algorithm or the sequential binary segmentation algorithm (Wang and Su, 2020) to the estimated factor loadings to identify the group membership. Alternatively, one can consider the regression of f_{it} on the estimated factors and apply the classifier-Lasso (Su et al., 2016) or other methods to recover the group membership. The advantage of ℓ_2 -relaxation is that it directly works with the sample moments and hence bypasses the factor structure and the group membership.

Remark 3.2. It is worth mentioning that Hsiao and Wan (2014) also assume that the forecast errors exhibit a multi-factor structure. But they do not assume the presence of K latent groups in the N factor loadings and write $\boldsymbol{\lambda}_i$ in place of $\boldsymbol{\lambda}_{g_i}$ in our Example 3. In the absence of the latent group structures among the factor loadings $\{\boldsymbol{\lambda}_i\}_{i \in [N]}$ in the true DGP, we can follow the lead of Bonhomme et al. (2017) and consider their discretization. For simplicity, if $\boldsymbol{\lambda}_i$ lies in a compact parameter space \mathbb{R}^q , then for each $i \in [N]$ there exists K , $g_i \in [K]$ and $\boldsymbol{\lambda}_{g_i} \in \mathbb{R}^q$ such that $\|\boldsymbol{\lambda}_i - \boldsymbol{\lambda}_{g_i}\|_2 \leq \delta_K$. As K diverges to infinity, the approximation error δ_K can be made as small as possible. As a result, we can continue to decompose the i -th forecast error in Example 3 as $e_{i,t} = \boldsymbol{\lambda}_{g_i}^{\dagger'} \boldsymbol{\eta}_t^\dagger - u_{i,t}^a$, where $\boldsymbol{\lambda}_{g_i}^\dagger = ((\boldsymbol{\lambda}_y - \boldsymbol{\lambda}_{g_i})', 1)'$ and $u_{i,t}^a$ now contains the discretization error.

Latent groups are present not only in approximate factor models, as in the above two motivating examples, but also in many forecast problems in which multi-factor structures are hidden implicitly. Here follows such an example.

Example 4 Suppose that the outcome variable y_{t+1} is generated via the process

$$y_{t+1} = \mathbf{x}_t' \boldsymbol{\theta}^0 + u_{t+1} \quad \text{for } t = -T_0, \dots, -1, 0, 1, \dots$$

where $\mathbf{x}_t = (x_{j,t})_{j=1}^p$ is a $p \times 1$ vector of potential predictive variables, $\boldsymbol{\theta}^0 = (\theta_j^0)_{j=1}^p$ is a $p \times 1$ vector of regression coefficients. Due to costly variable collection or ignorance, the forecaster i utilizes only a subset $\mathbf{x}_{S_i,t}$ of \mathbf{x}_t , where $S_i \subset [p]$, to exercise prediction with the OLS estimate. Let $\widehat{\boldsymbol{\theta}}_{S_i,t} = (\sum_{l=-T_0+1}^t \mathbf{x}_{S_i,l-1} \mathbf{x}_{S_i,l-1}')^{-1} \sum_{l=-T_0+1}^t \mathbf{x}_{S_i,l-1} y_l$, and $\widehat{\boldsymbol{\theta}}_{i,t}$ be the sparse $p \times 1$ vector that embeds the corresponding $\widehat{\boldsymbol{\theta}}_{S_i,t}$ so that $(\widehat{\boldsymbol{\theta}}_{i,t})_{S_i} = \widehat{\boldsymbol{\theta}}_{S_i,t}$ and $(\widehat{\boldsymbol{\theta}}_{i,t})_{[p] \setminus S_i} = 0$. We consider two forecasting schemes: fixed window and rolling window.

(i) In the case of fixed estimation window for $t \in [-T_0, \dots, 0]$, the i -th forecast of y_{t+1} is given by $\widehat{f}_{i,t} := \mathbf{x}_{S_i,t}' \widehat{\boldsymbol{\theta}}_{S_i,0}$ for $t \geq 1$. The associated forecast error

$$e_{it} = y_{t+1} - \widehat{f}_{i,t} = y_{t+1} - \mathbf{x}_t' \widehat{\boldsymbol{\theta}}_{i,0} = u_{t+1} + \mathbf{x}_t' (\boldsymbol{\theta}^0 - \widehat{\boldsymbol{\theta}}_{i,0}).$$

Apparently, this is an exact $(p+1)$ -factor model with factors $(\mathbf{x}_t', u_{t+1})'$ and factor loadings $((\boldsymbol{\theta}^0 - \widehat{\boldsymbol{\theta}}_{i,0})', 1)'$.

(ii) In the case of rolling window of extending length for $t \in [-T_0, \dots, t]$, the forecast error is

$$e_{it} = y_{t+1} - \mathbf{x}'_{S_i,t} \widehat{\boldsymbol{\theta}}_{S_i,t} = u_{t+1} + \mathbf{x}'_t(\boldsymbol{\theta}^0 - \widehat{\boldsymbol{\theta}}_{i,t}) = u_{t+1} + \mathbf{x}'_t(\boldsymbol{\theta}^0 - \boldsymbol{\theta}_i^0) + \epsilon_{i,t}$$

where $\widehat{\boldsymbol{\theta}}_{i,t} \xrightarrow{p} \boldsymbol{\theta}_i^0$ as $T_0 \rightarrow \infty$ is assumed to hold uniformly in (i, t) under some regularity conditions that include the covariance stationarity, and $\epsilon_{i,t} := \mathbf{x}'_t(\widehat{\boldsymbol{\theta}}_{i,t} - \boldsymbol{\theta}_i^0)$. Therefore, we have an approximate $(p+1)$ -factor model with factors $(\mathbf{x}'_t, u_{t+1})'$ and factor loadings $((\boldsymbol{\theta}^0 - \boldsymbol{\theta}_i^0)', 1)'$. Similar analysis applies to the rolling window of fixed length L , in which the forecaster i estimates the coefficient by $\widehat{\boldsymbol{\theta}}_{S_i,t}^L = (\sum_{l=t-L+1}^t \mathbf{x}_{S_i,t-1} \mathbf{x}'_{S_i,t-1})^{-1} \sum_{l=t-L+1}^t \mathbf{x}_{S_i,t-1} \mathbf{y}_t$.

In either case, if $p = 10$ then the total number of potential combinations amounts to $N = 2^{10} = 1024$ forecasting models under consideration. It is reasonable to model the forecast errors with a multi-factor structure as above.

3.2 The Oracle Problem

We consider the oracle problem by assuming away the idiosyncratic component $\widehat{\boldsymbol{\Sigma}}^e$ in $\widehat{\boldsymbol{\Sigma}}$. Given the latent group structure, an oracle version of (1.2) is

$$\min_{\mathbf{w} \in \mathbb{R}^N} \frac{1}{2} \mathbf{w}' \widehat{\boldsymbol{\Sigma}}^* \mathbf{w} \quad \text{subject to} \quad \mathbf{w}' \mathbf{1}_N = 1, \quad (3.5)$$

where the infeasible block equicorrelation covariance matrix $\widehat{\boldsymbol{\Sigma}}^*$ replaces the sample covariance $\widehat{\boldsymbol{\Sigma}}$ in (1.2). In the oracle problem, the group structure can be identified by inspecting the values of the entries in $\widehat{\boldsymbol{\Sigma}}^*$. The rank of $\widehat{\boldsymbol{\Sigma}}^*$ is at most K due to the latent group pattern, leading to multiple and in fact an infinite number of solutions. Despite this, each solution yields the same optimal value for the primal objective function. Therefore it suffices to identify only one solution. The oracle counterpart of the ℓ_2 -relaxation primal problem in (2.5) is

$$\min_{(\mathbf{w}, \gamma) \in \mathbb{R}^{N+1}} \frac{1}{2} \|\mathbf{w}\|_2^2 \quad \text{subject to} \quad \mathbf{w}' \mathbf{1}_N = 1 \text{ and } \|\widehat{\boldsymbol{\Sigma}}^* \mathbf{w} + \gamma \mathbf{1}_N\|_\infty \leq \tau. \quad (3.6)$$

Denote the solution of the weights in the above problem as \mathbf{w}_τ^* , which in general does not accommodate an explicit form due to the presence of sup-norm in the inequality constraint. However, in the special case of $\tau = 0$, the inequality constraint can be equivalently written as K equality constraints so that we can solve for $\mathbf{w}_0^* := \mathbf{w}_{\tau=0}^*$ in closed-form. Define $r_k := N_k/N$ as the fraction of the k -th group members on the cross section. Let $\mathbf{r} = (r_k)_{k \in [K]}$ and $\mathbf{r}^{1/2} = (\sqrt{r_k})_{k \in [K]}$. Let “ \circ ” be the Hadamard product.

Lemma 2 (a) *The solution to (3.6) must take within-group equal values in the form*

$$\mathbf{w}_\tau^* = (N^{-1} b_{\tau 1}^* \mathbf{1}'_{N_1}, \dots, N^{-1} b_{\tau K}^* \mathbf{1}'_{N_K})', \quad (3.7)$$

where $\mathbf{b}_\tau^* := (b_{\tau k}^*)_{k \in [K]}$ solves

$$\min_{(\mathbf{b}, \gamma) \in \mathbb{R}^{K+1}} \frac{1}{2N} \left\| \mathbf{b} \circ \mathbf{r}^{1/2} \right\|_2^2 \quad \text{subject to} \quad \mathbf{1}'_K (\mathbf{b} \circ \mathbf{r}) = 1 \text{ and } \|\widehat{\boldsymbol{\Sigma}}^{\text{co}} (\mathbf{b} \circ \mathbf{r}) + \gamma \mathbf{1}_K\|_\infty \leq \tau. \quad (3.8)$$

(b) In the special case of $\tau = 0$, the solution $\mathbf{b}_0^* := \mathbf{b}_{\tau=0}^*$ to (3.8) is given by

$$\mathbf{b}_0^* = \mathbf{r}^{-1} \circ \frac{(\widehat{\boldsymbol{\Sigma}}^{\text{co}})^{-1} \mathbf{1}_K}{\mathbf{1}'_K (\widehat{\boldsymbol{\Sigma}}^{\text{co}})^{-1} \mathbf{1}_K}, \quad (3.9)$$

where $\mathbf{r}^{-1} = (r_k^{-1})_{k \in [K]}$.

Lemma 2(a) shows that the squared ℓ_2 -norm objective function produces the *within-group equally weighted solution* \mathbf{w}_τ^* . The high dimensional oracle problem with $(N + 1)$ free parameters in (3.6) is reduced to only $(K + 1)$ free parameters by the ‘‘core primal problem’’ in (3.8). When $\tau = 0$, we can obtain the explicit solution of \mathbf{w}_0^* by inserting (3.9) into (3.7).

As in Lemma 1, we proceed with the dual problem. The dual of (3.6) is

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^N} \left\{ \frac{1}{2} \boldsymbol{\alpha}' \widehat{\mathbf{A}}^* \widehat{\mathbf{A}}^* \boldsymbol{\alpha} + \frac{1}{N} \mathbf{1}'_N \widehat{\boldsymbol{\Sigma}}^* \boldsymbol{\alpha} + \tau \|\boldsymbol{\alpha}\|_1 - \frac{1}{2N} \right\} \quad \text{subject to } \mathbf{1}'_N \boldsymbol{\alpha} = 0, \quad (3.10)$$

where $\widehat{\mathbf{A}}^* = (\mathbf{I}_N - N^{-1} \mathbf{1}_N \mathbf{1}'_N) \widehat{\boldsymbol{\Sigma}}^*$. For any $\boldsymbol{\alpha}^* = \boldsymbol{\alpha}_\tau^*$ that solves (3.10), Lemma 1 implies that the solution to (3.10) is not unique, and the unique \mathbf{w}_τ^* and the non-unique $\boldsymbol{\alpha}_\tau^*$ are connected via

$$\mathbf{w}_\tau^* = \widehat{\mathbf{A}}^* \boldsymbol{\alpha}_\tau^* + \frac{\mathbf{1}_N}{N}. \quad (3.11)$$

To develop the counterpart of Lemma 2 for the dual problem, we need some extra notations. For a generic N -vector $\boldsymbol{\alpha} = (\alpha_i)_{i \in [N]}$, denote $a_k = \sum_{i \in \mathcal{G}_k} \alpha_i$ as the k -th within-group summation of $\boldsymbol{\alpha}$. Let $\mathbf{a} = (a_k)_{k \in [K]}$. Let $\widehat{\mathbf{A}}^{\text{co}} = \mathbf{R}^{1/2} (\mathbf{I}_K - \mathbf{1}_K \mathbf{r}') \widehat{\boldsymbol{\Sigma}}^{\text{co}}$, where $\mathbf{R} = \text{diag}(\mathbf{r})$ is the $K \times K$ diagonal matrix that stacks the elements of \mathbf{r} along the diagonal line. Note that $\widehat{\mathbf{A}}^{\text{co}}$ is the *weighted demeaned core* of $\widehat{\mathbf{A}}^*$ with the weights depending on the relative group size \mathbf{r} . The following lemma characterizes the features of $\boldsymbol{\alpha}^* = \boldsymbol{\alpha}_\tau^*$.

Lemma 3 (a) If $\tau > 0$, any solution $\boldsymbol{\alpha}^* = \boldsymbol{\alpha}_\tau^*$ to (3.10) must be of similar signs, i.e., $\boldsymbol{\alpha}^* \in \mathbb{S}^{\text{all}}$.

(b) The low dimensional core dual problem for (3.10) is

$$\min_{\mathbf{a} \in \mathbb{R}^K} \left\{ \frac{N}{2} \mathbf{a}' \widehat{\mathbf{A}}^{\text{co}'} \widehat{\mathbf{A}}^{\text{co}} \mathbf{a} + \mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{a} + \tau \|\mathbf{a}\|_1 - \frac{1}{2N} \right\} \quad \text{subject to } \mathbf{1}'_K \mathbf{a} = 0. \quad (3.12)$$

(c) In the special case of $\tau = 0$, the solution $\mathbf{a}_0^* := \mathbf{a}_{\tau=0}^*$ to (3.12) is

$$\mathbf{a}_0^* = N^{-1} (\tilde{\mathbf{A}}^{\text{co}'} \tilde{\mathbf{A}}^{\text{co}})^{-1} \tilde{\mathbf{A}}^{\text{co}'} ((\mathbf{b}_0^* \circ \mathbf{r} - \mathbf{r})' \mathbf{R}^{-1/2} \quad 0) ', \quad (3.13)$$

whereas $\tilde{\mathbf{A}}^{\text{co}} := (\widehat{\mathbf{A}}^{\text{co}'} \quad \mathbf{1}_K)'$ is a $(K + 1) \times K$ matrix of full column rank.

Remark 3.4. Lemma 3(a) shows that the ℓ_1 -norm penalization in (3.10) precludes opposite signs of the estimates $\boldsymbol{\alpha}_\tau^*$ within a group, which implies $\|\boldsymbol{\alpha}_\tau^*\|_1 = \|\mathbf{a}_\tau^*\|_1$ for any $\tau > 0$. Lemma 3(b) reduces the high dimensional oracle dual problem in $\boldsymbol{\alpha} \in \mathbb{R}^N$ to the low dimensional oracle dual one in $\mathbf{a} \in \mathbb{R}^K$. Lemma 3(c) is the counterpart of (3.9) for the dual, which involves the *augmented (by a row of 1's) weighted demeaned core* $\tilde{\mathbf{A}}^{\text{co}}$. A numerical lower bound and a stochastic lower bound of $\tilde{\mathbf{A}}^{\text{co}}$ will be established in Lemma 4(b) and Lemma 5(a) in Section A.3 of the Appendix.

3.3 Numerical Properties

In this section, we derive the (finite sample) numerical properties of the estimates $\widehat{\mathbf{w}}$ and $\widehat{\boldsymbol{\alpha}}$ under a finite N . These properties are prepared for the development of the asymptotic theory in Section 4. The sample estimator $\widehat{\mathbf{w}}$ deviates from the oracle \mathbf{w}_τ^* due to the presence of the idiosyncratic shock and the tuning parameter τ . We will show that the effect of the idiosyncratic shock is embodied by the quantity

$$\phi_e := \|\widehat{\boldsymbol{\Sigma}}^e\|_{e2},$$

which can be viewed as a measurement of the noise level or contamination level of $\widehat{\boldsymbol{\Sigma}}^*$ in the model. Theorem 1 below reports the numerical properties of the sample estimator, where the condition $\tau > \phi_e \|\mathbf{b}_0^*\|_\infty / \sqrt{N}$ can be satisfied w.p.a.1 in the asymptotic analysis.

Theorem 1 *Suppose that $\tau > \phi_e \|\mathbf{b}_0^*\|_\infty / \sqrt{N}$. Then*

- (a) $\|\widehat{\mathbf{w}}\|_2 \leq \|\mathbf{w}_0^*\|_2 \leq \|\mathbf{b}_0^*\|_\infty / \sqrt{N}$;
- (b) $\widehat{\boldsymbol{\alpha}} \in \mathbb{S}^{\text{all}}$.

Remark 3.5. Theorem 1(a) reports the upper bound for $\|\widehat{\mathbf{w}}\|_2$, which is used in establishing part (b). If the ratio between the tolerance τ and the noise level ϕ_e is sufficiently large in that it is larger than $\|\mathbf{b}_0^*\|_\infty / \sqrt{N}$, the estimator $\widehat{\boldsymbol{\alpha}}$ must be of similar sign within each group. This result is proved by exploiting the KKT conditions associated with the Lagrangian of (2.6). The intuition is that when the specified τ is large, for any $i, j \in \mathcal{G}_k$, the column-wise difference in the noise, i.e., $\widehat{\boldsymbol{\Sigma}}_{\cdot i}^e - \widehat{\boldsymbol{\Sigma}}_{\cdot j}^e$, is unable to push the two associated KKT conditions to be satisfied simultaneously for the pair of $\widehat{\alpha}_i$ and $\widehat{\alpha}_j$ of opposite signs.

Remark 3.6. The result in Theorem 1(b) is important. It reminds us of the grouping effect of elastic net (Zou and Hastie, 2005). A regression method exhibits the grouping effect if the regression coefficients of a group of highly correlated regressors in the design matrix \mathbf{X} tend to be equal (up to a change of sign if negatively correlated). It is well-known that while Lasso yields sparse solutions in many cases, it does not have the grouping effect. In contrast, the elastic net penalty, as a convex combination of the Lasso (ℓ_1) and ridge (ℓ_2) penalties, encourages the grouping effect and has the advantage of including automatically all the highly correlated variables in the group. In the presence of a latent group structure in the dominant component $\widehat{\boldsymbol{\Sigma}}^*$ of $\widehat{\boldsymbol{\Sigma}}$, $\widehat{\boldsymbol{\Sigma}}$ plays the role as the *Gram matrix* $\mathbf{X}'\mathbf{X}$. Then $\widehat{\boldsymbol{\Sigma}}_{\cdot i}$ and $\widehat{\boldsymbol{\Sigma}}_{\cdot j}$ are asymptotically collinear if the i -th and j -th forecasts come from the same group (e.g., $i, j \in \mathcal{G}_k$ for some $k \in [K]$), a feature similar to highly correlated regressors in the regression framework. As a result, the ℓ_2 -relaxation estimator of the weights enjoys similar properties as the elastic net estimator.

Remark 3.7. The ℓ_1 penalized form in (2.6) also draws similarity to high dimensional Lasso estimation under sparsity. The consistency of Lasso requires that the correlation among the columns of the design matrix should not be too strong; otherwise, various versions of restricted eigenvalue conditions break down (Bickel et al., 2009; Van De Geer and Bühlmann, 2009; Belloni et al., 2012). When $\widehat{\boldsymbol{\Sigma}}$ is treated as a design matrix (or more precisely, the Gram matrix) in the regression framework, its columns are highly correlated, or asymptotically perfectly collinear. Consider the extreme case where $\widehat{\boldsymbol{\Sigma}} = \widehat{\boldsymbol{\Sigma}}^*$ so that $\widehat{\boldsymbol{\Sigma}}^*$ is not contaminated by the noise component $\widehat{\boldsymbol{\Sigma}}^e$. If we try to solve (2.6) in this case, the estimated $\widehat{\boldsymbol{\alpha}}$ will be numerically very unstable due to perfect collinearity, and we cannot expect it to converge to a fixed $\boldsymbol{\alpha}_0^*$ under the ℓ_1 norm. Therefore we must seek a compatibility condition tailored for the group structure. Due to its technical nature, we relegate it to Lemma 4 in Section A.3 of the Appendix.

4 Asymptotic Theory

We study the asymptotic properties of our ℓ_2 -relaxed estimator of the combination weights in this section. To this end, we impose some conditions and study the asymptotic properties of the estimator in the dual problem first.

4.1 Assumptions

We consider a triangular array of models indexed by T and N , both of which pass to infinity. Let $\phi_{NT} := \sqrt{(\log N)/(T \wedge N)} \rightarrow 0$. Note that we allow both $N \gg T$ (as in standard high dimensional problems) and $T \gg N$ or $T \asymp N$ in the definition of ϕ_{NT} . But we rule out the traditional case of “fixed N and large T ”, which can be trivially handled by (1.2). Furthermore, let $\Sigma_0^e = E[\widehat{\Sigma}^e]$, $\Delta^e = \widehat{\Sigma}^e - \Sigma_0^e$, $\Sigma_0^{\text{co}} = E[\widehat{\Sigma}^{\text{co}}]$, and $\Delta^{\text{co}} = \widehat{\Sigma}^{\text{co}} - \Sigma_0^{\text{co}}$. We impose the following assumption.

Assumption 1

- (a) $\phi_{\max}(\Sigma_0^e) = O(\sqrt{N}\phi_{NT})$, $\|\Sigma_0^e\|_{c2} \leq C_{e0} \cdot \phi_{\max}(\Sigma_0^e)$, and $\|\Delta^e\|_{\infty} = O_p((T/\log N)^{-1/2})$;
(b) $\underline{c} \leq \phi_{\min}(\Sigma_0^{\text{co}}) \leq \phi_{\max}(\Sigma_0^{\text{co}}) \leq \bar{c}$, and $\|\Delta^{\text{co}}\|_{\infty} = O_p((T/\log N)^{-1/2})$;

where \underline{c} , \bar{c} , and C_{e0} are positive finite constants.

The first condition in Assumption 1(a) allows the maximal eigenvalue of the $N \times N$ matrix Σ_0^e to diverge to infinity, but at a limited rate $\sqrt{N}\phi_{NT}$. The second condition in (a) is similar to but weaker than the absolute row-sum condition that is frequently used to model weak cross-sectional dependence; see, e.g., Fan et al. (2013). The third condition in (a) requires that the sampling error of Δ_{ij}^e be controlled by $(T/\log N)^{-1/2}$ uniformly over i and j so that each element of $\widehat{\Sigma}^e$ should not deviate too much from its population mean Σ_0^e . This condition can be established under some low-level assumptions; see, e.g., Chapter 6 in Wainwright (2019). Assumption 1(b) bounds all eigenvalues of the population core from 0 and infinity, and impose similar stochastic order on Δ^{co} as that on Δ^e in Assumption 1(a).

Example 5 (Example 3, cont.) Following the notation of Example 3, we can decompose the population variance-covariance matrix $\Sigma := E[\mathbf{e}_t \mathbf{e}_t']$ as $\Sigma = \Sigma^* + \Sigma_0^e$, where $\Sigma^* = \Lambda^\dagger E[\boldsymbol{\eta}_t^\dagger \boldsymbol{\eta}_t^\dagger'] \Lambda^\dagger$, and $\Sigma_0^e = \Omega_x = \{\Omega_{x,ij}\}$. The corresponding sampling error is

$$\Delta_{ij}^e = \{\mathbb{E}_T[\epsilon_{i,t} \epsilon_{j,t}] - \Omega_{x,ij}\} - \sum_{l \in \{i,j\}} \{(\lambda_y - \lambda_{g_l})' \mathbb{E}_T[\boldsymbol{\eta}_t(u_{y,t+1} - u_{l,t})] + \mathbb{E}_T[u_{y,t+1} u_{l,t}]\}.$$

Then the first part of Assumption 1(a) is satisfied as long as $\phi_{\max}(\Omega_x) = O(\sqrt{N}\phi_{NT})$. For the sampling error matrix, if

$$\max_{i,j \in [N]} \{|\mathbb{E}_T[\boldsymbol{\eta}_t(u_{y,t+1} - u_{i,t})]| + |\mathbb{E}_T[u_{y,t} u_{i,t}]| + |\mathbb{E}_T[u_{i,t} u_{i,j}] - \Omega_{ij,x}|\} = O_p((T/\log N)^{-1/2}),$$

then $\|\Delta^e\|_{\infty} = O_p((T/\log N)^{-1/2})$ is satisfied as well.

The extent of relaxation in (2.5) is controlled by the tuning parameter τ , which is to be chosen by cross validations (CV) in simulations and applications. We spell out admissible range of τ in Assumption 2(a) below. Assumption 2(b) restricts $\underline{r} := \min_{k \in [K]} r_k$ relative to K .

Assumption 2 As $(N, T) \rightarrow \infty$,

(a) $\sqrt{K}\phi_{NT}/\tau + K^{5/2}\tau \rightarrow 0$;

(b) $\underline{\tau} \asymp K^{-1}$.

In order to meet the condition $\sqrt{K}\phi_{NT}/\tau \rightarrow 0$ in Assumption 2(a), it is suffice to specify

$$\tau = D_\tau \sqrt{K} \phi_{NT}$$

for some slowly diverging sequence D_τ as $(N, T) \rightarrow \infty$, for example $\log \log (N \wedge T)$. If K is finite, this specification implies that τ should shrink to zero at a rate slightly slower than ϕ_{NT} . We allow $K \rightarrow \infty$, provided $K^{5/2}\tau \rightarrow 0$ so that the sampling error in $\widehat{\Sigma}^e$ would not offset the dominant grouping effect of the ℓ_2 -relaxation in the presence of latent groups in $\widehat{\Sigma}^*$. The particular rate $K^{5/2}\tau$ will appear as the order of convergence in Theorem 3 below. Though the exact number of groups K is usually unknown in reality, if the researcher believes that K is asymptotically dominated by some explicit rate function $\bar{K}_{N,T}$ of N and T in that $\limsup K/\bar{K}_{N,T} < 1$, say $\bar{K}_{N,T} = (N \wedge T)^{1/7}$, then all the following theoretical results still hold if K is replaced by $\bar{K}_{N,T}$ and τ is replaced by $\tau_{\bar{K}} = C_\tau \bar{K}_{N,T}^{1/2} \phi_{NT}$ for some positive constant C_τ , as long as Assumption 2(a) is replaced by $\bar{K}_{N,T}^{1/2} \phi_{NT}/\tau + \bar{K}_{N,T}^{5/2} \tau \rightarrow 0$ accordingly.

Assumption 2(b) requires that the smallest relative group size $\underline{\tau}$ be proportional to the reciprocal of K . Had a group included too few members, the weight of the group would be too small to matter and the corresponding coefficient b_{0k}^* in (3.9) is inversely affected by the group size. Indeed Assumption 2(b) is a simplifying assumption for notational conciseness. If we drop it, $\underline{\tau}$ will appear in the rates of convergence in all the following results, which complicates the expressions but adds no new insight.

4.2 Asymptotic Properties of $\widehat{\alpha}$ in the Dual

We start with the dual problem (2.6) under Assumption 1. Following the discussion in Section 3.2, there are multiple solutions to the oracle dual in (3.10) due to the rank deficiency of Σ^* . But if we want to establish convergence in probability, we must declare a target to which the estimator will converge. We construct such a desirable α_0^* in (4.1) below, denoted as $\widehat{\alpha}^*$ where the ‘‘hat’’ signifies its dependence on the realization of $\widehat{\alpha}$ and ‘‘star’’ indicates its validity as an oracle estimator. Define $\widehat{\alpha}^* = (\widehat{\alpha}_{\mathcal{G}_k}^*)_{k \in [K]}$, where

$$\widehat{\alpha}_{\mathcal{G}_k}^* = a_{0k}^* \left(\frac{\widehat{\alpha}_{\mathcal{G}_k}}{\widehat{a}_k} \cdot 1 \{ \widehat{a}_k a_{0k}^* > 0 \} + \frac{\mathbf{1}_{N_k}}{N_k} \cdot 1 \{ \widehat{a}_k a_{0k}^* \leq 0 \} \right), \quad (4.1)$$

where $\widehat{a}_k = \sum_{i \in \mathcal{G}_k} \widehat{\alpha}_i$ is the sum of the $\widehat{\alpha}_i$ in the k -th group and $1 \{ \cdot \}$ is the usual indicator function. The above $\widehat{\alpha}_{\mathcal{G}_k}^*$ is designed such that the k -th oracle group weight a_{0k}^* is distributed across the k th group members proportionally to $\widehat{\alpha}_{\mathcal{G}_k}/\widehat{a}_k$ when \widehat{a}_k and a_{0k}^* share the same sign, whereas a_{0k}^* is distributed equally across k -th group members when they take opposite signs. When $\widehat{\alpha} \in \widetilde{\mathbb{S}}^{\text{all}}$, which holds w.p.a.1 in view of Theorem 1 and Lemma 5(b) in the Appendix, it is easy to verify that

$$(i) \widehat{\alpha}^* \in \widetilde{\mathbb{S}}^{\text{all}}, \quad (ii) \|\alpha_0^*\|_1 = \|\widehat{\alpha}^*\|_1 \quad \text{and} \quad (iii) \widehat{\alpha} - \widehat{\alpha}^* \in \widetilde{\mathbb{S}}^{\text{all}}. \quad (4.2)$$

For example, (i) in (4.2) holds because by construction $\widehat{\alpha}^* \in \mathbb{S}^{\text{all}}$ as long as $\widehat{\alpha} \in \widetilde{\mathbb{S}}^{\text{all}}$, and $\mathbf{1}'_N \widehat{\alpha}^* = \sum_{k=1}^K \mathbf{1}'_{N_k} \widehat{\alpha}_{\mathcal{G}_k}^* = \sum_{k=1}^K a_{0k}^* = \mathbf{1}'_K \mathbf{a}_0^* = 0$. The following theorem shows that the solution to the Lasso-type ℓ_1 -penalization problem (2.6) is close to the desirable oracle estimator $\widehat{\alpha}^*$.

Theorem 2 *Suppose that Assumptions 1 and 2 hold. Then*

$$\|\widehat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}^*\|_1 = O_p(N^{-1}K^3\tau) \quad \text{and} \quad \|\widehat{\mathbf{A}}(\widehat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}^*)\|_2 = O_p(N^{-1/2}K^2\tau).$$

Remark 4.1. Theorem 2 is a key result that characterizes the convergence rate of the high-dimensional parameter $\widehat{\boldsymbol{\alpha}}$ in the dual problem to its oracle group counterpart $\boldsymbol{\alpha}_0^*$, represented by the constructed unique solution $\boldsymbol{\alpha}^*$. Although our ultimate interest lies in the weight estimate $\widehat{\mathbf{w}}$ in the primal problem, in theoretical analysis we first work with $\widehat{\boldsymbol{\alpha}}$ in the dual problem instead. This detour is taken because the dual is an ℓ_1 -penalized optimization which resembles Lasso. The intensive study of Lasso in statistics and econometrics offers a set of inequalities involving the ℓ_1 -norms of $\widehat{\boldsymbol{\alpha}}$, $\boldsymbol{\alpha}^*$ and their difference ($\widehat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}^*$) at our disposal to analyze its asymptotic behavior.

Remark 4.2. All high dimensional estimation problems require certain notion of sparsity to reduce dimensionality. It is helpful to compare our setting of latent group structures with the sparse regression estimated by Lasso. For Lasso estimation, the complexity of the problem is governed by the total number of regressors (p in Example 1) while under sparsity those non-zero coefficients control the effective number of parameters, which is assumed to be far fewer than the sample size. For the ℓ_1 penalized dual in (2.6), the complexity is the number of forecasters N whereas under the group structures the number of groups K determines the effective number of parameters.

Remark 4.3. A critical technical step in proving the consistency of high-dimensional Lasso problems is the *compatibility condition* (Bühlmann and van de Geer, 2011, Ch 6.13) in conjunction with the *restricted eigenvalue condition* (Bickel et al., 2009). In our paper, the compatibility condition is established in Lemma 4(a) and the restricted eigenvalue is represented by $\phi_A := 1 \wedge \phi_{\min}(\widetilde{\mathbf{A}}^{\text{co}'}\widetilde{\mathbf{A}}^{\text{co}})$. In particular, instead of *assuming* a lower bound for the restricted eigenvalues as most high dimensional Lasso papers do, we *derive* a finite sample lower bound for ϕ_A in Lemma 4(b) as well as its convergence rate in Lemma 5(a) under our latent group structures and primitive Assumptions 1 and 2. We deem these developments as original contributions to the literature, though we relegate them to Section A.3 due to the technical nature.

4.3 Asymptotic Properties of $\widehat{\mathbf{w}}$ in the Primal

Given Theorem 2, we proceed to handle the estimator $\widehat{\mathbf{w}}$ in the ℓ_2 -relaxation primal problem. Notice that for the simple average weight $\widehat{\mathbf{w}}^{\text{SA}} = \mathbf{1}_N/N$, although $\widehat{\mathbf{w}}^{\text{SA}}$ does not mimic \mathbf{w}_0^* in general, the ℓ_2 -distance

$$\|\widehat{\mathbf{w}}^{\text{SA}} - \mathbf{w}_0^*\|_2 \leq \|\widehat{\mathbf{w}}^{\text{SA}}\|_2 + \|\mathbf{w}_0^*\|_2 \leq (1 + \|\mathbf{b}_0^*\|_\infty)/\sqrt{N}.$$

can shrink to zero in the limit, for example $\|\widehat{\mathbf{w}}^{\text{SA}} - \mathbf{w}_0^*\|_2 \asymp 1/\sqrt{N}$ in the simplest case when K is finite. It is thus only non-trivial if we manage to show $\|\widehat{\mathbf{w}} - \mathbf{w}_0^*\|_1 = o_p(1)$ and $\|\widehat{\mathbf{w}} - \mathbf{w}_0^*\|_2 = o_p(N^{-1/2})$, which is stated in the following Corollary 1.

Corollary 1 *Under the assumptions in Theorem 2, we have*

$$\|\widehat{\mathbf{w}} - \mathbf{w}_0^*\|_2 = O_p(N^{-1/2}K^2\tau) = o_p(N^{-1/2}) \quad \text{and} \quad \|\widehat{\mathbf{w}} - \mathbf{w}_0^*\|_1 = O_p(K^2\tau) = o_p(1).$$

Remark 4.4. To connect the primal problem with the dual problem, we consider the decomposition:

$$\widehat{\mathbf{w}} - \mathbf{w}_0^* = \widehat{\mathbf{A}}(\widehat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}^*) + (\widehat{\mathbf{A}} - \mathbf{A}^*)\boldsymbol{\alpha}^*;$$

see (A.39) in the appendix. The ℓ_2 -norm of the first term on the right-hand side is bounded by Theorem 2. Moreover, the first term dominates the ℓ_2 -norm of the second term in which the magnitude of $(\widehat{\mathbf{A}} - \mathbf{A}^*)$ is well controlled under Assumption 1.

Corollary 1 establishes meaningful convergence of the norms of the $\widehat{\mathbf{w}}$ to its oracle counterpart \mathbf{w}_0^* . The convergence further implies a desirable oracle inequality in Theorem 3 below, which shows that the empirical risk under $\widehat{\mathbf{w}}$ is asymptotically as small as if we knew the oracle object $\widehat{\Sigma}^*$.

Theorem 3 (Oracle inequalities) *Under the assumptions in Theorem 2, we have*

$$(a) \quad \widehat{\mathbf{w}}' \widehat{\Sigma} \widehat{\mathbf{w}} \leq \mathbf{w}_0^{*'} \widehat{\Sigma}^* \mathbf{w}_0^* + O_p(\tau K^{5/2}).$$

Furthermore, let $\widehat{\Sigma}^{\text{new}}$ and $\widehat{\Sigma}^{*\text{new}}$ be the counterparts of $\widehat{\Sigma}$ and $\widehat{\Sigma}^*$ from a new (testing) sample, which can be dependent or independent of the training dataset used to estimate $\widehat{\mathbf{w}}$ and \mathbf{w}_0^* . If the testing dataset is generated by the same data generating process as that of the training dataset, then

$$(a) \quad \widehat{\mathbf{w}}' \widehat{\Sigma}^{\text{new}} \widehat{\mathbf{w}} \leq \mathbf{w}_0^{*'} \widehat{\Sigma}^{*\text{new}} \mathbf{w}_0^* + O_p(\tau K^{5/2});$$

$$(b) \quad \widehat{\mathbf{w}}' \widehat{\Sigma} \widehat{\mathbf{w}} \leq \widehat{\mathbf{w}}' \widehat{\Sigma}^{\text{new}} \widehat{\mathbf{w}} \leq Q(\Sigma_0) + O_p(\tau K^{5/2}), \text{ where } \Sigma_0 = \Sigma_0^* + \Sigma_0^e \text{ and } Q(\Sigma_0) := \min_{\mathbf{w}' \mathbf{1}_N = 1} \mathbf{w}' \Sigma_0 \mathbf{w}.$$

Remark 4.5. Theorem 3(a) is an in-sample oracle inequality, and (b) is an out-of-sample oracle inequality. The proof is a term-by-term analysis of the difference between $\widehat{\mathbf{w}}' \widehat{\Sigma} \widehat{\mathbf{w}}$ and $\mathbf{w}_0^{*'} \widehat{\Sigma}^* \mathbf{w}_0^*$ caused by the idiosyncratic shock. Again, because the magnitude of the idiosyncratic shock is controlled by Assumption 1, the convergence of the weight estimator in Corollary 1 allows the sample risk $\widehat{\mathbf{w}}' \widehat{\Sigma} \widehat{\mathbf{w}}$ to approximate the oracle risk $\mathbf{w}_0^{*'} \widehat{\Sigma}^* \mathbf{w}_0^*$. The approximation is nontrivial by noting that $\mathbf{w}_0^{*'} \widehat{\Sigma}^* \mathbf{w}_0^*$ and $\mathbf{w}_0^{*'} \widehat{\Sigma}^{*\text{new}} \mathbf{w}_0^*$ are bounded away from 0 given the low rank structure of $\widehat{\Sigma}^{*\text{new}}$ and that $\tau K^{5/2} \rightarrow 0$ under Assumption 2(a). In other words, the risk of our sample estimator would be as low as if we were informed of the infeasible oracle group membership, up to an asymptotically negligible term.

Remark 4.6. While our ℓ_2 -relaxation regularizes the combination weights, there is another line of literature of regularizing the high dimensional VC estimation or its inverse (the precision matrix); see Bickel and Levina (2008), Fan et al. (2013), and the overview by Fan et al. (2016). Theorem 3(c) implies that the in-sample and out-of-sample risks coming out of the ℓ_2 -relaxation are comparable with the resultant risk from estimating the high dimensional VC matrix. Specifically, $Q(\Sigma_0)$ is Bates and Granger (1969)'s optimal risk under the population VC Σ_0 , and Σ_0 is the target of high dimensional VC estimation or precision matrix estimation. The population VC Σ_0 takes into account both the low rank component Σ_0^* and the high rank component Σ_0^e . Even if Σ_0 can be estimated so well that the estimation error is completely eliminated, our out-of-sample risk $\widehat{\mathbf{w}}' \widehat{\Sigma}^{\text{new}} \widehat{\mathbf{w}}$ is within an $O_p(\tau K^{5/2})$ tolerance level of $Q(\Sigma_0)$.

In summary, we have shown that under the high dimensional asymptotic framework where $N/T \rightarrow \infty$ is allowed as $(N, T) \rightarrow \infty$, we can construct a unique oracle target $\widehat{\alpha}^*$ that satisfies a set of desirable properties. Since the dual problem (2.6) is an ℓ_1 -penalized optimization, we establish in Theorem 1 the convergence of $\widehat{\alpha}$ to $\widehat{\alpha}^*$ by statistical techniques that deal with the ℓ_1 -regularization, thanks to the amenable comparability condition and the derived restricted eigenvalue. Then Theorem 2 can be extended to the convergence of the weight $\widehat{\mathbf{w}}$ in Corollary 1 and furthermore the convergence of the sample risk to the oracle risks in Theorem 3.

5 Monte Carlo Simulations

We illustrate the performance of the proposed ℓ_2 -relaxation method via Monte Carlo simulations. For simplicity, the simulated DGP follows a group pattern, with an equal number of members in each group, i.e., $N_k = N/K$ for each $k \in [K]$. We start with a baseline model of independent factors, and then extend it to allow dynamic factors and approximate factors.

5.1 Simulation Design

Let Ψ^{co} be a $K \times K$ symmetric positive definite matrix, and $\Psi = \Psi^{\text{co}} \otimes (\mathbf{1}_{N_1} \mathbf{1}'_{N_1})$ be its $N \times N$ equicorrelation matrix. The N forecasters are generated by

$$\mathbf{f}_t = \Psi^{1/2} \boldsymbol{\eta}_t + \mathbf{u}_t, \quad (5.1)$$

where $\Psi^{1/2} = N_1^{-1/2} (\Psi^{\text{co}})^{1/2} \otimes (\mathbf{1}_{N_1} \mathbf{1}'_{N_1})$, $\boldsymbol{\eta}_t \sim N(\mathbf{0}, \mathbf{I}_N)$, and the idiosyncratic noise $\mathbf{u}_t \sim N(\mathbf{0}, \Omega_u)$ are independent across t , and $\boldsymbol{\eta}_t$ and \mathbf{u}_t are independent. The target variable is generated as $y_{t+1} = \mathbf{w}_\psi^* \Psi^{1/2} \boldsymbol{\eta}_t + u_{y,t+1}$, where $u_{y,t+1} \sim N(0, \sigma_y^2)$ is independent of $\boldsymbol{\eta}_t$ and \mathbf{u}_t , and $\mathbf{w}_\psi^* = [(\Psi^{\text{co}})^{-1} \mathbf{1}_K] \otimes \mathbf{1}_{N_1} / [N_1 \mathbf{1}'_K (\Psi^{\text{co}})^{-1} \mathbf{1}_K]$. The forecast error vector is

$$\mathbf{e}_t := y_{t+1} \mathbf{1}_N - \mathbf{f}_t = (\mathbf{1}_N \mathbf{w}_\psi^{*'} - \mathbf{I}_N) \Psi^{1/2} \boldsymbol{\eta}_t + u_{y,t+1} \mathbf{1}_N - \mathbf{u}_t,$$

and its VC can be written as

$$E[\mathbf{e}_t \mathbf{e}_t'] = \underbrace{(\mathbf{I}_N - \mathbf{1}_N \mathbf{w}_\psi^{*'}) \Psi (\mathbf{I}_N - \mathbf{w}_\psi^* \mathbf{1}'_N)}_{\Sigma_0} + \underbrace{\sigma_y^2 \mathbf{1}_N \mathbf{1}'_N + \Omega_u}_{\Omega_0}.$$

By construction \mathbf{w}_ψ^* solves $\min_{\mathbf{w}' \mathbf{1}_N = 1} \mathbf{w}' \Sigma_0 \mathbf{w}$.

We compare the following estimators of \mathbf{w} , all subject to the restriction $\mathbf{w}' \mathbf{1}_N = 1$: (i) the group-membership oracle estimator; (ii) simple averaging (SA); (iii) Lasso; (iv) Ridge; (v) the principle component (PC) grouping estimator; and (vi) ℓ_2 -relaxation. All the tuning parameters are obtained by the conventional 5-fold cross validation (CV) through a grid search from 0.1 to 1 with increment 0.1. We have also tried the 3-fold CV and the 10-fold CV, which yield similar results.

We elaborate the two estimators in which the group identity is explicitly involved. The oracle estimator takes advantage of the true group membership in the DGP. Given information about the group membership, we reduce the N forecasters to K forecasters $f_{(g_k),t} = |\mathcal{G}_k|^{-1} \sum_{i \in \mathcal{G}_k} f_{it}$ for $k \in [K]$, and use the low-dimensional (1.3) to find the optimal weights. We estimate the group membership in PC as follows. We compute the $T \times N$ in-sample forecasters' error matrix $\hat{\mathbf{E}} = (\hat{\mathbf{e}}_1, \dots, \hat{\mathbf{e}}_T)'$, save the associated $N \times N$ factor loading matrix $\hat{\mathbf{\Gamma}}$ of the singular decomposition $\hat{\mathbf{E}} = \hat{\mathbf{U}} \hat{\mathbf{D}} \hat{\mathbf{\Gamma}}'$, where $\hat{\mathbf{D}}$ is the 'diagonal' matrix of the singular values in descending order. We extract the the first q columns of $\hat{\mathbf{\Gamma}}$, and perform the standard k-means clustering algorithm to partition the factor loading vectors into K estimated groups $\hat{\mathcal{G}}_k$, $k \in [K]$. In the simulation we use the true K and try $q = 5, 10, 20$ to avoid tuning on these hyperparameters in this PC grouping procedure.

We estimated the weights $\hat{\mathbf{w}}$ with the training sample $\{(y_{t+1}, \mathbf{f}_t), t \in [T]\}$, and then cast the 1-step-ahead prediction $\hat{\mathbf{w}}' \mathbf{f}_{T+1}$ for y_{T+2} . The above exercise is repeated to evaluate the mean squared forecast error (MSFE) $E[(y_{T+2} - \hat{\mathbf{w}}' \mathbf{f}_{T+1})^2] - \sigma_y^2$ of each estimator, where the unpredictable components σ_y^2 in the MSFE is subtracted and the mathematical expectations are approximated by empirical averages of 1000 simulation replications. In addition, we report the mean absolute forecast error (MAFE), which is not covered by our theory, in Appendix B.4 for reference.

We experiment with three training sample sizes $T = 50, 100, 200$ with the corresponding $K =$

2, 4, 6 and $N = 100, 200, 300$, respectively. We specify

$$\Psi^{\text{co}} = \begin{bmatrix} 1 & 0.1 & 0 & \cdots & 0 \\ 0.1 & \frac{3}{2} & 0.1 & \cdots & 0 \\ 0 & 0.1 & 2 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & \frac{K+1}{2} \end{bmatrix}.$$

While $\Omega_u = \sigma_u^2 \mathbf{I}_N$ with σ_u fixed at $\sigma_u = 5$. To highlight the effect of the signal-to-noise ratio (SNR) on the forecast accuracy, we specify $\sigma_y = 1$ as the low-signal design (with SNR around 3:7) and $\sigma_y = 0.1$ as the high-signal design (with SNR around 7:3). Appendix B.2 details the formula of the SNR for our setting.

5.2 Simulation Results

DGP1. The baseline model sets $\eta_t \sim N(\mathbf{0}, \mathbf{I}_N)$ i.i.d. across t . Figure 2 illustrates the estimated weights of a typical replication. The four rows of sub-figures correspond to the oracle, the ℓ_2 -relaxation, Lasso, and ridge, respectively; the three columns represent the results under $K = 2, 4, 6$, respectively. For each sub-figure, the estimated weights are plotted against N . Figure 2 shows clearly that ℓ_2 -relaxation is capable of capturing the grouping patterns, although it does not explicitly classify individuals into groups. Such patterns are observed in neither Lasso nor Ridge.

Table 1: MSFE for DGP1: the Baseline Model

T	N	K	Oracle	SA	Lasso	Ridge	PC			ℓ_2 -relax
							$q = 5$	$q = 10$	$q = 20$	
<i>Panel A: Low SNR</i>										
50	100	2	0.189	1.039	0.722	1.389	0.785	0.801	0.848	0.252
100	200	4	0.179	3.259	0.403	1.072	1.189	1.267	1.503	0.251
200	300	6	0.070	4.615	0.171	0.324	1.399	1.323	1.363	0.116
<i>Panel B: High SNR</i>										
50	100	2	0.158	0.999	0.440	0.470	0.739	0.768	0.731	0.230
100	200	4	0.137	3.134	0.185	0.252	1.108	1.239	1.303	0.161
200	300	6	0.093	4.420	0.120	0.126	1.306	1.437	1.391	0.108

Table 1 reports the out-of-sample prediction accuracy under DGP1. The first three columns show the settings of T , N and K , and Columns 4–11 contain the MSFEs of the labeled estimators. Almost all estimators have stronger performance under a high SNR than that under a low SNR (except for the PC estimator with large q and T). The infeasible grouping information helps the oracle estimator to prevail in all cases. Our ℓ_2 -relaxation estimator is always the best feasible estimator, with MSFE close to that of the oracle estimator. Off-the-shelf shrinkage estimators Lasso and Ridge are in general better than the PC estimator. With the group pattern in the DGP, SA in general lags far behind the other feasible estimators that learn the combination weights from the data. Notice that the MSFEs by Oracle, Lasso, Ridge, and the ℓ_2 -relaxation decrease as T grows along with N . However, the results by SA and PC under all values of q diverge as (T, N) increases.

DGP2. Time dependence is an essential feature in time series forecast. We extend the baseline i.i.d. model by allowing for temporal serial dependence in $\{\eta_t\}$. Specifically, for each i , we generate

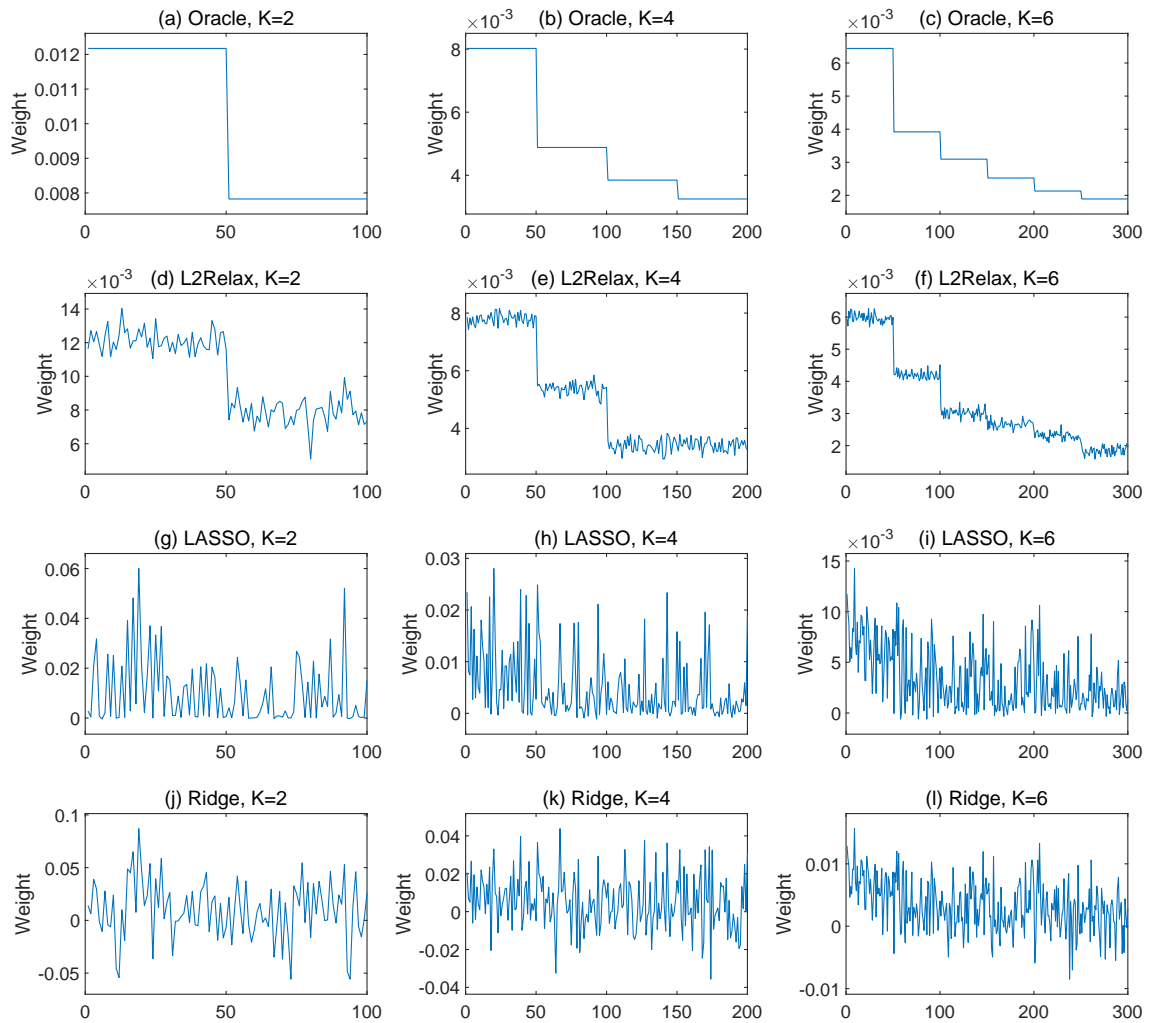


Figure 2: Illustration of Estimated Weights in DGP1

η_{it} from an AR(1) model

$$\eta_{it} = \rho_i \eta_{i,t-1} + \epsilon_{it}^\eta,$$

where $\rho_i \sim \text{Uniform}(0, 0.9)$ is a random autoregressive coefficient, the noise $\epsilon_{it}^\eta \sim \text{i.i.d. } N(0, 1 - \rho_i^2)$, and the initial values $\eta_{i0} \sim \text{iid } N(0, 1)$. This design allows for heterogeneity in $\{\rho_i\}$ while maintains strictly stationarity over t .

Table 2: MSFE for DGP2: the Dynamic Factor Model

T	N	K	Oracle	SA	Lasso	Ridge	PC			ℓ_2 -relax
							$q = 5$	$q = 10$	$q = 20$	
<i>Panel A: Low SNR</i>										
50	100	2	0.259	1.129	0.644	1.174	0.834	0.828	0.939	0.276
100	200	4	0.161	2.980	0.397	1.178	1.138	1.257	1.218	0.168
200	300	6	0.117	4.311	0.241	0.425	1.368	1.430	1.339	0.159
<i>Panel B: High SNR</i>										
50	100	2	0.263	1.078	0.476	0.485	0.769	0.838	0.820	0.267
100	200	4	0.149	3.101	0.200	0.261	1.081	1.202	1.291	0.162
200	300	6	0.096	4.452	0.130	0.130	1.265	1.367	1.459	0.112

Intuitively, the conventional 5-fold CV that randomly permutes the data may not be suitable for time series data, since it accounts for neither the chronological order nor the serial correlation of the data. Practitioners usually resort to the out-of-sample (OOS) evaluation instead; see Bergmeir and Benítez (2012) and Mirakyan et al. (2017), among others. See also Arlot and Celisse (2010) for a survey of cross-validation procedures for model selection.

We adopt the following procedure among many variations of OOS. The full dataset, in the chronological order, is divided into 5 blocks of equal sized (up to rounding of decimals). In the first iteration, block 1 is the training data for model estimation and blocks 2 make the evaluation data for computing the MFSE. In the second iteration, blocks 1-2 are taken as the training data and blocks 3 serve as the evaluation data, and so on for 4 iterations in total. Over a pre-specified grid of the potential values of the tuning parameter, the one that minimizes the average MSFE over the 4 iterations is selected.

The MSFEs for DGP2 are presented in Table 2. Apparently, the rankings of relative performance among the six estimators are similar to that in Table 1. ℓ_2 -relaxation remains the best performer among the feasible estimator in all cases, only second to the infeasible oracle estimator.

Our results are not sensitive to different evaluation methods. Bergmeir et al. (2018) argue the standard 5-fold CV is valid in purely autoregressive models with uncorrelated errors. Simulation results for DGP2 by the conventional 5-fold CV are reported in Appendix B.3.

DGP3. Acknowledging that equal factor loadings within a group can be an approximation of more general factor loading configurations, we experiment with a second extension to the baseline model in this DGP. We define $\tilde{\Psi}^{1/2} = \Psi^{1/2} + \text{iid } N(0, N_1^{-1/2})$ as a perturbed factor loading matrix to replace $\Psi^{1/2}$ in (5.1). The results are reported in Table 3. Although the additional noises in the factor loadings enlarge the MSFEs of all estimators relative to their counterparts in Table 1, the ranking pattern of the estimators stays the same as in DGP1.

In summary, we observe consistently robust performance of ℓ_2 relaxation superior to the other feasible estimators across the DGP designs, signal strength, and CV methods.

Table 3: MSFE for DGP3: Approximate Factor Model

T	N	K	Oracle	SA	Lasso	Ridge	PC			ℓ_2 -relax
							$q = 5$	$q = 10$	$q = 20$	
<i>Panel A: Low SNR</i>										
50	100	2	0.372	1.054	1.076	1.668	0.848	0.949	0.868	0.504
100	200	4	0.340	3.470	0.745	1.580	1.621	1.754	1.806	0.445
200	300	6	0.335	4.635	0.552	0.920	2.277	2.283	2.288	0.420
<i>Panel B: High SNR</i>										
50	100	2	0.392	1.059	0.722	0.726	0.895	0.909	0.915	0.422
100	200	4	0.280	3.335	0.403	0.508	1.641	1.791	1.956	0.306
200	300	6	0.274	5.226	0.347	0.343	2.174	2.242	2.260	0.309

6 Empirical Applications

We apply our proposed method to three empirical examples. The first exercise is a microeconomic study of forecasting box office. The second one is a macroeconomic application to the survey of professional forecasters. The last one conducts a one-period-ahead forecast of the realized volatility of the S&P500 index. In all three cases, we employ MSFE as the criterion to assess the forecasting performance. The results under MAFE are available in Appendix B.5. The ℓ_2 -relaxation enjoys excellent performance in these examples.

6.1 Box Office

The motion picture industry devotes enormous resources to marketing in order to influence consumer sentiment toward their products. These resources intend to reduce the supply-demand friction on the market. On the supply side, movie making is an expensive business; on the demand side, however, the audience’s taste is notoriously hard to catch. Accurate prediction of box office is financially crucial for motion picture investors.

Based on the data of Hollywood movies released in North America between October 1, 2010 and June 30, 2012, Lehrer and Xie (2017) demonstrate the sound OOS performance of the *prediction model averaging* (PMA). We revisit their dataset of 94 cross-sectional observations (movies), 28 non-constant explanatory variables and 95 candidate forecasters according to a multitude of model specifications. Guided by the intuition that the input variables capturing similar characteristics are “closer” to one another, Lehrer and Xie (2017) cluster input variables into six groups in their Appendix D.1:

Key variables : Constant, Animation, Family, Weeks, Screens, VOL : T - 1 / - 3

Twitter Volume : T - 21 / - 27, T - 14 / - 20, T - 7 / - 13, T - 4 / - 6

Twitter Sentiment : T - 21 / - 27, T - 14 / - 20, T - 7 / - 13, T - 4 / - 6, T - 1 / - 3

Rating Related : PG, PG13, R, Budget

Male Genre : Action, Adventure, Crime, Fantasy, Sci - Fi, Thriller

Female Genre : Comedy, Drama, Mystery, Romance

Since the 95 forecasters are generated based on these input variables, the potential group patterns may help the ℓ_2 -relaxation achieve more accurate forecasts than other off-the-shelf machine learning

shrinkage methods in this setting.

Following Lehrer and Xie (2017), we randomly rearrange the full sample with $n = 94$ movies into a training set of size n_{tr} and an evaluation set of size $n_{ev} = n - n_{tr}$. We consider evaluation sizes $n_{ev} = 10, 20, 30,$ and 40 . We repeat this procedure for 1,000 times and evaluate the OOS MSFE of PMA, Lasso, Ridge, and ℓ_2 -relaxation. Movies are viewed as independent observations and thus the tuning parameters are chosen by the conventional 5-fold CV. We conduct grid search from 0 to 5.

Table 4: Relative MSFE of Movie Forecasting

n_{ev}	PMA	Lasso	Ridge	ℓ_2 -relax
10	1.000	1.010	1.101	0.915
20	1.000	1.032	1.153	0.893
30	1.000	1.017	1.192	0.833
40	1.000	1.030	1.185	0.766

Note: The MSFE of PMA is normalized as 1.

Since the magnitude of MSFEs varies on the evaluation sizes n_{ev} , for convenience of comparison we report the mean risk relative to that of PMA in Table 4. Entries smaller than 1 indicate better performance relative to that of PMA. The results show that ℓ_2 -relaxation yields the lowest risk in each case, and it is followed by PMA which is known to outperform Lasso and Ridge in Lehrer and Xie (2017). The improvement by the ℓ_2 -relaxation from PMA increases with values of n_{ev} . In particular, when $n_{ev} = 40$ the MSFE of ℓ_2 -relaxation is 23.4% smaller than that of PMA.

The averaged weights across the 1000 repetitions generated by the ℓ_2 -relaxation method tell small variations in the relative importance amongst the 95 candidate models. As each of these forecast models are ensembles of several underlying input variables, it is more interesting to look at the partial weights of the 28 input characteristics. Figure 3 plots the averaged partial weights from the ℓ_2 -relaxation associated with these 28 variables under $n_{ev} = 10$, and the pattern is similar for other values of n_{ev} . The horizontal axis of Figure 3 shows all the 28 non-constant variables from the one with the highest partial weight (GENRE: Adventure) to the one with the lowest (GENRE: Thriller). The vertical axis represents the averaged weight. We shade the variables with partial weights either 1 or 0 exactly. These variables are either included or excluded in all candidate models following Lehrer and Xie (2017)’s pre-screening procedure.

We focus on the variables with averaged partial weights lie strictly between 0 and 1. Variables such as RATE-R, VOL: T-4/-6 and VOL: T-7/-13 yield similar averaged partial weights, and this is robust under other values of n_{ev} . The variable groups revealed by ℓ_2 -relaxation offer accurate prediction, whereas they do not coincide the manually selected groups. These findings open interesting possibilities for alternative categorization of the features of movies as well as their social media popularity.

6.2 Inflation

Firms, consumers as well as monetary policy authorities count on the outlook of inflation to make rational economic decisions. Besides model-based inflation forecasts published by government and research institutes, surveys of professional forecasters (SPF) report experts’ perceptions about the price level movement in the future. A long-standing myth of forecast combination lies in the robustness of the simple average, documented by Ang et al. (2007) which extract the mean or median as a predictor in a simple linear regression. Recent research shows modern machine learning methods can assist by assigning data-driven weights to individual forecasters to gather disaggregated information,

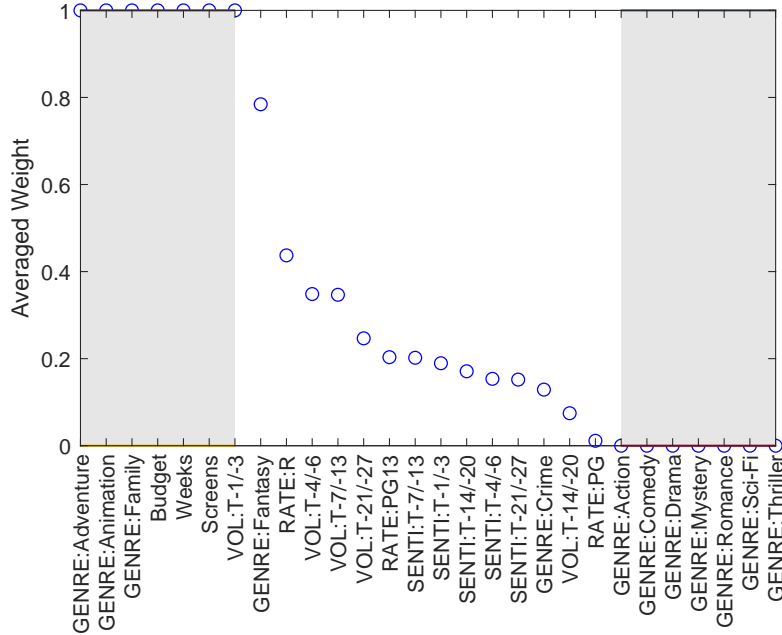


Figure 3: Averaged Partial Weights for All Variables by ℓ_2 -relaxation under $n_E = 10$

for example Diebold and Shin (2018).

The European Central Bank’s SPF inquires many professional institutions for their expectations of the euro-zone macroeconomic outlook. We revisit Genre et al. (2013)’s harmonized index of consumer prices (HICP) dataset, which covers 1999Q1–2018Q4. The experts were asked about their one-year- and two-year-ahead predictions. The raw data record 119 forecasters in total, but are highly unbalanced with many missing values, mainly due to entry and exit in the long time span. We follow Genre et al. (2013) to obtain 30 qualified forecasters by first filtering out irregular respondents if he or she missed more than 50% of the observations, and then using a simple AR(1) regression to interpolate the missing values in the middle.

Our benchmark is the simple averaging (SA) with equal weights on all 30 forecasters. We compare the forecast errors of SA, Lasso, Ridge, and ℓ_2 -relaxation. We use a rolling window of 40 quarters for estimation. The tuning parameters are selected by the OOS CV approach described the simulation of DGP2.

Table 5: Relative MSFE of HICP Forecasting

Horizon	SA	Lasso	Ridge	ℓ_2 -relax
One-year-ahead	1.000	0.954	0.964	0.940
Two-year-ahead	1.000	0.887	0.969	0.518

Note: The MSFE of SA is normalized as 1.

The results of relative risks are presented in Table 5, with the MSFEs of SA standardized as 1. ℓ_2 -relaxation attains the best performance in both horizons. While Lasso and Ridge yield close results to SA under one-year-ahead horizon, ℓ_2 -relaxation improves the relative MSFE by a moderate 6%. Moreover, under two-year-ahead horizon ℓ_2 -relaxation beats SA by a whopping gain of 48.2%.

Since we do not directly observe the underlying factors based upon which the forecasters make

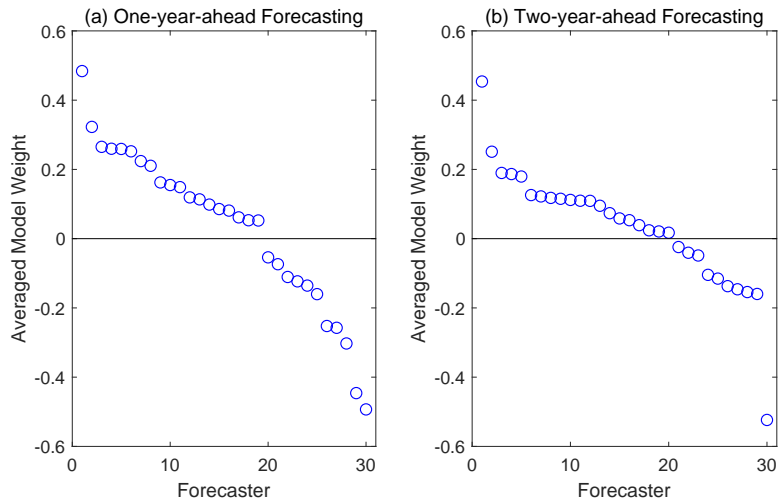


Figure 4: Averaged Model Weights for Different Forecast Horizons by ℓ_2 -relaxation

decisions, we illustrate in Figure 4 the weights associated with the 30 forecasters, averaged over the rolling windows. The figure suggests that the 30 forecasters roughly fall into a few groups by their weight estimates. All weights lie within the range of $(-1, 1)$, though some are negative. Forecasters are faced with more uncertainty in a longer horizon. As a result, the weights in Subplot (b) are more concentrated around zero than those in Subplot (a), reflecting the diminishing disparity in predictability.

6.3 Stock Volatility

Modeling the realized volatility (RV) of financial assets is of great interest to practitioners in risk management and portfolio allocation. To capture the main characteristics of this variance time series, a plethora of models have been proposed, such as Andersen et al. (2001)'s *fractionally integrated autoregressive moving average model* and Corsi (2009)'s *heterogeneous autoregressive model* (HAR). In particular, HAR owes its popularity to the simple linear regression structure. HAR linearly projects the variance of discretely sampled returns onto the linear space spanned by the lagged squared return over the identical return horizon in combination with the squared returns over longer and/or shorter return horizons. The standard HAR model writes the h -step-ahead daily y_{t+h} in the linear regression form

$$y_{t+h} = \beta_0 + \beta_d y_t^{(1)} + \beta_w y_t^{(5)} + \beta_m y_t^{(22)} + \epsilon_{t+h}, \quad (6.1)$$

where $y_t^{(l)} := l^{-1} \sum_{s=1}^l y_{t-s}$ is the averages of the previous l periods of RV from period t . A typical choice for the index vector \mathbf{l} is $[1, 5, 22]$, which allows the model to mirror the daily, weekly, and monthly components of the volatility process. The somewhat arbitrary index vector $\mathbf{l} = [1, 5, 22]$ demonstrates strong OOS performance in practice. Such a phenomenon implies that the RV data may exhibit daily, weekly, and monthly grouping patterns, which coincide with the trading behavior in the stock market for short, medium, and long return horizons.

In this exercise, we explore the OOS performance of ℓ_2 -relaxation with a full index vector $[1, 2, \dots, 22]$ relative to the benchmark HAR with $\mathbf{l} = [1, 5, 22]$. The sample period for the S&P500 RV spans from January 2, 2018 to December 31, 2018, with 250 observations in total. We consider a simple h -step-ahead rolling window forecasting with the window length 100. The daily, weekly, bi-weekly, and monthly forecasts with $h = 1, 5, 10$, and 22 are analyzed. For comparison, we also run the

h -step ahead forecasts by the simple OLS, Lasso and Ridge when all the 22 variables $\{y_t^{(1)}, \dots, y_t^{(22)}\}$ are used. We do not expect each individual regressor to be an unbiased estimate of the dependent variable, and hence for all estimators we drop the restriction that the coefficients add up to 1. All tuning parameters are determined using the OOS CV as for the simulation in DGP2.

Table 6: Relative Performance of Volatility Forecasting

horizon h	HAR	OLS	Lasso	Ridge	ℓ_2 -relax
1 (daily)	1.0000	1.4401	0.9500	0.9555	0.8654
5 (weekly)	1.0000	1.3691	0.9240	0.8564	0.8207
10 (bi-weekly)	1.0000	1.5419	1.0034	0.9474	0.8793
22 (monthly)	1.0000	1.7485	1.1737	1.1622	0.9012

Note: The MSFE of HAR is normalized as 1.

Table 6 shows the results of out-of-sample MSFEs relative to the benchmark HAR. Entries smaller than 1 indicate better performance relative to HAR. There are several interesting findings. (i) OLS suffers the worst performance in all cases, due to overfitting by the 22-variable full index vector in the 100-day rolling windows. (ii) Shrinkage estimators mitigate overfitting. While HAR outperforms OLS by fusing the daily, weekly and monthly lagged observations to reduce the number of regressors from 22 to 3, the pre-selected index vector $[1, 5, 22]$ cannot beat the shrinkage estimators in the short horizons. Ridge outperforms Lasso in general, indicating the underlying coefficients may take small values but may not be exactly zero. (iii) ℓ_2 -relaxation is more accurate than the other methods under all horizons. In particular, as noises accumulate under the prolonged horizon $h = 22$, HAR becomes more stable than Lasso and Ridge, whereas ℓ_2 -relaxation still prevails.

Panels (a)–(d) in Figure 5 represent the averaged (over the rolling windows) coefficients for each forecast horizon, respectively. The solid line, circles, and dashed line indicate the benchmark HAR, ℓ_2 -relaxation and OLS, respectively. Since HAR incorporates a fixed lag index $\mathbf{l} = [1, 5, 22]$, the estimated coefficients naturally follow a step function that consists of three segments: Lag 1, Lags 2–5 and Lags 6–22. On the other hand, OLS yields volatile estimated leading to poor performance. ℓ_2 -relaxation exhibits its own patterns distinctive from those of HAR and OLS. The curve of averaged coefficients are smooth, balancing the flexibility of OLS and the parsimony of HAR. Furthermore, in short horizons the estimated coefficients monotonically shrink to zero as the lags increases, reflecting the dampening effects of past observations.

7 Conclusion

This paper presents a new machine learning algorithm, the ℓ_2 -relaxation for forecast combination in the presence of many forecasts. When the forecast error VC matrix can be approximated by a block equicorrelation structure, we establish its asymptotic optimality in the high dimensional context when N is potentially larger than T . Simulations and real data applications demonstrate the excellent performance of the ℓ_2 -relaxation method relative to Lasso, Ridge and SA estimators.

The present work raises several interesting issues for further research. First, additional forms of restrictions can be imposed to accompany the ℓ_2 -relaxation. For example, if sparsity is also desirable, we may consider adding another constraint $\|\mathbf{w}\|_1 \leq \tau_1$ for some tuning parameter τ_1 , which is similar to the idea of mixed ℓ_2 - ℓ_1 penalty of elastic net (Zou and Hastie, 2005). For another example, if non-negative weights are desirable, we can add the constraints $w_i \geq 0$ for all i , similar to Jagannathan and Ma (2003). Second, as mentioned in the introduction, the idea of ℓ_2 -relaxation can also be applied to portfolio optimization problems, and it is interesting to investigate how this

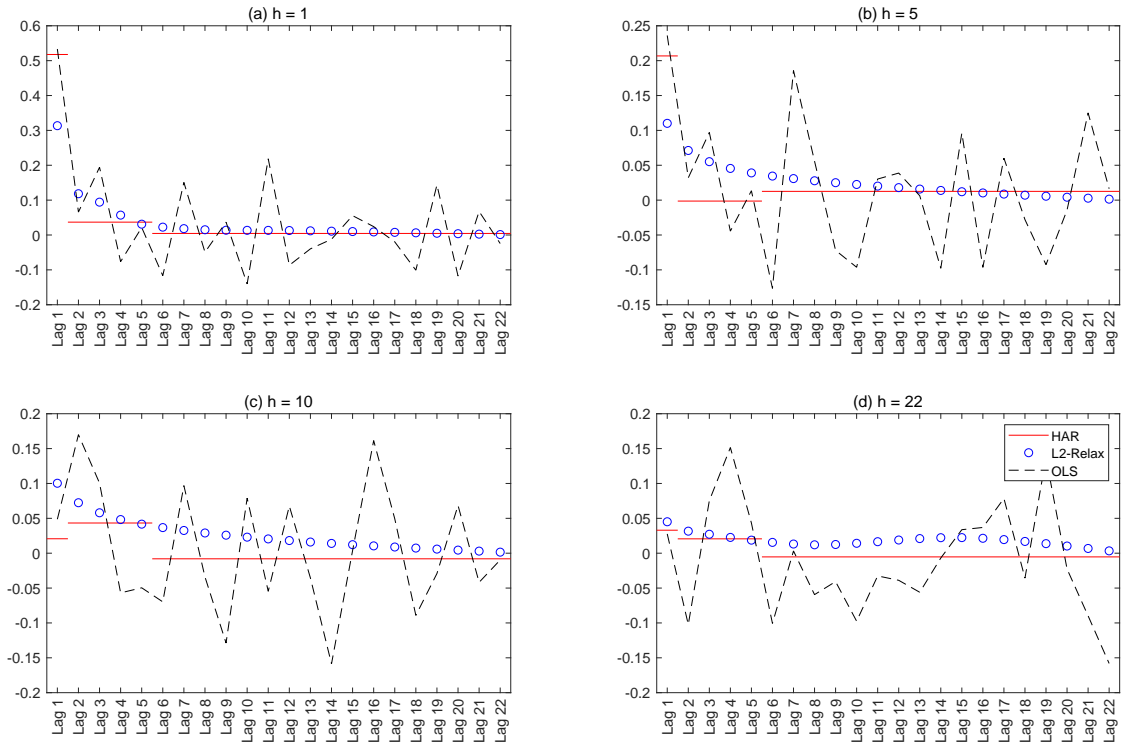


Figure 5: Averaged Coefficients by ℓ_2 -relaxation

approach works in comparison with shrinkage or regularization methods used to estimate the optimal portfolios. Third, our ℓ_2 -relaxation is motivated from the MSFE loss function, it is possible to consider other forms of relaxation if the other forms of loss functions (e.g., MAFE) are under investigation. We shall explore these topics in future work.

References

- Andersen, T. G., T. Bollerslev, F. X. Diebold, and P. Labys (2001). The distribution of realized exchange rate volatility. *Journal of the American Statistical Association* 96(453), 42–55.
- Ang, A., G. Bekaert, and M. Wei (2007). Do macro variables, asset markets, or surveys forecast inflation better? *Journal of Monetary Economics* 54(4), 1163–1212.
- Ao, M., L. Yingying, and X. Zheng (2019). Approaching mean-variance efficiency for large portfolios. *The Review of Financial Studies* 32(7), 2890–2919.
- Arlot, S. and A. Celisse (2010). A survey of cross-validation procedures for model selection. *Statistics Surveys* 4, 40–79.
- Bates, J. M. and C. W. Granger (1969). The combination of forecasts. *Operational Research Quarterly*, 451–468.
- Bayer, S. (2018). Combining value-at-risk forecasts using penalized quantile regressions. *Econometrics and Statistics* 8, 56–77.

- Belloni, A., D. Chen, V. Chernozhukov, and C. Hansen (2012). Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica* 80, 2369–2429.
- Bergmeir, C. and J. M. Benítez (2012). On the use of cross-validation for time series predictor evaluation. *Information Sciences* 191, 192 – 213. Data Mining for Software Trustworthiness.
- Bergmeir, C., R. J. Hyndman, and B. Koo (2018). A note on the validity of cross-validation for evaluating autoregressive time series prediction. *Computational Statistics & Data Analysis* 120, 70–83.
- Bickel, P., Y. Ritov, and A. Tsybakov (2009). Simultaneous analysis of Lasso and Dantzig selector. *The Annals of Statistics* 37(4), 1705–1732.
- Bickel, P. J. and E. Levina (2008). Regularized estimation of large covariance matrices. *The Annals of Statistics* 36(1), 199–227.
- Bonhomme, S., T. Lamadon, and E. Manresa (2017). Discretizing unobserved heterogeneity. *University of Chicago, Becker Friedman Institute for Economics Working Paper* (2019-16).
- Bonhomme, S. and E. Manresa (2015). Grouped patterns of heterogeneity in panel data. *Econometrica* 83(3), 1147–1184.
- Boyd, S. and L. Vandenberghe (2004). *Convex optimization*. Cambridge University Press.
- Bühlmann, P. and S. van de Geer (2011). *Statistics for High-Dimensional Data: Methods, Theory and Applications*. Springer.
- Candes, E. and T. Tao (2007). The Dantzig selector: Statistical estimation when p is much larger than n . *The Annals of Statistics* 35(6), 2313–2351.
- Chan, F. and L. L. Pauwels (2018). Some theoretical results on forecast combinations. *International Journal of Forecasting* 34(1), 64–74.
- Claeskens, G., J. R. Magnus, A. L. Vasnev, and W. Wang (2016). The forecast combination puzzle: A simple theoretical explanation. *International Journal of Forecasting* 32(3), 754–762.
- Clemen, R. T. (1989). Combining forecasts: A review and annotated bibliography. *International Journal of forecasting* 5(4), 559–583.
- Conflitti, C., C. De Mol, and D. Giannone (2015). Optimal combination of survey forecasts. *International Journal of Forecasting* 31(4), 1096–1103.
- Corsi, F. (2009). A simple approximate long-memory model of realized volatility. *Journal of Financial Econometrics* 7(2), 174–196.
- Coulombe, P. G., M. Leroux, D. Stevanovic, and S. Surprenant (2020). How is machine learning useful for macroeconomic forecasting? Technical report, CIRANO.
- Cowles, A. (1933). Can stock market forecasters forecast? *Econometrica*, 309–324.
- Diebold, F. X. and M. Shin (2018). Machine learning for regularized survey forecast combination: Partially-egalitarian lasso and its derivatives. *International Journal of Forecasting*.

- Diebold, F. X. and M. Shin (2019). Machine learning for regularized survey forecast combination: Partially-egalitarian lasso and its derivatives. *International Journal of Forecasting* 35(4), 1679–1691.
- Disatnik, D. and S. Katz (2012). Portfolio optimization using a block structure for the covariance matrix. *Journal of Business Finance & Accounting* 39(5-6), 806–843.
- Elliott, G., A. Gargano, and A. Timmermann (2013). Complete subset regressions. *Journal of Econometrics* 177(2), 357–373.
- Engle, R. and B. Kelly (2012). Dynamic equicorrelation. *Journal of Business & Economic Statistics* 30(2), 212–228.
- Fan, J., A. Furger, and D. Xiu (2016). Incorporating global industrial classification standard into portfolio allocation: A simple factor-based large covariance matrix estimator with high-frequency data. *Journal of Business & Economic Statistics* 34(4), 489–503.
- Fan, J. and R. Li (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* 96(456), 1348–1360.
- Fan, J., Y. Liao, and H. Liu (2016). An overview of the estimation of large covariance and precision matrices. *The Econometrics Journal* 19(1), C1–C32.
- Fan, J., Y. Liao, and M. Mincheva (2013). Large covariance estimation by thresholding principal orthogonal complements. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 75(4), 603–680.
- Fan, J., J. Zhang, and K. Yu (2012). Vast portfolio selection with gross-exposure constraints. *Journal of the American Statistical Association* 107(498), 592–606.
- Gao, Z. and Z. Shi (2020). Implementing convex optimization in r: Two econometric examples. *Computational Economics*.
- Genre, V., G. Kenny, A. Meyler, and A. Timmermann (2013). Combining expert forecasts: Can anything beat the simple average? *International Journal of Forecasting* 29(1), 108 – 121.
- Granger, C. W. and R. Ramanathan (1984). Improved methods of combining forecasts. *Journal of Forecasting* 3(2), 197–204.
- Hsiao, C. and S. K. Wan (2014). Is there an optimal forecast combination? *Journal of Econometrics* 178, 294–309.
- Jagannathan, R. and T. Ma (2003). Risk reduction in large portfolios: Why imposing the wrong constraints helps. *The Journal of Finance* 58(4), 1651–1683.
- Konzen, E. and F. A. Ziegelmann (2016). Lasso-type penalties for covariate selection and forecasting in time series. *Journal of Forecasting* 35(7), 592–612.
- Kotchoni, R., M. Leroux, and D. Stevanovic (2019). Macroeconomic forecast accuracy in a data-rich environment. *Journal of Applied Econometrics* 34(7), 1050–1072.
- Ledoit, O. and M. Wolf (2004). Honey, i shrunk the sample covariance matrix. *The Journal of Portfolio Management* 30(4), 110–119.

- Ledoit, O. and M. Wolf (2017). Nonlinear shrinkage of the covariance matrix for portfolio selection: Markowitz meets goldilocks. *The Review of Financial Studies* 30(12), 4349–4388.
- Lehrer, S. F. and T. Xie (2017). Box office buzz: Does social media data steal the show from model uncertainty when forecasting for hollywood? *The Review of Economics and Statistics* 99(5), 749–755.
- Li, J. and W. Chen (2014). Forecasting macroeconomic time series: Lasso-based approaches and their forecast combinations with dynamic factor models. *International Journal of Forecasting* 30(4), 996–1015.
- Mirakyan, A., M. Meyer-Renschhausen, and A. Koch (2017). Composite forecasting approach, application for next-day electricity price forecasting. *Energy Economics* 66, 228 – 237.
- Roccazzella, F., P. Gambetti, and F. Vrina (2020). Optimal and robust combination of forecasts via constrained optimization and shrinkage. Technical report, LFIN Working Paper Series, 2020/6, 1–2.
- Shi, Z. (2016). Econometric estimation with high-dimensional moment equalities. *Journal of Econometrics* 195(1), 104–119.
- Smith, J. and K. F. Wallis (2009). A simple explanation of the forecast combination puzzle. *Oxford Bulletin of Economics and Statistics* 71(3), 331–355.
- Stasinakis, C., G. Sermpinis, K. Theofilatos, and A. Karathanasopoulos (2016). Forecasting us unemployment with radial basis neural networks, kalman filters and support vector regressions. *Computational Economics* 47(4), 569–587.
- Stock, J. H. and M. W. Watson (2004). Combination forecasts of output growth in a seven-country data set. *Journal of Forecasting* 23(6), 405–430.
- Su, L. and G. Ju (2018). Identifying latent grouped patterns in panel data models with interactive fixed effects. *Journal of Econometrics* 206(2), 554–573.
- Su, L., Z. Shi, and P. C. Phillips (2016). Identifying latent structures in panel data. *Econometrica* 84(6), 2215–2264.
- Su, L., X. Wang, and S. Jin (2019). Sieve estimation of time-varying panel data models with latent structures. *Journal of Business & Economic Statistics* 37(2), 334–349.
- Theil, H. (1992). *Who forecasts best?*, Chapter International Economic Papers 5 (1955) 194–199, pp. 1115–1120. Springer.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 267–288.
- Tikhonov (1977). *Solutions of Ill-Posed Problems*. Winston and Sons, Washington, DC.
- Timmermann, A. (2006). Forecast combinations. *Handbook of Economic Forecasting* 1, 135–196.
- Van De Geer, S. A. and P. Bühlmann (2009). On the conditions used to prove oracle results for the lasso. *Electronic Journal of Statistics* 3, 1360–1392.

- Vogt, M. and O. Linton (2017). Classification of non-parametric regression functions in longitudinal data models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 79(1), 5–27.
- Vogt, M. and O. Linton (2020). Multiscale clustering of nonparametric regression curves. *Journal of Econometrics*.
- Wainwright, M. J. (2019). *High-dimensional statistics: A non-asymptotic viewpoint*, Volume 48. Cambridge University Press.
- Wang, W. and L. Su (2020). Identifying latent group structures in nonlinear panels. *Journal of Econometrics*.
- Wilms, I., J. Rombouts, and C. Croux (2018). Multivariate lasso-based forecast combinations for stock market volatility.
- Zou, H. and T. Hastie (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: series B (Statistical Methodology)* 67(2), 301–320.

APPENDIX

This appendix is composed of two sections. Appendix A contains the proofs of the theoretical statements in the paper. Appendix B contains further explanation and additional results on the simulation and empirical exercises.

A Proofs of the Main Results

A.1 Proof of the Result in Section 2

Proof of Lemma 1. First, we can rewrite the minimization problem in (2.5) in terms of linear constraints:

$$\begin{aligned} & \min_{(\mathbf{w}, \gamma) \in \mathbb{R}^{N+1}} \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t. } & \mathbf{w}'\mathbf{1}_N - 1 = 0, \quad \left(\widehat{\Sigma} \quad \mathbf{1}_N \right) \left(\mathbf{w}' \quad \gamma' \right)' \leq \tau \mathbf{1}_N, \text{ and } - \left(\widehat{\Sigma} \quad \mathbf{1}_N \right) \left(\mathbf{w}' \quad \gamma' \right)' \leq \tau \mathbf{1}_N \end{aligned} \quad (\text{A.1})$$

where “ \leq ” holds elementwise hereafter. Define the Lagrangian function as

$$\begin{aligned} \mathcal{L}(\mathbf{w}, \gamma; \boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \alpha_3) &= \frac{1}{2} \mathbf{w}'\mathbf{w} + \boldsymbol{\alpha}'_1 \left(\left(\widehat{\Sigma} \quad \mathbf{1}_N \right) \begin{pmatrix} \mathbf{w} \\ \gamma \end{pmatrix} - \tau \mathbf{1}_N \right) \\ &\quad - \boldsymbol{\alpha}'_2 \left(\left(\widehat{\Sigma} \quad \mathbf{1}_N \right) \begin{pmatrix} \mathbf{w} \\ \gamma \end{pmatrix} + \tau \mathbf{1}_N \right) + \alpha_3 (\mathbf{w}'\mathbf{1}_N - 1) \end{aligned} \quad (\text{A.2})$$

and the associated Lagrangian dual function as $g(\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \alpha_3) = \inf_{\mathbf{w}, \gamma} \mathcal{L}(\mathbf{w}, \gamma; \boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \alpha_3)$, where $\boldsymbol{\alpha}_1 \geq 0$, $\boldsymbol{\alpha}_2 \geq 0$, and α_3 are the Lagrangian multipliers for the three constraints in (A.1), respectively.

Let $\varphi(\mathbf{w}, \gamma) = \frac{1}{2} \|\mathbf{w}\|_2^2$, the objective function in (A.1). Define its conjugate function as

$$\varphi^*(\mathbf{a}, b) = \sup_{\mathbf{w}, \gamma} \left\{ \mathbf{a}'\mathbf{w} + b\gamma - \frac{1}{2} \|\mathbf{w}\|_2^2 \right\} = \begin{cases} \frac{1}{2} \|\mathbf{a}\|_2^2 & \text{if } b = 0 \\ \infty & \text{otherwise} \end{cases}.$$

The linear constraints indicate an explicit dual function (See Boyd and Vandenberghe, 2004, p.221):

$$\begin{aligned} & g(\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \alpha_3) \\ &= -\tau \mathbf{1}'_N (\boldsymbol{\alpha}_1 + \boldsymbol{\alpha}_2) - \alpha_3 - \varphi^* \left(\widehat{\Sigma} (\boldsymbol{\alpha}_2 - \boldsymbol{\alpha}_1) - \alpha_3 \mathbf{1}_N, \mathbf{1}'_N (\boldsymbol{\alpha}_2 - \boldsymbol{\alpha}_1) \right) \\ &= \begin{cases} -\tau \mathbf{1}'_N (\boldsymbol{\alpha}_1 + \boldsymbol{\alpha}_2) - \alpha_3 - \frac{1}{2} \left\| \widehat{\Sigma} (\boldsymbol{\alpha}_2 - \boldsymbol{\alpha}_1) - \alpha_3 \mathbf{1}_N \right\|_2^2 & \text{if } \mathbf{1}'_N (\boldsymbol{\alpha}_2 - \boldsymbol{\alpha}_1) = 0 \\ \infty & \text{otherwise} \end{cases}. \end{aligned}$$

Let $\boldsymbol{\alpha} = \boldsymbol{\alpha}_2 - \boldsymbol{\alpha}_1$. When $\tau > 0$, the two inequalities $\widehat{\Sigma}_i \mathbf{w} + \gamma \leq \tau$ and $-\widehat{\Sigma}_i \mathbf{w} - \gamma \leq \tau$ cannot be binding simultaneously. The associated Lagrangian multipliers α_{1i} and α_{2i} must satisfy $\alpha_{1i} \cdot \alpha_{2i} = 0$ for all $i \in [N]$. This implies that $\|\boldsymbol{\alpha}\|_1 = \mathbf{1}'_N \boldsymbol{\alpha}_1 + \mathbf{1}'_N \boldsymbol{\alpha}_2$ so that the dual problem can be simplified as

$$\max_{\boldsymbol{\alpha}, \alpha_3} \left\{ -\frac{1}{2} \left\| \widehat{\Sigma} \boldsymbol{\alpha} - \alpha_3 \mathbf{1}_N \right\|_2^2 - \alpha_3 - \tau \|\boldsymbol{\alpha}\|_1 \right\} \quad \text{s.t. } \mathbf{1}'_N \boldsymbol{\alpha} = 0. \quad (\text{A.3})$$

Taking the partial derivative of the above criterion function with respect to α_3 yields

$$(\widehat{\Sigma}\alpha - \alpha_3\mathbf{1}_N)'\mathbf{1}_N - 1 = 0,$$

or equivalently, $\alpha_3 = \frac{1}{N}(\mathbf{1}'_N\widehat{\Sigma}\alpha - 1)$. Then

$$\left\|\widehat{\Sigma}\alpha - \alpha_3\mathbf{1}_N\right\|_2^2 = \left\|\widehat{\mathbf{A}}\alpha + \frac{\mathbf{1}_N}{N}\right\|_2^2 = \alpha'\widehat{\mathbf{A}}'\widehat{\mathbf{A}}\alpha + \frac{1}{N}$$

where $\widehat{\mathbf{A}} = (\mathbf{I}_N - N^{-1}\mathbf{1}_N\mathbf{1}'_N)\widehat{\Sigma}$ as defined in the main text. We conclude that the dual problem in (A.3) is equivalent to

$$\min_{\alpha \in \mathbb{R}^N} \left\{ \frac{1}{2}\alpha'\widehat{\mathbf{A}}'\widehat{\mathbf{A}}\alpha + \frac{1}{N}\mathbf{1}'_N\widehat{\Sigma}\alpha + \tau\|\alpha\|_1 - \frac{1}{2N} \right\} \text{ subject to } \mathbf{1}'_N\alpha = 0, \quad (\text{A.4})$$

where we keep the constant $-\frac{1}{2N}$ which is irrelevant to the optimization.

When $\widehat{\alpha} = \widehat{\alpha}_2 - \widehat{\alpha}_1$ is the solution to (A.4), the solution of α_3 in (A.3) is $\widehat{\alpha}_3 = \frac{1}{N}(\mathbf{1}'_N\widehat{\Sigma}\widehat{\alpha} - 1)$. The first order condition of (A.2) with respect to \mathbf{w} evaluated at the solution gives

$$\mathbf{0}_N = \widehat{\mathbf{w}} + \widehat{\Sigma}(\widehat{\alpha}_1 - \widehat{\alpha}_2) + \widehat{\alpha}_3\mathbf{1}_N = \widehat{\mathbf{w}} - \widehat{\Sigma}\widehat{\alpha} + \frac{1}{N}(\mathbf{1}'_N\widehat{\Sigma}\widehat{\alpha} - 1)\mathbf{1}_N = \widehat{\mathbf{w}} - \widehat{\Sigma}\widehat{\alpha} - \frac{1}{N}\mathbf{1}_N.$$

as $\mathbf{1}'_N\widehat{\alpha} = 0$. The result in (2.7) follows. ■

A.2 Proofs of the Results in Section 3

Proof of Lemma 2. Part (a). Suppose \mathbf{w} is a feasible point. Due to the block equicorrelation structure, the N rows in $\widehat{\Sigma}^*\mathbf{w}$ contain only K distinct rows:

$$\widehat{\Sigma}^*\mathbf{w} = \begin{pmatrix} \widehat{\Sigma}_1^{\text{co}} \left(\sum_{i \in \mathcal{G}_1} w_i, \dots, \sum_{i \in \mathcal{G}_K} w_i \right)' \cdot \mathbf{1}_{N_1} \\ \vdots \\ \widehat{\Sigma}_K^{\text{co}} \left(\sum_{i \in \mathcal{G}_1} w_i, \dots, \sum_{i \in \mathcal{G}_K} w_i \right)' \cdot \mathbf{1}_{N_K} \end{pmatrix}.$$

We prove the result by contradiction. Suppose the elements in the k -th group of an optimizer $\mathbf{w}_\tau^* = (w_{\tau,1}^*, \dots, w_{\tau,N}^*)'$ are unequal for some $k \in [K]$. As before, we frequently suppress the dependence of \mathbf{w}_τ^* and $w_{\tau,i}^*$ on τ and write them as \mathbf{w}^* and w_i^* , respectively. Then we can construct an alternative estimator $\check{\mathbf{w}}^* = (\check{w}_1^*, \dots, \check{w}_N^*)'$ such that

$$\check{w}_i^* = N_k^{-1} \sum_{j \in \mathcal{G}_k} w_j^* \text{ if } i \in \mathcal{G}_k \text{ for each } k \in [K].$$

It is easy to verify that $\check{\mathbf{w}}^{*'}\mathbf{1}_N = \sum_{i=1}^N \check{w}_i^* = \sum_{k=1}^K N_k \check{w}_i^* = \sum_{k=1}^K \sum_{j \in \mathcal{G}_k} w_j^* = \sum_{i=1}^N w_i^* = 1$ and $\left\|\widehat{\Sigma}^*\check{\mathbf{w}}^* + \gamma\mathbf{1}_N\right\|_\infty = \left\|\widehat{\Sigma}^*\mathbf{w}^* + \gamma\mathbf{1}_N\right\|_\infty \leq \tau$ by the fact that \mathbf{w}^* satisfies the constraints. That is, $\check{\mathbf{w}}^*$

is also feasible. On the other hand, by the Jensen inequality,

$$\begin{aligned}\|\tilde{\mathbf{w}}^*\|_2^2 &= \sum_{i=1}^N (\tilde{w}_i^*)^2 = \sum_{k=1}^K N_k \left(N_k^{-1} \sum_{j \in \mathcal{G}_k} w_j^* \right)^2 \\ &\leq \sum_{k=1}^K \sum_{j \in \mathcal{G}_k} (w_j^*)^2 = \|\mathbf{w}^*\|_2^2,\end{aligned}$$

where the inequality becomes a strict inequality as long as w_j^* , $j \in \mathcal{G}_k$, are not all equal for some $k \in [K]$. But this contradicts the presumption that \mathbf{w}^* is the minimizer. Therefore, $\mathbf{w}^* = \mathbf{w}_\tau^*$ must be in the form of (3.7) with $b_{\tau k}^* = r_k^{-1} \sum_{i \in \mathcal{G}_k} w_{\tau, i}^*$. Substituting (3.7) into (3.6), the minimization problem in (3.6) can be equivalently written as that in (3.8).

Part (b). Under $\tau = 0$, the restriction of (3.6) can be written as

$$\begin{pmatrix} \widehat{\Sigma}^* & \mathbf{1}_N \\ \mathbf{1}'_N & 0 \end{pmatrix} \begin{pmatrix} \mathbf{w} \\ \gamma \end{pmatrix} = \begin{pmatrix} \mathbf{0}_N \\ 1 \end{pmatrix}.$$

Since the rank of $\widehat{\Sigma}^*$ is K and there are an infinite number of solutions of (\mathbf{w}, γ) to the above system of $N + 1$ equations. However, the i -th equation and the j -th equation are exactly the same if i and j are in the same group, and the $(N + 1)$ -equation system can be reduced to a system of $K + 1$ equations:

$$\begin{pmatrix} \widehat{\Sigma}^{\text{co}} & \mathbf{1}_K \\ \mathbf{1}'_K & 0 \end{pmatrix} \begin{pmatrix} \sum_{i \in \mathcal{G}_1} w_i, \dots, \sum_{i \in \mathcal{G}_K} w_i, \gamma \end{pmatrix}' = \begin{pmatrix} \mathbf{0}_K \\ 1 \end{pmatrix}.$$

As $\widehat{\Sigma}^{\text{co}}$ is of full rank, the solution to the above $(K + 1)$ -equation system is unique:

$$\mathbf{b}_0^* \circ \mathbf{r} = \begin{pmatrix} \sum_{i \in \mathcal{G}_1} w_{0, i}^*, \dots, \sum_{i \in \mathcal{G}_K} w_{0, i}^* \end{pmatrix}' = \frac{(\widehat{\Sigma}^{\text{co}})^{-1} \mathbf{1}_K}{\mathbf{1}'_K (\widehat{\Sigma}^{\text{co}})^{-1} \mathbf{1}_K}.$$

where we use the fact that $\sum_{i \in \mathcal{G}_k} w_{\tau, i}^* = b_{\tau, k}^* r_k$ for each $k \in [K]$ and all $\tau \geq 0$ by the result in part (a). It follows that $\mathbf{b}_0^* = \mathbf{r}^{-1} \circ \frac{(\widehat{\Sigma}^{\text{co}})^{-1} \mathbf{1}_K}{\mathbf{1}'_K (\widehat{\Sigma}^{\text{co}})^{-1} \mathbf{1}_K}$. ■

Proof of Lemma 3. Part (a). We prove the result by contradiction. Suppose that there exists some $k \in [K]$ such that the k -th group of an optimizer $\boldsymbol{\alpha}^* = (\alpha_1^*, \dots, \alpha_N^*)'$ has elements of opposite signs, viz, there are $i, j \in \mathcal{G}_k$ such that $\alpha_i^* \alpha_j^* < 0$. Construct an alternative estimator $\check{\boldsymbol{\alpha}}^* = (\check{\alpha}_1^*, \dots, \check{\alpha}_N^*)'$, where

$$\check{\alpha}_i^* = N_k^{-1} a_k^* \text{ for } i \in \mathcal{G}_k \text{ and all } k \in [K]$$

where $a_k^* = \sum_{j \in \mathcal{G}_k} \alpha_j^*$. By construction, $\check{\boldsymbol{\alpha}}^* \in \mathbb{S}^{\text{all}}$ as it replaces each α_i^* with $i \in \mathcal{G}_k$ by the groupwise average. It is obvious that $\boldsymbol{\alpha}^{*'} \widehat{\mathbf{A}}^* \widehat{\mathbf{A}}^* \boldsymbol{\alpha} = \check{\boldsymbol{\alpha}}^{*'} \widehat{\mathbf{A}}^* \widehat{\mathbf{A}}^* \check{\boldsymbol{\alpha}}^*$ and $\mathbf{1}'_N \widehat{\Sigma}^* \boldsymbol{\alpha} = \mathbf{1}'_N \widehat{\Sigma}^* \check{\boldsymbol{\alpha}}^*$. On the other hand, by the triangle inequality

$$\|\check{\boldsymbol{\alpha}}^*\|_1 = \sum_{i=1}^N |\check{\alpha}_i^*| = \sum_{k=1}^K N_k \left| N_k^{-1} \sum_{j \in \mathcal{G}_k} \alpha_j^* \right| < \sum_{i=1}^N |\alpha_i^*| = \|\boldsymbol{\alpha}^*\|_1,$$

where the strict inequality follows from the fact that the elements in $\{\alpha_j^*, j \in \mathcal{G}_k\}$ change signs for some $k \in [K]$. As a result, the objective function of the dual problem in (3.10) is strictly larger when evaluated at α^* than that at $\check{\alpha}^*$. This contradicts the presumption that α^* is an optimizer of (3.10).

Part (b). For any $\alpha \in \mathbb{R}^N$, we have

$$\left\| \widehat{\mathbf{A}}^* \alpha \right\|_2^2 = \alpha' \widehat{\Sigma}^* \left(\mathbf{I}_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}'_N \right) \widehat{\Sigma}^* \alpha = \alpha' \widehat{\Sigma}^* \widehat{\Sigma}^* \alpha - \frac{1}{N} \left(\mathbf{1}'_N \widehat{\Sigma}^* \alpha \right)^2. \quad (\text{A.5})$$

The group structure in $\widehat{\Sigma}^*$ implies $\widehat{\Sigma}^* \alpha = \left(\widehat{\Sigma}_{1 \cdot}^{\text{co}} \mathbf{a} \cdot \mathbf{1}'_{N_1}, \dots, \widehat{\Sigma}_{K \cdot}^{\text{co}} \mathbf{a} \cdot \mathbf{1}'_{N_K} \right)'$. Therefore, we have

$$\begin{aligned} \alpha' \widehat{\Sigma}^* \widehat{\Sigma}^* \alpha &= \sum_{k=1}^K N_k \left(\widehat{\Sigma}_{k \cdot}^{\text{co}} \mathbf{a} \right)^2 = N \sum_{k=1}^K r_k \left(\widehat{\Sigma}_{k \cdot}^{\text{co}} \mathbf{a} \right)^2 = N \mathbf{a}' \widehat{\Sigma}^{\text{co}} \mathbf{R} \widehat{\Sigma}^{\text{co}} \mathbf{a}, \\ \mathbf{1}'_N \widehat{\Sigma}^* \alpha &= \sum_{k=1}^K N_k \widehat{\Sigma}_{k \cdot}^{\text{co}} \mathbf{a} = N \sum_{k=1}^K r_k \widehat{\Sigma}_{k \cdot}^{\text{co}} \mathbf{a} = N \mathbf{r}' \widehat{\Sigma}^{\text{co}} \mathbf{a}. \end{aligned} \quad (\text{A.6})$$

Substituting these two equations to (A.5) yields

$$\left\| \widehat{\mathbf{A}}^* \alpha \right\|_2^2 = N \mathbf{a}' \widehat{\Sigma}^{\text{co}} \mathbf{R} \widehat{\Sigma}^{\text{co}} \mathbf{a} - N \left(\mathbf{r}' \widehat{\Sigma}^{\text{co}} \mathbf{a} \right)^2 = N \mathbf{a}' \widehat{\Sigma}^{\text{co}} (\mathbf{R} - \mathbf{r} \mathbf{r}') \widehat{\Sigma}^{\text{co}} \mathbf{a}.$$

On the other hand, noticing $\mathbf{R} \mathbf{1}_K = \mathbf{r}$ and $\mathbf{1}'_K \mathbf{R} \mathbf{1}_K = \mathbf{1}'_K \mathbf{r} = 1$, we have

$$\begin{aligned} N \left\| \widehat{\mathbf{A}}^{\text{co}} \mathbf{a} \right\|_2^2 &= N \mathbf{a}' \widehat{\Sigma}^{\text{co}} (\mathbf{I}_K - \mathbf{r} \mathbf{1}'_K) \mathbf{R} (\mathbf{I}_K - \mathbf{1}_K \mathbf{r}') \widehat{\Sigma}^{\text{co}} \mathbf{a} \\ &= N \mathbf{a}' \widehat{\Sigma}^{\text{co}} (\mathbf{R} - \mathbf{R} \mathbf{1}_K \mathbf{r}' - \mathbf{r} \mathbf{1}'_K \mathbf{R} + \mathbf{r} \mathbf{1}'_K \mathbf{R} \mathbf{1}_K \mathbf{r}') \widehat{\Sigma}^{\text{co}} \mathbf{a} \\ &= N \mathbf{a}' \widehat{\Sigma}^{\text{co}} (\mathbf{R} - \mathbf{r} \mathbf{r}') \widehat{\Sigma}^{\text{co}} \mathbf{a}. \end{aligned}$$

Therefore we obtain

$$\left\| \widehat{\mathbf{A}}^* \alpha \right\|_2^2 = N \left\| \widehat{\mathbf{A}}^{\text{co}} \mathbf{a} \right\|_2^2 = N \mathbf{a}' \widehat{\mathbf{A}}^{\text{co}'} \widehat{\mathbf{A}}^{\text{co}} \mathbf{a}. \quad (\text{A.7})$$

In the objective function (3.10), by (A.7) we can use $N \mathbf{a}' \widehat{\mathbf{A}}^{\text{co}'} \widehat{\mathbf{A}}^{\text{co}} \mathbf{a}$ to replace $\alpha' \widehat{\mathbf{A}}^* \widehat{\mathbf{A}}^* \alpha$, by (A.6) we can use $\mathbf{r}' \widehat{\Sigma}^{\text{co}} \mathbf{a}$ to replace $\frac{1}{N} \mathbf{1}'_N \widehat{\Sigma}^* \alpha$, and by Part (a) its solution must be of similar sign in that $\|\alpha^*\|_1 = \|\mathbf{a}^*\|_1$. Consequently, the problem in (3.10) is equivalent to that in (3.12).

Part (c). We first show $\widehat{\mathbf{A}}^{\text{co}}$ is of full column rank. The first K -rows of $\widehat{\mathbf{A}}^{\text{co}}$ is $\widehat{\mathbf{A}}^{\text{co}} = \mathbf{R}^{1/2} (\mathbf{I}_K - \mathbf{1}_K \mathbf{r}') \widehat{\Sigma}^{\text{co}}$. Notice that $\mathbf{I}_K - \mathbf{r} \mathbf{1}'_K$ is idempotent and $\mathbf{R}^{1/2}$ and $\widehat{\Sigma}^{\text{co}}$ are both of full rank, $\text{rank}(\widehat{\mathbf{A}}^{\text{co}}) = \text{rank}(\mathbf{I}_K - \mathbf{r} \mathbf{1}'_K) = \text{trace}(\mathbf{I}_K - \mathbf{r} \mathbf{1}'_K) = K - 1$. In other words, $\widehat{\mathbf{A}}^{\text{co}}$ is rank deficient and its null space is one-dimensional. The null space of $\widehat{\mathbf{A}}^{\text{co}}$ is

$$\ker(\widehat{\mathbf{A}}^{\text{co}}) = \left\{ c (\widehat{\Sigma}^{\text{co}})^{-1} \mathbf{1}_K : c \in \mathbb{R} \setminus \{0\} \right\},$$

as $\widehat{\mathbf{A}}^{\text{co}} (\widehat{\Sigma}^{\text{co}})^{-1} \mathbf{1}_K = \mathbf{R}^{1/2} (\mathbf{I}_K - \mathbf{1}_K \mathbf{r}') \mathbf{1}_K = \mathbf{0}_N$. Moreover, since $\mathbf{1}'_K (\widehat{\Sigma}^{\text{co}})^{-1} \mathbf{1}_K \neq 0$, $(\widehat{\Sigma}^{\text{co}})^{-1} \mathbf{1}_K$ is not in the null space of $\mathbf{1}'_K$. In other words, $\ker(\widehat{\mathbf{A}}^{\text{co}}) \cap \ker(\mathbf{1}'_K)$ is empty and we must have $\text{rank}(\widehat{\mathbf{A}}^{\text{co}}) = \text{rank}(\widehat{\mathbf{A}}^{\text{co}}) + \text{rank}(\mathbf{1}'_K) = (K - 1) + 1 = K$.

Setting $\tau = 0$ in (3.11), we have

$$\widehat{\mathbf{A}}^* \boldsymbol{\alpha}_0^* = \mathbf{w}_0^* - \frac{\mathbf{1}_N}{N}. \quad (\text{A.8})$$

Premultiplying both sides of the above equation by the $K \times N$ block diagonal matrix $\text{diag}(r_1^{-1/2} \mathbf{1}'_{N_1}, \dots, r_K^{-1/2} \mathbf{1}'_{N_K})$, we obtain

$$N \widehat{\mathbf{A}}^{\text{co}} \mathbf{a}_0^* = \mathbf{R}^{-1/2} (\mathbf{b}_0^* \circ \mathbf{r} - \mathbf{r}). \quad (\text{A.9})$$

As $\widehat{\mathbf{A}}^{\text{co}}$ is a submatrix of $\widetilde{\mathbf{A}}^{\text{co}}$, the above equation implies

$$N \widetilde{\mathbf{A}}^{\text{co}} \mathbf{a}_0^* = \begin{pmatrix} N \widehat{\mathbf{A}}^{\text{co}} \mathbf{a}_0^* \\ N \mathbf{1}'_K \mathbf{a}_0^* \end{pmatrix} = \begin{pmatrix} \mathbf{R}^{-1/2} (\mathbf{b}_0^* \circ \mathbf{r} - \mathbf{r}) \\ 0 \end{pmatrix},$$

where we use the restriction $N \mathbf{1}'_K \mathbf{a}_0^* = 0$. Since $\widetilde{\mathbf{A}}^{\text{co}}$ is of full column rank, we explicitly solve $\mathbf{a}_0^* = \left(\widetilde{\mathbf{A}}^{\text{co}'} \widetilde{\mathbf{A}}^{\text{co}} \right)^{-1} \widetilde{\mathbf{A}}^{\text{co}'} \widetilde{\mathbf{b}}^{\text{co}} / N$. ■

Proof of Theorem 1. Part (a). Substituting $(\mathbf{w}_0^*, \gamma_0^*)$ into the constraint in (2.5), we obtain

$$\begin{aligned} \left\| \widehat{\boldsymbol{\Sigma}} \mathbf{w}_0^* + \gamma_0^* \mathbf{1}_N \right\|_\infty &= \left\| \widehat{\boldsymbol{\Sigma}}^* \mathbf{w}_0^* + \widehat{\boldsymbol{\Sigma}}^e \mathbf{w}_0^* + \gamma_0^* \mathbf{1}_N \right\|_\infty = \left\| \widehat{\boldsymbol{\Sigma}}^e \mathbf{w}_0^* \right\|_\infty \\ &= \max_i \left\| \widehat{\boldsymbol{\Sigma}}_i^e \mathbf{w}_0^* \right\|_\infty \leq \left\| \widehat{\boldsymbol{\Sigma}}^e \right\|_{c2} \left\| \mathbf{w}_0^* \right\|_2 \leq \phi_e \left\| \mathbf{b}_0^* \right\|_\infty / \sqrt{N}. \end{aligned} \quad (\text{A.10})$$

where the second equality follows by the KKT condition $\widehat{\boldsymbol{\Sigma}}^* \mathbf{w}_0^* + \gamma_0^* \mathbf{1}_N = 0$. The presumption $\phi_e \left\| \mathbf{b}_0^* \right\|_\infty / \sqrt{N} < \tau$ in the statement makes sure that $\left\| \widehat{\boldsymbol{\Sigma}} \mathbf{w}_0^* + \gamma_0^* \mathbf{1}_N \right\|_\infty < \tau$ holds with strict inequality. This strict inequality means that $(\mathbf{w}_0^*, \gamma_0^*)$ lies in the interior of the feasible set of (2.5). Because $\widehat{\mathbf{w}}$ is the minimizer of the problem in (2.5), its ℓ_2 -norm is no greater than any other feasible solution. Thus $\left\| \widehat{\mathbf{w}} \right\|_2 \leq \left\| \mathbf{w}_0^* \right\|_2$ and furthermore $\left\| \mathbf{w}_0^* \right\|_2$ is bounded by $\left\| \mathbf{b}_0^* \right\|_\infty / \sqrt{N}$ by (3.7).

Part (b). The Lagrangian of (2.6) can be written as

$$\mathcal{L}(\boldsymbol{\alpha}, \gamma) = \frac{1}{2} \boldsymbol{\alpha}' \widehat{\mathbf{A}}' \widehat{\mathbf{A}} \boldsymbol{\alpha} + \frac{1}{N} \mathbf{1}'_N \widehat{\boldsymbol{\Sigma}} \boldsymbol{\alpha} + \tau \left\| \boldsymbol{\alpha} \right\|_1 + \gamma \mathbf{1}'_N \boldsymbol{\alpha} - \frac{1}{2N},$$

where γ is the Lagrangian multiplier for the constraint $\mathbf{1}'_N \boldsymbol{\alpha} = 0$. Consider the subgradient of $\mathcal{L}(\widehat{\boldsymbol{\alpha}}, \widehat{\gamma})$ with respect to α_i for any $i \in [N]$, where $(\widehat{\boldsymbol{\alpha}}, \widehat{\gamma})$ is the optimizer. Noting that $\widehat{\mathbf{A}}' \widehat{\mathbf{A}} = \widehat{\mathbf{A}}' \widehat{\boldsymbol{\Sigma}}$ due to the fact that $\mathbf{I}_N - N^{-1} \mathbf{1}_N \mathbf{1}'_N$ is a projection matrix, the KKT conditions imply that

$$\left| \widehat{\boldsymbol{\alpha}}' \widehat{\mathbf{A}}' \widehat{\mathbf{A}}_{\cdot i} + \frac{1}{N} \mathbf{1}'_N \widehat{\boldsymbol{\Sigma}}_{\cdot i} + \widehat{\gamma} \right| = \left| (\widehat{\mathbf{A}} \widehat{\boldsymbol{\alpha}} + \frac{1}{N} \mathbf{1}_N)' \widehat{\boldsymbol{\Sigma}}_{\cdot i} + \widehat{\gamma} \right| = \left| \widehat{\mathbf{w}}' \widehat{\boldsymbol{\Sigma}}_{\cdot i} + \widehat{\gamma} \right| \leq \tau \text{ for all } i \in [N]$$

and furthermore

$$\widehat{\mathbf{w}}' \widehat{\boldsymbol{\Sigma}}_{\cdot i} + \widehat{\gamma} = \tau \text{sign}(\widehat{\alpha}_i) \text{ for all } \widehat{\alpha}_i \neq 0. \quad (\text{A.11})$$

Suppose $\widehat{\boldsymbol{\alpha}} \notin \mathbb{S}^{\text{all}}$. Without loss of generality, let $\widehat{\alpha}_i > 0$ and $\widehat{\alpha}_j < 0$ for some $i, j \in \mathcal{G}_k$, $i \neq j$. (A.11) indicates $\widehat{\mathbf{w}}' \widehat{\boldsymbol{\Sigma}}_{\cdot i} + \widehat{\gamma} = \tau$ and $\widehat{\mathbf{w}}' \widehat{\boldsymbol{\Sigma}}_{\cdot j} + \widehat{\gamma} = -\tau$. Subtracting these two equations on both sides

yields

$$\begin{aligned} 2\tau &= \left| \widehat{\mathbf{w}}'(\widehat{\boldsymbol{\Sigma}}_{.i} - \widehat{\boldsymbol{\Sigma}}_{.j}) \right| = \left| \widehat{\mathbf{w}}'[(\widehat{\boldsymbol{\Sigma}}_{.i}^e + \widehat{\boldsymbol{\Sigma}}_{.i}^*) - (\widehat{\boldsymbol{\Sigma}}_{.j}^e + \widehat{\boldsymbol{\Sigma}}_{.j}^*)] \right| = \left| \widehat{\mathbf{w}}'(\widehat{\boldsymbol{\Sigma}}_{.i}^e - \widehat{\boldsymbol{\Sigma}}_{.j}^e) \right| \\ &\leq \|\widehat{\boldsymbol{\Sigma}}_{.i}^e - \widehat{\boldsymbol{\Sigma}}_{.j}^e\|_2 \|\widehat{\mathbf{w}}\|_2 \leq 2\|\widehat{\boldsymbol{\Sigma}}^e\|_{c2} \|\widehat{\mathbf{w}}\|_2 \leq 2\phi_e \|\mathbf{b}_0^*\|_\infty / \sqrt{N}, \end{aligned} \quad (\text{A.12})$$

where the third equality holds as $\widehat{\boldsymbol{\Sigma}}_{.i}^* = \widehat{\boldsymbol{\Sigma}}_{.j}^*$ for i and j in the same group k , and the last inequality by Part (a) which bounds $\|\widehat{\mathbf{w}}\|_2$. The above inequality (A.12) violates the presumption $\tau > \phi_e \|\mathbf{b}_0^*\|_\infty / \sqrt{N}$. We thus conclude $\widehat{\boldsymbol{\alpha}} \in \mathbb{S}^{\text{all}}$. ■

A.3 Lemmas Prepared for the Proof of Theorem 2

Lemma 4(a) below provides a *compatibility inequality* that links $\|\boldsymbol{\delta}\|_1$ and $\|\widehat{\mathbf{A}}\boldsymbol{\delta}\|_2^2$ for any $\boldsymbol{\delta} \in \widetilde{\mathbb{S}}^{\text{all}}$. Instead of imposing a restricted eigenvalue condition as an assumption, we establish our comparability inequality under a primitive condition about ϕ_e . Recall that $\phi_A := \phi_{\min}(\widehat{\mathbf{A}}^{\text{co}} \widehat{\mathbf{A}}^{\text{co}}) \wedge 1$ and $\underline{r} := \min_{k \in [K]} r_k$.

Lemma 4 (a) If $\phi_e \leq \frac{1}{2} \sqrt{N\phi_A/K}$, we have $\|\boldsymbol{\delta}\|_1 \leq 2\sqrt{K/(N\phi_A)} \|\widehat{\mathbf{A}}\boldsymbol{\delta}\|_2$ for any $\boldsymbol{\delta} \in \widetilde{\mathbb{S}}^{\text{all}}$.

$$(b) \phi_A^{-1} \leq 2\underline{r}^{-1} \phi_{\min}^{-2}(\widehat{\boldsymbol{\Sigma}}^{\text{co}}) + K^{-1} \phi_{\max}(\widehat{\boldsymbol{\Sigma}}^{\text{co}}) / \phi_{\min}(\widehat{\boldsymbol{\Sigma}}^{\text{co}}).$$

$$(c) \|\mathbf{a}_0^*\|_1 \leq N^{-1} \sqrt{K/\phi_A} (\|\mathbf{b}_0^*\|_\infty + 1).$$

Remark A.1. The constants 1/2 and 2 in Lemma 4(a) are not important in the asymptotic analysis. It means if the magnitude of the idiosyncratic shock, represented by ϕ_e , is controlled by the order $\sqrt{N\phi_A/K}$, then the ℓ_1 -norm of $\boldsymbol{\delta}$ can be controlled by the ℓ_2 -norm of $\|\widehat{\mathbf{A}}\boldsymbol{\delta}\|_2$ multiplied by a factor involving K/ϕ_A , which is the ratio between the number of groups K and the square of the minimal non-trivial singular value of the augmented weighted demeaned core $\widehat{\mathbf{A}}^{\text{co}}$. In the proof of Lemma 4, we introduce an original self-defined semi-norm (A.14) to take advantage of the group pattern.

Remark A.2. Another necessary condition for Lasso to achieve reasonable performance is that the ℓ_1 -norm of the true coefficients cannot be too large. In our context, since Lemma 3 implies $\|\boldsymbol{\alpha}^*\|_1 = \|\mathbf{a}^*\|_1$, Part (c) sets an upper bound for the ℓ_1 -norm of the true coefficient value for (3.10).

Proof of Lemma 4. Part (a). For a generic vector $\boldsymbol{\delta} \in \widetilde{\mathbb{S}}^{\text{all}}$, we have

$$\|\widehat{\mathbf{A}}\boldsymbol{\delta}\|_2 \geq \|\widehat{\mathbf{A}}^* \boldsymbol{\delta}\|_2 - \|(\mathbf{I}_N - N^{-1} \mathbf{1}_N \mathbf{1}'_N) \widehat{\boldsymbol{\Sigma}}^e \boldsymbol{\delta}\|_2 \geq \|\widehat{\mathbf{A}}^* \boldsymbol{\delta}\|_2 - \|\widehat{\boldsymbol{\Sigma}}^e \boldsymbol{\delta}\|_2, \quad (\text{A.13})$$

where the first inequality holds by the triangle inequality, and the second follows because $\mathbf{I}_N - N^{-1} \mathbf{1}_N \mathbf{1}'_N$ is a projection matrix. We will bound the two terms on the right hand side.

To take advantage of the group structure to handle collinearity, we introduce a novel groupwise semi-norm and establish a corresponding version of compatibility condition. Let $d_k = \sum_{i \in \mathcal{G}_k} \delta_i$ for $k \in [K]$ and $\mathbf{d} = (d_1, \dots, d_K)'$. Define a groupwise ℓ_2 semi-norm $\|\cdot\|_{2\mathcal{G}} : \mathbb{R}^N \mapsto \mathbb{R}^+$ as

$$\|\boldsymbol{\delta}\|_{2\mathcal{G}} = \|\mathbf{d}\|_2. \quad (\text{A.14})$$

The definition of the semi-norm depends on the true group membership, which is infeasible in reality. We introduce this semi-norm only for theoretical development. In the estimation we do not need to

know the true group membership. This semi-norm $\|\boldsymbol{\delta}\|_{2\mathcal{G}}$ allows $\|\boldsymbol{\delta}\|_{2\mathcal{G}} = 0$ even if $\boldsymbol{\delta} \neq \mathbf{0}_N$, while it remains homogeneous, sub-additive, and non-negative — all other desirable properties of a norm. Moreover, if $\boldsymbol{\delta} \in \mathbb{S}^{\text{all}}$ it is obvious

$$\|\boldsymbol{\delta}\|_1 = \sum_{k \in [K]} \left| \sum_{i \in \mathcal{G}_k} \delta_i \right| = \sum_{k \in [K]} |d_k| \leq \sqrt{K} \|\boldsymbol{\delta}\|_{2\mathcal{G}} \quad (\text{A.15})$$

by either the Cauchy-Schwarz or Jensen's inequality.

For any $\boldsymbol{\delta} \in \tilde{\mathbb{S}}^{\text{all}}$, we have

$$\begin{aligned} \|\widehat{\mathbf{A}}^* \boldsymbol{\delta}\|_2 &= \sqrt{N} \|\widehat{\mathbf{A}}^{\text{co}} \mathbf{d}\|_2 = \sqrt{N} \left\| \begin{pmatrix} \widehat{\mathbf{A}}^{\text{co}} \mathbf{d} \\ 0 \end{pmatrix} \right\|_2 = \sqrt{N} \|\tilde{\mathbf{A}}^{\text{co}} \mathbf{d}\|_2 \\ &\geq \sqrt{N\phi_A} \|\mathbf{d}\|_2 = \sqrt{N\phi_A} \|\boldsymbol{\delta}\|_{2\mathcal{G}} \geq \sqrt{\frac{N\phi_A}{K}} \|\boldsymbol{\delta}\|_1, \end{aligned} \quad (\text{A.16})$$

where the first equality follows by (A.7), the third equality by $\mathbf{1}'_K \mathbf{d} = \mathbf{1}'_N \boldsymbol{\delta} = 0$, and the last inequality by (A.15). We have found a lower bound for the first term on the right hand side of (A.13). For the second term on the right hand side of (A.13), we have

$$\|\widehat{\boldsymbol{\Sigma}}^e \boldsymbol{\delta}\|_2 \leq \|\widehat{\boldsymbol{\Sigma}}\|_{c2} \|\boldsymbol{\delta}\|_1 \leq \phi_e \|\boldsymbol{\delta}\|_1 \quad (\text{A.17})$$

by (A.57) and the definition of ϕ_e .

Under the presumption $\phi_e \leq \frac{1}{2} \sqrt{N\phi_A/K}$, (A.13) and (A.16)-(A.17) together imply

$$\|\widehat{\mathbf{A}} \boldsymbol{\delta}\|_2 \geq (\sqrt{N\phi_A/K} - \phi_e) \|\boldsymbol{\delta}\|_1 \geq \frac{1}{2} \sqrt{N\phi_A/K} \|\boldsymbol{\delta}\|_1.$$

Then the result in (a) follows.

Part (b). Notice

$$\begin{aligned} \tilde{\mathbf{A}}^{\text{co}'} \tilde{\mathbf{A}}^{\text{co}} &= \widehat{\mathbf{A}}^{\text{co}'} \widehat{\mathbf{A}}^{\text{co}} + \mathbf{1}_K \mathbf{1}'_K = \widehat{\boldsymbol{\Sigma}}^{\text{co}} (\mathbf{I}_K - \mathbf{r} \cdot \mathbf{1}'_K) \mathbf{R} (\mathbf{I}_K - \mathbf{1}_K \cdot \mathbf{r}') \widehat{\boldsymbol{\Sigma}}^{\text{co}} + \mathbf{1}_K \mathbf{1}'_K \\ &= \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{R} \widehat{\boldsymbol{\Sigma}}^{\text{co}} + \mathbf{1}_K \mathbf{1}'_K - \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{r} \mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}}. \end{aligned}$$

By the Sherman-Morrison formula in (A.59),

$$(\tilde{\mathbf{A}}^{\text{co}'} \tilde{\mathbf{A}}^{\text{co}})^{-1} = \mathbf{A}_1^{-1} + \frac{\mathbf{A}_1^{-1} \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{r} \mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{A}_1^{-1}}{1 - \mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{A}_1^{-1} \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{r}}, \quad (\text{A.18})$$

where $\mathbf{A}_1 = \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{R} \widehat{\boldsymbol{\Sigma}}^{\text{co}} + \mathbf{1}_K \mathbf{1}'_K$ and moreover

$$\mathbf{A}_1^{-1} = \mathbf{A}_2^{-1} - \frac{\mathbf{A}_2^{-1} \mathbf{1}_K \mathbf{1}'_K \mathbf{A}_2^{-1}}{1 + \mathbf{1}'_K \mathbf{A}_2^{-1} \mathbf{1}_K} \quad (\text{A.19})$$

by (A.58), where $\mathbf{A}_2 = \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{R} \widehat{\boldsymbol{\Sigma}}^{\text{co}}$. Obviously

$$\phi_{\max}(\mathbf{A}_1^{-1}) \leq [\phi_{\min}(\mathbf{A}_2)]^{-1} \leq \underline{r}^{-1} \phi_{\min}^{-2}(\widehat{\boldsymbol{\Sigma}}^{\text{co}}) \quad (\text{A.20})$$

and

$$\mathbf{1}'_K \mathbf{A}_2^{-1} \mathbf{1}_K \leq \phi_{\min}^{-1}(\mathbf{R}) \phi_{\min}^{-1}(\widehat{\boldsymbol{\Sigma}}^{\text{co}}) \mathbf{1}'_K (\widehat{\boldsymbol{\Sigma}}^{\text{co}})^{-1} \mathbf{1}_K. \quad (\text{A.21})$$

The denominator of the second term on the right hand side of (A.18) is

$$\begin{aligned} 1 - \mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{A}_1^{-1} \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{r} &= 1 - \mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \left(\mathbf{A}_2^{-1} - \frac{\mathbf{A}_2^{-1} \mathbf{1}_K \mathbf{1}'_K \mathbf{A}_2^{-1}}{1 + \mathbf{1}'_K \mathbf{A}_2^{-1} \mathbf{1}_K} \right) \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{r} \\ &= \mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \frac{\mathbf{A}_2^{-1} \mathbf{1}_K \mathbf{1}'_K \mathbf{A}_2^{-1}}{1 + \mathbf{1}'_K \mathbf{A}_2^{-1} \mathbf{1}_K} \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{r} = \frac{[\mathbf{1}'_K (\widehat{\boldsymbol{\Sigma}}^{\text{co}})^{-1} \mathbf{1}_K]^2}{1 + \mathbf{1}'_K \mathbf{A}_2^{-1} \mathbf{1}_K} > 0, \end{aligned} \quad (\text{A.22})$$

where the second equality follows by $\mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{A}_2^{-1} \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{r} = \mathbf{r}' \mathbf{R}^{-1} \mathbf{r} = \mathbf{1}'_K \mathbf{r} = 1$, and the third equality by $\mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{A}_2^{-1} \mathbf{1}_K = \mathbf{1}'_K (\widehat{\boldsymbol{\Sigma}}^{\text{co}})^{-1} \mathbf{1}_K$. The numerator of the second term on the right hand side of (A.18) has rank 1, and thus

$$\begin{aligned} \phi_{\max} \left(\mathbf{A}_1^{-1} \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{r} \mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{A}_1^{-1} \right) &= \text{trace} \left(\mathbf{A}_1^{-1} \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{r} \mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{A}_1^{-1} \right) = \mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{A}_1^{-2} \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{r} \\ &\leq \mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{A}_2^{-2} \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{r} = \mathbf{r}' \mathbf{R}^{-1} (\widehat{\boldsymbol{\Sigma}}^{\text{co}})^{-2} \mathbf{R}^{-1} \mathbf{r} = \mathbf{1}'_K (\widehat{\boldsymbol{\Sigma}}^{\text{co}})^{-2} \mathbf{1}_K. \end{aligned} \quad (\text{A.23})$$

Combine (A.22) and (A.23),

$$\begin{aligned} \frac{\phi_{\max} \left(\mathbf{A}_1^{-1} \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{r} \mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{A}_1^{-1} \right)}{1 - \mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{A}_1^{-1} \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{r}} &\leq \mathbf{1}'_K (\widehat{\boldsymbol{\Sigma}}^{\text{co}})^{-2} \mathbf{1}_K \times \frac{1 + \mathbf{1}'_K \mathbf{A}_2^{-1} \mathbf{1}_K}{[\mathbf{1}'_K (\widehat{\boldsymbol{\Sigma}}^{\text{co}})^{-1} \mathbf{1}_K]^2} \\ &\leq \phi_{\min}^{-1}(\widehat{\boldsymbol{\Sigma}}^{\text{co}}) \left[\left(\mathbf{1}'_K (\widehat{\boldsymbol{\Sigma}}^{\text{co}})^{-1} \mathbf{1}_K \right)^{-1} + \phi_{\min}^{-1}(\mathbf{R}) \phi_{\min}^{-1}(\widehat{\boldsymbol{\Sigma}}^{\text{co}}) \right] \\ &\leq \frac{\phi_{\min}^{-1}(\widehat{\boldsymbol{\Sigma}}^{\text{co}})}{K \phi_{\max}^{-1}(\widehat{\boldsymbol{\Sigma}}^{\text{co}})} + \underline{r}^{-1} \phi_{\min}^{-2}(\widehat{\boldsymbol{\Sigma}}^{\text{co}}), \end{aligned} \quad (\text{A.24})$$

where the second inequality holds by (A.21).

Applying the spectral norm to (A.18) yields

$$\begin{aligned} \phi_A^{-1} &= \phi_{\max} \left((\tilde{\mathbf{A}}^{\text{co}'} \tilde{\mathbf{A}}^{\text{co}})^{-1} \right) \leq \phi_{\max} \left(\mathbf{A}_1^{-1} \right) + \frac{\phi_{\max} \left(\mathbf{A}_1^{-1} \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{r} \mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{A}_1^{-1} \right)}{1 - \mathbf{r}' \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{A}_1^{-1} \widehat{\boldsymbol{\Sigma}}^{\text{co}} \mathbf{r}} \\ &\leq \underline{r}^{-1} \phi_{\min}^{-2}(\widehat{\boldsymbol{\Sigma}}^{\text{co}}) + \frac{\phi_{\max}(\widehat{\boldsymbol{\Sigma}}^{\text{co}})}{K \phi_{\min}(\widehat{\boldsymbol{\Sigma}}^{\text{co}})} + \underline{r}^{-1} \phi_{\min}^{-2}(\widehat{\boldsymbol{\Sigma}}^{\text{co}}) \end{aligned}$$

by (A.20) and (A.24).

Part (c). Given the expression of \mathbf{a}_0^* in Lemma 3, its ℓ_2 -norm is bounded by

$$\begin{aligned} \|\mathbf{a}_0^*\|_2 &\leq \|(\tilde{\mathbf{A}}^{\text{co}'} \tilde{\mathbf{A}}^{\text{co}})^{-1} \tilde{\mathbf{A}}^{\text{co}'}\|_{\text{sp}} \|(\mathbf{b}_0^* \circ \mathbf{r} - \mathbf{r})' \mathbf{R}^{-1/2}\|_2 / N \\ &\leq \frac{1}{N \sqrt{\phi_A}} \left(\|\mathbf{b}_0^* \circ \mathbf{r}^{1/2}\|_2 + \|\mathbf{r}^{1/2}\|_2 \right) \leq \frac{1}{N \sqrt{\phi_A}} (\|\mathbf{b}_0^*\|_{\infty} + 1), \end{aligned} \quad (\text{A.25})$$

where $\mathbf{r}^{1/2} = (r_1^{1/2}, \dots, r_k^{1/2})'$. In addition, the Cauchy-Schwarz inequality entails

$$\|\mathbf{a}_0^*\|_1 \leq \sqrt{K} \|\mathbf{a}_0^*\|_2 = N^{-1} \sqrt{K/\phi_A} (\|\mathbf{b}_0^*\|_\infty + 1) \quad (\text{A.26})$$

as stated in the lemma. ■

Lemma 5 collects the implications of Assumptions 1 and 2 for some building blocks of our asymptotic theory. Lemma 5(a) provides the magnitude ϕ_e , ϕ_A^{-1} and $\|\mathbf{b}_0^*\|_\infty$, and part (b) shows that the key condition in the numerical properties in Theorem 1 is satisfied.

Lemma 5 *Under Assumptions 1 and 2,*

- (a) $\phi_e = O_p(\sqrt{N}\phi_{NT})$, $\phi_A^{-1} = O_p(K)$, and $\|\mathbf{b}_0^*\|_\infty = O_p(\sqrt{K})$;
(b) The event $\{\phi_e \|\mathbf{b}_0^*\|_\infty / \sqrt{N} < \tau\}$ occurs w.p.a.1.

Proof of Lemma 5. Part (a). By the definition of ϕ_e and the triangle inequality,

$$\begin{aligned} \phi_e &= \|\widehat{\boldsymbol{\Sigma}}^e\|_{c2} \leq \|\boldsymbol{\Sigma}_0^e\|_{c2} + \|\boldsymbol{\Delta}^e\|_{c2} \\ &\leq C_{e0} \phi_{\max}(\boldsymbol{\Sigma}_0^e) + \sqrt{N} \|\boldsymbol{\Delta}^e\|_\infty = O_p(\sqrt{N}\phi_{NT}), \end{aligned} \quad (\text{A.27})$$

where the second inequality and the last equality follow by Assumption 1(a).

Noting that $\widehat{\boldsymbol{\Sigma}}^{\text{co}} = \boldsymbol{\Sigma}_0^{\text{co}} + \boldsymbol{\Delta}^{\text{co}}$, and then w.p.a.1

$$\begin{aligned} \phi_{\min}(\widehat{\boldsymbol{\Sigma}}^{\text{co}}) &\geq \phi_{\min}(\boldsymbol{\Sigma}_0^{\text{co}}) - \|\boldsymbol{\Delta}^{\text{co}}\|_{\text{sp}} \geq \phi_{\min}(\boldsymbol{\Sigma}_0^{\text{co}}) - K \|\boldsymbol{\Delta}^{\text{co}}\|_\infty \\ &\geq \underline{c} - O_p(K(T/\log N)^{-1/2}) \geq \underline{c}/2 \end{aligned} \quad (\text{A.28})$$

where the first inequality follows by Weyl inequality, the second inequality by the Gershgorin circle theorem, the third inequality by Assumption 1(b), and the last inequality holds when the sample size is sufficiently large. Similarly,

$$\phi_{\max}(\widehat{\boldsymbol{\Sigma}}^{\text{co}}) \leq \phi_{\max}(\boldsymbol{\Sigma}_0^{\text{co}}) + \|\boldsymbol{\Delta}^{\text{co}}\|_{\text{sp}} \leq 2\bar{c} \quad (\text{A.29})$$

w.p.a.1. Suppose (A.28) and (A.29) occur. Given Assumption 2(b) about the rate of \underline{r} , Lemma 4(b) implies

$$\phi_A^{-1} \leq 8\underline{r}^{-1}\underline{c}^{-2} + \frac{4\bar{c}}{K\underline{c}} = O_p(K)$$

and (3.9) implies

$$\begin{aligned} \|\mathbf{b}_0^*\|_\infty &\leq \left\| (\widehat{\boldsymbol{\Sigma}}^{\text{co}})^{-1} \mathbf{1}_K \right\|_\infty / \left[\underline{r} \cdot \mathbf{1}'_K (\widehat{\boldsymbol{\Sigma}}^{\text{co}})^{-1} \mathbf{1}_K \right] \leq \left(\mathbf{1}'_K (\widehat{\boldsymbol{\Sigma}}^{\text{co}})^{-2} \mathbf{1}_K \right)^{1/2} / \left[\underline{r} \cdot \mathbf{1}'_K (\widehat{\boldsymbol{\Sigma}}^{\text{co}})^{-1} \mathbf{1}_K \right] \\ &\leq \underline{r}^{-1} \phi_{\min}^{-1}(\widehat{\boldsymbol{\Sigma}}^{\text{co}}) \left(\mathbf{1}'_K (\widehat{\boldsymbol{\Sigma}}^{\text{co}})^{-1} \mathbf{1}_K \right)^{-1/2} \leq \underline{r}^{-1} K^{-1/2} \phi_{\max}^{1/2}(\widehat{\boldsymbol{\Sigma}}^{\text{co}}) / \phi_{\min}(\widehat{\boldsymbol{\Sigma}}^{\text{co}}) \\ &\leq \underline{r}^{-1} K^{-1/2} \cdot O_p(\bar{c}^{1/2}/\underline{c}) = O_p(\sqrt{K}). \end{aligned} \quad (\text{A.30})$$

Part (b). The results of Part (a) gives $\phi_e \|\mathbf{b}_0^*\|_\infty / \sqrt{N} = O_p(K^{1/2}\phi_{NT})$. Since $K^{1/2}\phi_{NT}/\tau \rightarrow 0$ in Assumption 2(a), $\tau > \phi_e \|\mathbf{b}_0^*\|_\infty / \sqrt{N}$ when (N, T) are sufficiently large and thus the event $\{\phi_e \|\mathbf{b}_0^*\|_\infty / \sqrt{N} < \tau\}$ occurs w.p.a.1. ■

A.4 Proofs of the Results in Section 4

Proof of Theorem 2. When the sample size is sufficiently large we have $\{\phi_e \|\mathbf{b}_0^*\|_\infty / \sqrt{N} < \tau\}$ w.p.a.1 by Lemma 5(b). We can then construct the desirable $\hat{\boldsymbol{\alpha}}^*$ according to (4.1). Since $\hat{\boldsymbol{\alpha}}$ is the solution to (A.4),

$$\frac{1}{2} \hat{\boldsymbol{\alpha}}' \hat{\mathbf{A}}' \hat{\mathbf{A}} \hat{\boldsymbol{\alpha}} + \frac{1}{N} \mathbf{1}'_N \hat{\boldsymbol{\Sigma}} \hat{\boldsymbol{\alpha}} + \tau \|\boldsymbol{\alpha}\|_1 \leq \frac{1}{2} \hat{\boldsymbol{\alpha}}'^* \hat{\mathbf{A}}' \hat{\mathbf{A}} \hat{\boldsymbol{\alpha}}^* + \frac{1}{N} \mathbf{1}'_N \hat{\boldsymbol{\Sigma}} \hat{\boldsymbol{\alpha}}^* + \tau \|\hat{\boldsymbol{\alpha}}^*\|_1.$$

Define $\boldsymbol{\psi} = \hat{\boldsymbol{\alpha}} - \hat{\boldsymbol{\alpha}}^*$. Rearranging the above inequality yields

$$\boldsymbol{\psi}' \hat{\mathbf{A}}' \hat{\mathbf{A}} \boldsymbol{\psi} + 2\tau \|\hat{\boldsymbol{\alpha}}\|_1 \leq -2\boldsymbol{\psi}' \hat{\boldsymbol{\Sigma}} (\hat{\mathbf{A}} \hat{\boldsymbol{\alpha}}^* + \frac{\mathbf{1}_N}{N}) + 2\tau \|\hat{\boldsymbol{\alpha}}^*\|_1. \quad (\text{A.31})$$

Notice that

$$\begin{aligned} & \boldsymbol{\psi}' \hat{\boldsymbol{\Sigma}} \left(\hat{\mathbf{A}} \hat{\boldsymbol{\alpha}}^* + \frac{\mathbf{1}_N}{N} \right) \\ &= \boldsymbol{\psi}' \hat{\boldsymbol{\Sigma}} (\hat{\mathbf{A}}^* \hat{\boldsymbol{\alpha}}^* + \frac{\mathbf{1}_N}{N}) + \boldsymbol{\psi}' \hat{\boldsymbol{\Sigma}} (\hat{\mathbf{A}} - \hat{\mathbf{A}}^*) \hat{\boldsymbol{\alpha}}^* = \boldsymbol{\psi}' \hat{\boldsymbol{\Sigma}} \mathbf{w}_0^* + \boldsymbol{\psi}' \hat{\boldsymbol{\Sigma}}' (\mathbf{I}_N - N^{-1} \mathbf{1}_N \mathbf{1}_N') \hat{\boldsymbol{\Sigma}}^e \hat{\boldsymbol{\alpha}}^* \\ &= (\boldsymbol{\psi}' \hat{\boldsymbol{\Sigma}}^* \mathbf{w}_0^* + \boldsymbol{\psi}' \hat{\boldsymbol{\Sigma}}^e \mathbf{w}_0^*) + \boldsymbol{\psi}' \hat{\mathbf{A}}' \hat{\boldsymbol{\Sigma}}^e \hat{\boldsymbol{\alpha}}^* = (-\gamma_0^* \boldsymbol{\psi}' \mathbf{1}_N + \boldsymbol{\psi}' \hat{\boldsymbol{\Sigma}}^e \mathbf{w}_0^*) + \boldsymbol{\psi}' \hat{\mathbf{A}}' \hat{\boldsymbol{\Sigma}}^e \hat{\boldsymbol{\alpha}}^* \\ &= \boldsymbol{\psi}' \hat{\boldsymbol{\Sigma}}^e \mathbf{w}_0^* + \boldsymbol{\psi}' \hat{\mathbf{A}}' \hat{\boldsymbol{\Sigma}}^e \hat{\boldsymbol{\alpha}}^*, \end{aligned} \quad (\text{A.32})$$

where the fourth equality follows by the fact that $\hat{\boldsymbol{\Sigma}}^* \mathbf{w}_0^* = -\gamma_0^* \mathbf{1}_N$ implied by the KKT conditions in (2.2) with $\tau = 0$, and the last equality by $\boldsymbol{\psi}' \mathbf{1}_N = (\hat{\boldsymbol{\alpha}} - \hat{\boldsymbol{\alpha}}^*)' \mathbf{1}_N = 0$ as the dual problems entail both $\hat{\boldsymbol{\alpha}}$ and $\hat{\boldsymbol{\alpha}}^*$ sum up to 0. Plugging (A.32) into (A.31) to bound the right hand side of (A.31), we have

$$\boldsymbol{\psi}' \hat{\mathbf{A}}' \hat{\mathbf{A}} \boldsymbol{\psi} + 2\tau \|\hat{\boldsymbol{\alpha}}\|_1 \leq 2 \left| \boldsymbol{\psi}' \hat{\boldsymbol{\Sigma}}^e \mathbf{w}_0^* + \boldsymbol{\psi}' \hat{\mathbf{A}}' \hat{\boldsymbol{\Sigma}}^e \hat{\boldsymbol{\alpha}}^* \right| + 2\tau \|\hat{\boldsymbol{\alpha}}^*\|_1 \leq 2 \|\boldsymbol{\psi}\|_1 (\zeta_1 + \zeta_2) + 2\tau \|\hat{\boldsymbol{\alpha}}^*\|_1, \quad (\text{A.33})$$

where $\zeta_1 = \|\hat{\boldsymbol{\Sigma}}^e \mathbf{w}_0^*\|_\infty$ and $\zeta_2 = \|\hat{\mathbf{A}}' \hat{\boldsymbol{\Sigma}}^e \hat{\boldsymbol{\alpha}}^*\|_\infty$.

Now we bound ζ_1 and ζ_2 in turn. By (A.10), (A.27) and Lemma 5(a), we have

$$\zeta_1 \leq \|\hat{\boldsymbol{\Sigma}}^e\|_{c2} \|\mathbf{w}_0^*\|_2 \leq \phi_e \|\mathbf{b}_0^*\|_\infty / \sqrt{N} = O_p(K^{1/2} \phi_{NT}).$$

For ζ_2 , we have by (A.56),

$$\zeta_2 \leq \|\hat{\mathbf{A}}\|_{c2} \|\hat{\boldsymbol{\Sigma}}^e\|_{c2} \|\hat{\boldsymbol{\alpha}}^*\|_1 = \|\hat{\mathbf{A}}\|_{c2} \cdot \phi_e \|\hat{\boldsymbol{\alpha}}^*\|_1. \quad (\text{A.34})$$

Noting that $\|\mathbf{I}_N - N^{-1} \mathbf{1}_N \mathbf{1}_N'\|_{\text{sp}} = 1$ and by the triangle inequality, we have

$$\begin{aligned} \|\hat{\mathbf{A}}\|_{c2} &\leq \|\hat{\boldsymbol{\Sigma}}\|_{c2} \leq \|\hat{\boldsymbol{\Sigma}}^*\|_{c2} + \|\hat{\boldsymbol{\Sigma}}^e\|_{c2} \\ &\leq \sqrt{N} (\phi_{\max}(\boldsymbol{\Sigma}_0^{\text{co}}) + \|\boldsymbol{\Delta}^{\text{co}}\|_\infty) = O_p(\sqrt{N}), \end{aligned}$$

where the third inequality follows from Assumption 1(b). Noting that $\|\hat{\boldsymbol{\alpha}}^*\|_1 = \|\mathbf{a}_0^*\|_1 \leq (\|\mathbf{b}_0^*\|_\infty + 1)$

$\times \sqrt{K/\phi_A}/N$ by Lemma 4(c), we continue (A.34) to obtain

$$\zeta_2 = O_p(\sqrt{N})O_p(\sqrt{N}\phi_{NT})N^{-1}\sqrt{K/\phi_A}(\|\mathbf{b}_0^*\|_\infty + 1) = O_p(K^{1/2}\phi_{NT}\sqrt{K/\phi_A})$$

as $\|\mathbf{b}_0^*\|_\infty = O_p(K^{1/2})$ by Lemma 5(a). We thus obtain $\zeta_1 + \zeta_2 = O_p(K^{1/2}\phi_{NT}\sqrt{K/\phi_A})$ by the fact $K/\phi_A \geq 1$ according to the definition of ϕ_A .

Now, suppose that the sample size is sufficiently large so that $\zeta_1 + \zeta_2 \leq \tau\sqrt{K/\phi_A}/2$ in view of the rate of τ in Assumption 2(a). We push (A.33) further to attain

$$\boldsymbol{\psi}'\widehat{\mathbf{A}}'\widehat{\mathbf{A}}\boldsymbol{\psi} + 2\tau\|\widehat{\boldsymbol{\alpha}}\|_1 \leq \tau\sqrt{K/\phi_A}\|\boldsymbol{\psi}\|_1 + 2\tau\|\widehat{\boldsymbol{\alpha}}^*\|_1.$$

Then

$$\boldsymbol{\psi}'\widehat{\mathbf{A}}'\widehat{\mathbf{A}}\boldsymbol{\psi} \leq \tau\sqrt{K/\phi_A}\|\boldsymbol{\psi}\|_1 + 2\tau(\|\widehat{\boldsymbol{\alpha}}^*\|_1 - \|\widehat{\boldsymbol{\alpha}}\|_1) \leq \tau\left(\sqrt{K/\phi_A} + 2\right)\|\boldsymbol{\psi}\|_1,$$

where the last inequality follows by the triangle inequality: $\|\widehat{\boldsymbol{\alpha}}^*\|_1 - \|\widehat{\boldsymbol{\alpha}}\|_1 \leq \|\boldsymbol{\psi}\|_1$. Adding $\tau\sqrt{K/\phi_A}\|\boldsymbol{\psi}\|_1$ to both sides of the above inequality yields

$$\boldsymbol{\psi}'\widehat{\mathbf{A}}'\widehat{\mathbf{A}}\boldsymbol{\psi} + \tau\sqrt{K/\phi_A}\|\boldsymbol{\psi}\|_1 \leq 2\tau\left(\sqrt{K/\phi_A} + 1\right)\|\boldsymbol{\psi}\|_1 \leq 4\tau\sqrt{K/\phi_A}\|\boldsymbol{\psi}\|_1 \quad (\text{A.35})$$

where the last inequality follows the fact $K/\phi_A \geq 1$.

By $\phi_A^{-1} = O_p(K)$ in Lemma 5(a), we have

$$\sqrt{\frac{K/\phi_A}{N}}\phi_e = \sqrt{\frac{K/\phi_A}{N}}O_p(\sqrt{N}\phi_{NT}) = O_p\left(\sqrt{K/\phi_A}\phi_{NT}\right) = O_p(K\phi_{NT}) = O_p(K^{1/2}\tau) = o_p(1).$$

where the last two equalities hold by Assumption 2(a). This implies that the condition $\phi_e \leq \frac{1}{2}\sqrt{N\phi_A/K}$ in Lemma 4 is satisfied w.p.a.1. Moreover, $\boldsymbol{\psi} \in \mathfrak{S}^{\text{all}}$ by construction of $\widehat{\boldsymbol{\alpha}}^*$ in (4.1). We hence invoke Lemma 4 to continue (A.35):

$$4\tau\sqrt{K/\phi_A}\|\boldsymbol{\psi}\|_1 \leq 8\tau\frac{K/\phi_A}{\sqrt{N}}\|\widehat{\mathbf{A}}\boldsymbol{\psi}\|_2 \leq \frac{1}{2}\boldsymbol{\psi}'\widehat{\mathbf{A}}'\widehat{\mathbf{A}}\boldsymbol{\psi} + 32\tau^2\frac{(K/\phi_A)^2}{N} \quad (\text{A.36})$$

where the last inequality follows by $8ab \leq \frac{1}{2}a^2 + 32b^2$. Combining (A.35) and (A.36), we have

$$\frac{1}{2}\boldsymbol{\psi}'\widehat{\mathbf{A}}'\widehat{\mathbf{A}}\boldsymbol{\psi} + \tau\sqrt{K/\phi_A}\|\boldsymbol{\psi}\|_1 \leq 32\tau^2\frac{(K/\phi_A)^2}{N}. \quad (\text{A.37})$$

The above equality immediately implies

$$\|\boldsymbol{\psi}\|_1 = \|\widehat{\boldsymbol{\alpha}} - \widehat{\boldsymbol{\alpha}}^*\|_1 \leq 32\tau\frac{(K/\phi_A)^{3/2}}{N} = O_p\left(\frac{(K/\phi_A)^{3/2}\tau}{N}\right) = O_p\left(\frac{K^3\tau}{N}\right)$$

and

$$\sqrt{\boldsymbol{\psi}'\widehat{\mathbf{A}}'\widehat{\mathbf{A}}\boldsymbol{\psi}} = \left\|\widehat{\mathbf{A}}(\widehat{\boldsymbol{\alpha}} - \widehat{\boldsymbol{\alpha}}^*)\right\|_2 \leq 8\tau\frac{K/\phi_A}{\sqrt{N}} = O_p\left(\frac{(K/\phi_A)\tau}{\sqrt{N}}\right) = O_p\left(\frac{K^2\tau}{\sqrt{N}}\right), \quad (\text{A.38})$$

given $K/\phi_A = O_p(K^2)$ by Lemma 5(a). This completes the proof of the theorem. \blacksquare

Proof of Corollary 1. Recall that $\widehat{\mathbf{A}}^* = (\mathbf{I}_N - N^{-1}\mathbf{1}_N\mathbf{1}'_N)\widehat{\Sigma}^*$ and $\widehat{\mathbf{A}} = (\mathbf{I}_N - N^{-1}\mathbf{1}_N\mathbf{1}'_N)\widehat{\Sigma}$. Let $\widehat{\mathbf{A}}^e := (\mathbf{I}_N - N^{-1}\mathbf{1}_N\mathbf{1}'_N)\widehat{\Sigma}^e$. Then we have

$$\widehat{\mathbf{w}} - \mathbf{w}_0^* = \left(\widehat{\mathbf{A}}\widehat{\boldsymbol{\alpha}} + \frac{\mathbf{1}_N}{N} \right) - \left(\widehat{\mathbf{A}}^*\widehat{\boldsymbol{\alpha}}^* + \frac{\mathbf{1}_N}{N} \right) = \widehat{\mathbf{A}}\widehat{\boldsymbol{\alpha}} - \left(\widehat{\mathbf{A}} - \widehat{\mathbf{A}}^e \right) \widehat{\boldsymbol{\alpha}}^* = \widehat{\mathbf{A}}(\widehat{\boldsymbol{\alpha}} - \widehat{\boldsymbol{\alpha}}^*) + \widehat{\mathbf{A}}^e\widehat{\boldsymbol{\alpha}}^*. \quad (\text{A.39})$$

For the first term in (A.39), by (A.38) we have

$$\|\widehat{\mathbf{A}}(\widehat{\boldsymbol{\alpha}} - \widehat{\boldsymbol{\alpha}}^*)\|_2 = O_p\left(N^{-1/2}K^2\tau\right). \quad (\text{A.40})$$

For the second term in (A.39), we have $\|\widehat{\mathbf{A}}^e\widehat{\boldsymbol{\alpha}}^*\|_2 \leq \|\widehat{\Sigma}^e\widehat{\boldsymbol{\alpha}}^*\|_2 \leq \|\Sigma_0^e\widehat{\boldsymbol{\alpha}}^*\|_2 + \|\Delta^e\widehat{\boldsymbol{\alpha}}^*\|_2 := I_1 + I_2$ by the triangle inequality. Notice that

$$\begin{aligned} I_1 &\leq \phi_{\max}(\Sigma_0^e)\|\widehat{\boldsymbol{\alpha}}^*\|_2 \leq \phi_{\max}(\Sigma_0^e)\|\mathbf{a}_0^*\|_2 \\ &\leq O_p\left(\sqrt{N}\phi_{NT}\right) \frac{\|\mathbf{b}_0^*\|_\infty + 1}{N\sqrt{\phi_A}} = O_p\left(\frac{K^{1/2}\phi_{NT}}{\sqrt{N}\phi_A}\right), \end{aligned} \quad (\text{A.41})$$

where the second inequality follows as $\widehat{\boldsymbol{\alpha}}^* \in \widetilde{\mathbb{S}}^{\text{all}} \subset \mathbb{S}^{\text{all}}$ by construction, the third inequality by (A.25) and Assumption 1(a), and the last equality by and Lemma 5(a). Moreover,

$$\begin{aligned} I_2 &\leq \|\Delta^e\|_{c_2}\|\widehat{\boldsymbol{\alpha}}^*\|_1 \leq \sqrt{N}\|\Delta^e\|_\infty\|\widehat{\boldsymbol{\alpha}}^*\|_1 \\ &= \sqrt{N}\|\Delta^e\|_\infty\|\mathbf{a}_0^*\|_1 \leq \sqrt{N}O_p\left((T/\log N)^{-1/2}\right)N^{-1}\sqrt{K/\phi_A}(\|\mathbf{b}_0^*\|_\infty + 1) \\ &= O_p\left(K\phi_{NT}/\sqrt{N}\phi_A\right), \end{aligned} \quad (\text{A.42})$$

where the first inequality follows by (A.57), the first equality holds by the fact that $\|\widehat{\boldsymbol{\alpha}}^*\|_1 = \|\mathbf{a}_0^*\|_1 = \|\mathbf{a}_0^*\|_1$ as in (4.2), the third inequality by Assumption 1(a) and Lemma 4(c), and the last equality holds and Lemma 5(a).

Then collecting (A.39), (A.40), (A.41) and (A.42), we have

$$\|\widehat{\mathbf{w}} - \mathbf{w}_0^*\|_2 = O_p\left(N^{-1/2}\tau K^2\right) + O_p\left(K\phi_{NT}/\sqrt{N}\phi_A\right) = O_p\left(N^{-1/2}\tau K^2\right) = o_p\left(N^{-1/2}\right)$$

by Assumption 2(a) and Lemma 5(a). In addition, the Cauchy-Schwarz inequality immediately implies $\|\widehat{\mathbf{w}} - \mathbf{w}_0^*\|_1 \leq \sqrt{N}\|\widehat{\mathbf{w}} - \mathbf{w}_0^*\|_2 = O_p\left(K^2\tau\right) = o_p(1)$. ■

Proof of Theorem 3. Part (a). Denote $\boldsymbol{\psi}_w = \widehat{\mathbf{w}} - \mathbf{w}_0^*$. We first show the in-sample oracle inequality. Decompose

$$\begin{aligned} &\widehat{\mathbf{w}}'\widehat{\Sigma}\widehat{\mathbf{w}} - \mathbf{w}_0^{*'}\widehat{\Sigma}^*\mathbf{w}_0^* \\ &= (\mathbf{w}_0^{*'}\widehat{\Sigma}\mathbf{w}_0^* + 2\boldsymbol{\psi}_w'\widehat{\Sigma}\mathbf{w}_0^* + \boldsymbol{\psi}_w'\widehat{\Sigma}\boldsymbol{\psi}_w) - \mathbf{w}_0^{*'}\widehat{\Sigma}^*\mathbf{w}_0^* = \mathbf{w}_0^{*'}\widehat{\Sigma}^e\mathbf{w}_0^* + 2\boldsymbol{\psi}_w'\widehat{\Sigma}\mathbf{w}_0^* + \boldsymbol{\psi}_w'\widehat{\Sigma}\boldsymbol{\psi}_w \\ &= \mathbf{w}_0^{*'}\widehat{\Sigma}^e\mathbf{w}_0^* + 2\boldsymbol{\psi}_w'(\widehat{\Sigma}^* + \Delta^e)\mathbf{w}_0^* + 2\boldsymbol{\psi}_w'\Sigma_0^e\mathbf{w}_0^* + \boldsymbol{\psi}_w'(\widehat{\Sigma}^* + \Delta^e)\boldsymbol{\psi}_w + \boldsymbol{\psi}_w'\Sigma_0^e\boldsymbol{\psi}_w \\ &=: II_1 + 2II_2 + 2II_3 + II_4 + II_5. \end{aligned}$$

We bound II_1 by

$$\begin{aligned} |II_1| &\leq \phi_{\max}(\widehat{\Sigma}^e) \|\mathbf{w}_0^*\|_2^2 \leq (\phi_{\max}(\Sigma_0^e) + \phi_{\max}(\Delta^e)) \|\mathbf{w}_0^*\|_2^2 \\ &\leq (\phi_{\max}(\Sigma_0^e) + N\|\Delta^e\|_\infty) \|\mathbf{w}_0^*\|_2^2 \\ &\leq \left(O_p(\sqrt{N}\phi_{NT}) + NO_p((T/\log N)^{-1/2}) \right) \frac{\|\mathbf{b}_0^*\|_\infty^2}{N} = O_p(K\phi_{NT}), \end{aligned}$$

where the third inequality holds by the Gershgorin circle theorem, and the fourth by Assumption 1, and the last by Lemma 5(a). The second term II_2 is bounded by

$$\begin{aligned} |II_2| &\leq \|\widehat{\Sigma}^* + \Delta^e\|_\infty \|\psi_w\|_1 \|\mathbf{w}_0^*\|_1 \leq \left(\|\widehat{\Sigma}^*\|_\infty + \|\Delta^e\|_\infty \right) \|\psi_w\|_1 \|\mathbf{w}_0^*\|_1 \\ &\leq \left(\|\widehat{\Sigma}^{\text{co}}\|_\infty + \|\Delta^e\|_\infty \right) \|\psi_w\|_1 \sqrt{N} \|\mathbf{w}_0^*\|_2 \\ &= O_p\left(\bar{c} + (T/\log N)^{-1/2}\right) O_p(\tau K^2) \|\mathbf{b}_0^*\|_\infty = O_p\left(\tau K^{5/2}\right), \end{aligned}$$

where the first inequality follows by (A.54), the third inequality by the Cauchy-Schwarz inequality, the first equality holds by Assumptions 1, Corollary 1, and Theorem 1(a), and the last equality by Lemma 5(a). For II_3 , we have

$$\begin{aligned} |II_3| &\leq \phi_{\max}(\Sigma_0^e) \|\psi_w\|_2 \|\mathbf{w}_0^*\|_2 \\ &= O_p(\sqrt{N}\phi_{NT}) O_p\left(N^{-1/2}\tau K^2\right) \|\mathbf{b}_0^*\|_\infty N^{-1/2} = O_p\left(N^{-1/2}\phi_{NT}\tau K^{5/2}\right) \end{aligned}$$

by (A.55), Assumptions 1, and Corollary 1. Similarly,

$$\begin{aligned} |II_4| &\leq \|\widehat{\Sigma}^* + \Delta^e\|_\infty \|\psi_w\|_1^2 = O_p(\bar{c} + (T/\log N)^{-1/2}) O_p(\tau^2 K^4) = O_p(\tau^2 K^4), \text{ and} \\ |II_5| &\leq \phi_{\max}(\Sigma_0^e) \|\psi_w\|_2^2 = O_p(\sqrt{N}\phi_{NT}) (N^{-1}\tau^2 K^4) = O_p\left(N^{-1/2}\phi_{NT}\tau^2 K^4\right). \end{aligned}$$

Collecting all these five terms, and notice that $O_p(\tau K^{5/2})$ is the dominating order, we have

$$\left| \widehat{\mathbf{w}}' \widehat{\Sigma} \widehat{\mathbf{w}} - \mathbf{w}_0^{**'} \widehat{\Sigma}^* \mathbf{w}_0^* \right| = O_p\left(\tau K^{5/2}\right) = o_p(1)$$

under Assumption 2(a).

Part (b). The same argument goes through if we replace $\widehat{\Sigma}$ with $\widehat{\Sigma}^{\text{new}}$, and replace $\widehat{\Sigma}^*$ with $\widehat{\Sigma}^{*\text{new}}$, because in the above analysis of the in-sample oracle inequality we always bound the various quantities by separating the norms of vectors and the square matrices. We conclude the out-of-sample oracle inequality.

Part (c). This proof involves two steps: (i) Establish the closeness between $\widehat{\mathbf{w}}' \widehat{\Sigma}^{\text{new}} \widehat{\mathbf{w}}$ and $\widehat{\mathbf{w}}' \widehat{\Sigma} \widehat{\mathbf{w}}$ as shown in (A.47) below; (ii) Establish the closeness between $\widehat{\mathbf{w}}' \widehat{\Sigma} \widehat{\mathbf{w}}$ and $Q(\Sigma_0)$ as shown in (A.52) below, where $Q(\mathbf{S}) := \min_{\mathbf{w} \in \mathbb{R}^N, \mathbf{w}' \mathbf{1}_N = 1} \mathbf{w}' \mathbf{S} \mathbf{w}$ for a generic $N \times N$ positive semi-definite matrix \mathbf{S} .

Obviously $\widehat{\mathbf{w}}' \widehat{\Sigma}^{\text{new}} \widehat{\mathbf{w}} \geq \widehat{\mathbf{w}}' \widehat{\Sigma} \widehat{\mathbf{w}} = Q(\widehat{\Sigma})$. On the other hand,

$$\widehat{\mathbf{w}}' \widehat{\Sigma} \widehat{\mathbf{w}} = \widehat{\mathbf{w}}' \widehat{\Sigma}^{\text{new}} \widehat{\mathbf{w}} + \widehat{\mathbf{w}}' (\widehat{\Sigma} - \widehat{\Sigma}^{\text{new}}) \widehat{\mathbf{w}} \geq \widehat{\mathbf{w}}' \widehat{\Sigma}^{\text{new}} \widehat{\mathbf{w}} - \|\widehat{\Sigma} - \widehat{\Sigma}^{\text{new}}\|_{\text{sp}} \|\widehat{\mathbf{w}}\|_2^2 \quad (\text{A.43})$$

by the triangle inequality and (A.55). We focus on the term $\|\widehat{\Sigma} - \widehat{\Sigma}^{\text{new}}\|_{\text{sp}} \|\widehat{\mathbf{w}}\|_2^2$.

For the first factor, notice

$$\widehat{\Sigma} - \widehat{\Sigma}^{\text{new}} = \widehat{\Sigma}^* - \widehat{\Sigma}^{*,\text{new}} + \widehat{\Sigma}^e - \widehat{\Sigma}^{e,\text{new}} = (\widehat{\Sigma}^* - \Sigma_0^*) - (\widehat{\Sigma}^{*,\text{new}} - \Sigma_0^*) + \Delta^e - \Delta^{e,\text{new}}.$$

Under Assumption 1(b), $\|\widehat{\Sigma}^* - \Sigma_0^*\|_\infty = \|\widehat{\Sigma}^{\text{co}} - \Sigma_0^{\text{co}}\|_\infty = \|\Delta^{\text{co}}\|_\infty = O_p((T/\log N)^{-1/2})$ and therefore by Gershgorin circle theorem the spectral norm is bounded by

$$\|\widehat{\Sigma}^* - \Sigma_0^*\|_{\text{sp}} \leq N \|\Delta^{\text{co}}\|_\infty = O_p\left(N(T/\log N)^{-1/2}\right). \quad (\text{A.44})$$

Gershgorin circle theorem also implies $\|\Delta^e\|_{\text{sp}} \leq \phi_e = O_p(\sqrt{N}\phi_{NT})$, where the stochastic order follows by Lemma 5(a). Since the new testing data comes from the same data generating process as that of the training data, the same stochastic bounds is applicable to the terms involving the new data, and then

$$\begin{aligned} \|\widehat{\Sigma} - \widehat{\Sigma}^{\text{new}}\|_{\text{sp}} &\leq \|\widehat{\Sigma}^* - \Sigma_0^*\|_{\text{sp}} + \|\widehat{\Sigma}^{*,\text{new}} - \Sigma_0^*\|_{\text{sp}} + \|\Delta^e\|_{\text{sp}} + \|\Delta^{e,\text{new}}\|_{\text{sp}} \\ &= O_p\left(N(T/\log N)^{-1/2}\right) + O_p(\sqrt{N}\phi_{NT}) = O_p(N\phi_{NT}). \end{aligned} \quad (\text{A.45})$$

The second factor is bounded by

$$\|\widehat{\mathbf{w}}\|_2^2 \leq \|\mathbf{b}_0^*\|_\infty^2 / N = O_p(K/N) \quad (\text{A.46})$$

according to Theorem 1(a) and Lemma 5(a). Collecting (A.43), (A.45) and (A.46), we have

$$0 \leq \widehat{\mathbf{w}}' \widehat{\Sigma}^{\text{new}} \widehat{\mathbf{w}} - \widehat{\mathbf{w}}' \widehat{\Sigma} \widehat{\mathbf{w}} \leq \|\widehat{\Sigma} - \widehat{\Sigma}^{\text{new}}\|_{\text{sp}} \|\widehat{\mathbf{w}}\|_2^2 = O_p(N\phi_{NT}) O_p(K/N) = O_p(K\phi_{NT}). \quad (\text{A.47})$$

Next, consider the population matrices Σ_0 and Σ_0^* . Because $\Sigma_0 - \Sigma_0^* = \Sigma_0^e$ is positive semi-definite,

$$Q(\Sigma_0) \geq Q(\Sigma_0^*). \quad (\text{A.48})$$

Since $\text{rank}(\Sigma_0^*) = K \ll N$, the solution to $\min_{\mathbf{w} \in \mathbb{R}^N, \mathbf{w}' \mathbf{1}_N = 1} \mathbf{w}' \Sigma_0^* \mathbf{w}$ is not unique but all the solutions give the same minimum $Q(\Sigma_0^*)$. Thus in order to evaluate $Q(\Sigma_0^*)$ we can simply use the within-group equal weight optimizer $\mathbf{w}_0^\#$ which solves

$$\min_{(\mathbf{w}, \gamma) \in \mathbb{R}^{N+1}} \frac{1}{2} \|\mathbf{w}\|_2^2 \quad \text{subject to} \quad \mathbf{w}' \mathbf{1}_N = 1, \text{ and } \Sigma_0^* \mathbf{w} + \gamma = 0.$$

The only difference between $\mathbf{w}_0^\#$ and \mathbf{w}_0^* is that the former is associated with the population Σ_0^* and the latter associated with the sample $\widehat{\Sigma}^*$. Parallel to (3.7), (3.9) and (A.30), we can write $\mathbf{w}_0^\# = \left(\frac{b_{01}^\#}{N} \cdot \mathbf{1}'_{N_1}, \dots, \frac{b_{0K}^\#}{N} \cdot \mathbf{1}'_{N_K} \right)'$ where $\mathbf{b}_0^\# = \mathbf{r} \circ \frac{(\Sigma_0^{\text{co}})^{-1} \mathbf{1}_K}{\mathbf{1}'_K (\Sigma_0^{\text{co}})^{-1} \mathbf{1}_K}$, and it is bounded by

$$\begin{aligned} \|\mathbf{b}_0^\#\|_\infty &\leq \|(\Sigma_0^{\text{co}})^{-1} \mathbf{1}_K\|_\infty / [\underline{r} \cdot \mathbf{1}'_K (\Sigma_0^{\text{co}})^{-1} \mathbf{1}_K] \leq \underline{r}^{-1} K^{-1/2} \phi_{\max}^{1/2}(\Sigma_0^{\text{co}}) / \phi_{\min}(\Sigma_0^{\text{co}}) \\ &\leq \bar{c}^{1/2} / (\underline{r} c K^{1/2}) = O(\sqrt{K}) \end{aligned}$$

under Assumption 1(b) and furthermore

$$\|\mathbf{w}_0^\sharp\|_2^2 \leq N \|\mathbf{w}_0^\sharp\|_\infty^2 = N(\|\mathbf{b}_0^\sharp\|_\infty/N)^2 = O(K/N). \quad (\text{A.49})$$

We continue (A.48):

$$\begin{aligned} Q(\boldsymbol{\Sigma}_0^*) &= \mathbf{w}_0^\sharp \widehat{\boldsymbol{\Sigma}}^* \mathbf{w}_0^\sharp + \mathbf{w}_0^\sharp \left(\boldsymbol{\Sigma}_0^* - \widehat{\boldsymbol{\Sigma}}^* \right) \mathbf{w}_0^\sharp \geq \mathbf{w}_0^{*\prime} \widehat{\boldsymbol{\Sigma}}^* \mathbf{w}_0^* + \mathbf{w}_0^\sharp \left(\boldsymbol{\Sigma}_0^* - \widehat{\boldsymbol{\Sigma}}^* \right) \mathbf{w}_0^\sharp \\ &\geq \mathbf{w}_0^{*\prime} \widehat{\boldsymbol{\Sigma}}^* \mathbf{w}_0^* - \|\boldsymbol{\Sigma}_0^* - \widehat{\boldsymbol{\Sigma}}^*\|_{\text{sp}} \|\mathbf{w}_0^\sharp\|_2^2, \end{aligned} \quad (\text{A.50})$$

where the first inequality follows as \mathbf{w}_0^* is the optimizer associated with $\widehat{\boldsymbol{\Sigma}}^*$, and the second inequality is derived by the same reasoning as used to obtain (A.43). (A.50) and (A.48) imply

$$\mathbf{w}_0^{*\prime} \widehat{\boldsymbol{\Sigma}}^* \mathbf{w}_0^* \leq Q(\boldsymbol{\Sigma}_0) + \|\boldsymbol{\Sigma}_0^* - \widehat{\boldsymbol{\Sigma}}^*\|_{\text{sp}} \|\mathbf{w}_0^\sharp\|_2^2 \leq Q(\boldsymbol{\Sigma}_0) + O_p\left(K(T/\log N)^{-1/2}\right) \quad (\text{A.51})$$

in view of (A.44) and (A.49). Combine Part (a) and (A.51):

$$\widehat{\mathbf{w}} \widehat{\boldsymbol{\Sigma}} \widehat{\mathbf{w}} \leq Q(\boldsymbol{\Sigma}_0) + O_p(\tau K^{5/2}). \quad (\text{A.52})$$

In conjunction with (A.47) and notice $K\phi_{NT} = O(\tau K^{1/2})$ is of smaller order than $\tau K^{5/2}$, the conclusion follows. ■

A.5 Elementary Inequalities on Matrix Norms

We collect some elementary inequalities used in the proofs.

Lemma 6 *Let \mathbf{a} and \mathbf{b} be two vectors, and \mathbf{A} and \mathbf{B} be two matrices of compatible dimensions. Then we have*

$$\|\mathbf{A}\mathbf{b}\|_\infty \leq \|\mathbf{A}\|_\infty \|\mathbf{b}\|_1 \quad (\text{A.53})$$

$$|\mathbf{a}'\mathbf{A}\mathbf{b}| \leq \|\mathbf{A}\|_\infty \|\mathbf{a}\|_1 \|\mathbf{b}\|_1 \quad (\text{A.54})$$

$$|\mathbf{a}'\mathbf{A}\mathbf{b}| \leq \|\mathbf{A}\|_{\text{sp}} \|\mathbf{a}\|_2 \|\mathbf{b}\|_2 \quad (\text{A.55})$$

$$\|\mathbf{A}\mathbf{B}\mathbf{b}\|_\infty \leq \|\mathbf{A}'\|_{c_2} \|\mathbf{B}\|_{c_2} \|\mathbf{b}\|_1 \quad (\text{A.56})$$

If \mathbf{S} is a symmetric matrix,

$$\|\mathbf{S}\mathbf{b}\|_2 \leq \|\mathbf{S}\|_{c_2} \|\mathbf{b}\|_1. \quad (\text{A.57})$$

If $\boldsymbol{\Sigma}$ is positive definite,

$$(\boldsymbol{\Sigma} + \mathbf{a}\mathbf{a}')^{-1} = \boldsymbol{\Sigma}^{-1} - \frac{\boldsymbol{\Sigma}^{-1}\mathbf{a}\mathbf{a}'\boldsymbol{\Sigma}^{-1}}{1 + \mathbf{a}'\boldsymbol{\Sigma}^{-1}\mathbf{a}} \quad (\text{A.58})$$

$$(\boldsymbol{\Sigma} - \mathbf{a}\mathbf{a}')^{-1} = \boldsymbol{\Sigma}^{-1} + \frac{\boldsymbol{\Sigma}^{-1}\mathbf{a}\mathbf{a}'\boldsymbol{\Sigma}^{-1}}{1 - \mathbf{a}'\boldsymbol{\Sigma}^{-1}\mathbf{a}}. \quad (\text{A.59})$$

Proof. The first inequality follows because

$$\|\mathbf{A}\mathbf{b}\|_\infty \leq \max_i |\mathbf{A}_i \cdot \mathbf{b}| \leq \max_i \|\mathbf{A}_i\|_\infty \|\mathbf{b}\|_1 = \|\mathbf{A}\|_\infty \|\mathbf{b}\|_1.$$

It implies the second inequality $|\mathbf{a}'\mathbf{A}\mathbf{b}| \leq \|\mathbf{a}\|_1 \|\mathbf{A}\mathbf{b}\|_\infty \leq \|\mathbf{A}\|_\infty \|\mathbf{a}\|_1 \|\mathbf{b}\|_1$ and the fourth inequality

$$\|\mathbf{A}\mathbf{B}\mathbf{b}\|_\infty \leq \|\mathbf{A}\mathbf{B}\|_\infty \|\mathbf{b}\|_1 = \max_{i,j} |\mathbf{A}_i \cdot \mathbf{B}_j| \|\mathbf{b}\|_1 \leq \|\mathbf{A}'\|_{c2} \|\mathbf{B}\|_{c2} \|\mathbf{b}\|_1.$$

The third inequality follows by the Cauchy-Schwarz inequality $|\mathbf{a}'\mathbf{A}\mathbf{b}| \leq \|\mathbf{A}'\mathbf{a}\|_2 \|\mathbf{b}\|_2 \leq \|\mathbf{A}\|_{\text{sp}} \|\mathbf{a}\|_2 \|\mathbf{b}\|_2$. For the symmetric matrix \mathbf{S} ,

$$\|\mathbf{S}\mathbf{b}\|_2 = \sqrt{\mathbf{b}'\mathbf{S}\mathbf{S}\mathbf{b}} \leq \sqrt{\|\mathbf{S}\mathbf{S}\|_\infty} \|\mathbf{b}\|_1 \leq \sqrt{\max_i (\mathbf{S}\mathbf{S})_{ii}} \|\mathbf{b}\|_1 = \|\mathbf{S}\|_{c2} \|\mathbf{b}\|_1$$

where the first inequality follows by (A.54), and the second inequality and the last equality are due to the symmetry of \mathbf{Q} .

The Sherman-Morrison formula gives $(\mathbf{\Sigma} + \mathbf{a}\mathbf{b}')^{-1} = \mathbf{\Sigma}^{-1} - \mathbf{\Sigma}^{-1}\mathbf{a}\mathbf{b}'\mathbf{\Sigma}^{-1} / (1 + \mathbf{a}'\mathbf{\Sigma}^{-1}\mathbf{b})$ for any compatible vector \mathbf{a} and \mathbf{b} . (A.58) and (A.59) follow by setting $\mathbf{b} = \mathbf{a}$ and $\mathbf{b} = -\mathbf{a}$, respectively. ■

B Additional Results for the Numerical Work

In this appendix, we report additional designs and results for the numerical work in the paper.

B.1 Simulation of the Demonstrative Example in Figure 1

The demonstrative example in Figure 1 is generated from the DGP

$$y_{t+1} = \sum_{i=1}^{20} w_i f_{it} + u_{t+1}, \quad t = 1, \dots, 100.$$

The dependent variable y_{t+1} is a linear combination of two groups of input variables $\{f_{it}\}_{i=1}^{10}$ and $\{f_{it}\}_{i=11}^{20}$ with group weights $w_i = 0.09 \cdot 1\{1 \leq i \leq 10\} + 0.01 \cdot 1\{11 \leq i \leq 20\}$ and $1\{\cdot\}$ denoting the indicator function, so that $\sum_{i=1}^{20} w_i = 1$. We set $f_{it} \sim \text{i.i.d. } N(1, 1)$ for $1 \leq i \leq 10$, $f_{it} \sim \text{i.i.d. } N(0, 1)$ for $11 \leq i \leq 20$, and $u_{t+1} \sim \text{i.i.d. } N(0, 0.25)$. We estimate the weights by the ℓ_2 -relaxation method under different values of $\tau = 0, 0.005, 0.01, \dots, 0.1$. Let $\hat{w}_1 = \hat{w}_{1\tau}$ denote the first element of the ℓ_2 -relaxation estimator $\hat{\mathbf{w}}_\tau$. We report in Figure 1 the empirical squared bias, variance and MSE of \hat{w}_1 and those of the one-step-ahead forecast \hat{y}_{n+1} with $n = 100$ over 1000 replications as a function of τ .

B.2 Calculation of SNR in the Section 5

To obtain SNR in the simulations, we decompose $y_{t+1} = \mathbf{w}_\psi^{*'}(\mathbf{f}_t - \mathbf{u}_t) + u_{y,t+1}$ into

$$y_{t+1} = \underbrace{\mathbf{w}_\psi^{*'}(\mathbf{f}_t - \mathcal{P}[\mathbf{u}_t|\mathbf{f}_t])}_{\text{signal}} - \underbrace{\mathbf{w}_\psi^{*'}(\mathbf{u}_t - \mathcal{P}[\mathbf{u}_t|\mathbf{f}_t])}_{\text{noise}} + u_{y,t+1},$$

where $\mathcal{P}[\mathbf{u}_t|\mathbf{f}_t] = E[\mathbf{u}_t\mathbf{f}_t'] (E[\mathbf{f}_t\mathbf{f}_t']^{-1} \mathbf{f}_t = \mathbf{\Omega}_u (\mathbf{\Psi} + \mathbf{\Omega}_u)^{-1} \mathbf{f}_t$ is the projection of \mathbf{u}_t onto the linear space spanned by \mathbf{f}_t . By construction, the signal and noise terms are orthogonal in that

$$\begin{aligned} & E[\mathbf{w}_\psi^{*'} (\mathbf{f}_t - \mathcal{P}[\mathbf{u}_t|\mathbf{f}_t]) (\mathbf{u}_t - \mathcal{P}[\mathbf{u}_t|\mathbf{f}_t])' \mathbf{w}_\psi^*] \\ &= \mathbf{w}_\psi^{*'} E\left[\left(\mathbf{f}_t - \mathbf{\Omega}_u (\mathbf{\Psi} + \mathbf{\Omega}_u)^{-1} \mathbf{f}_t\right) \left(\mathbf{u}_t - \mathbf{\Omega}_u (\mathbf{\Psi} + \mathbf{\Omega}_u)^{-1} \mathbf{f}_t\right)'\right] \mathbf{w}_\psi^* = 0. \end{aligned}$$

Simple calculation shows that the variance of the noise component is

$$\text{var}[\mathbf{w}_\psi^{*'} (\mathbf{u}_t - \mathcal{P}[\mathbf{u}_t|\mathbf{f}_t]) + u_{y,t+1}] = \mathbf{w}_\psi^{*'} \left[\mathbf{\Omega}_u - \mathbf{\Omega}_u (\mathbf{\Psi} + \mathbf{\Omega}_u)^{-1} \mathbf{\Omega}_u \right] \mathbf{w}_\psi^* + \sigma_y^2.$$

This, along with the fact that $\text{var}[y_{t+1}] = \mathbf{w}_\psi^{*'} \mathbf{\Psi} \mathbf{w}_\psi^* + \sigma_y^2$, implies that SNR here can be defined as

$$\text{SNR} = \frac{\text{var}[y_{t+1}] - \text{var}[\mathbf{w}_\psi^{*'} (\mathbf{u}_t - \mathcal{P}[\mathbf{u}_t|\mathbf{f}_t]) + u_{y,t+1}]}{\text{var}[\mathbf{w}_\psi^{*'} (\mathbf{u}_t - \mathcal{P}[\mathbf{u}_t|\mathbf{f}_t]) + u_{y,t+1}]} = \frac{\mathbf{w}_\psi^{*'} \left[\mathbf{\Psi} - \mathbf{\Omega}_u + \mathbf{\Omega}_u (\mathbf{\Psi} + \mathbf{\Omega}_u)^{-1} \mathbf{\Omega}_u \right] \mathbf{w}_\psi^*}{\mathbf{w}_\psi^{*'} \left[\mathbf{\Omega}_u - \mathbf{\Omega}_u (\mathbf{\Psi} + \mathbf{\Omega}_u)^{-1} \mathbf{\Omega}_u \right] \mathbf{w}_\psi^* + \sigma_y^2}.$$

B.3 Simulation Results under DGP2 by Conventional 5-fold CV

The MSFEs under DGP2 by the conventional 5-fold CV are close to those in Table 2 and all patterns remain.

Table 7: MSFE for DGP2 by Conventional 5-fold CV

T	N	K	Oracle	SA	Lasso	Ridge	PC			ℓ_2 -relax
							$q = 5$	$q = 10$	$q = 20$	
<i>Panel A: Low SNR</i>										
50	100	2	0.259	1.129	0.708	1.484	0.834	0.828	0.939	0.272
100	200	4	0.161	2.980	0.379	1.169	1.138	1.257	1.218	0.196
200	300	6	0.117	4.311	0.206	0.419	1.368	1.430	1.339	0.118
<i>Panel B: High SNR</i>										
50	100	2	0.263	1.078	0.483	0.512	0.769	0.838	0.820	0.270
100	200	4	0.149	3.101	0.200	0.262	1.081	1.202	1.291	0.155
200	300	6	0.096	4.452	0.129	0.139	1.265	1.367	1.459	0.121

B.4 MAFE for the Simulations

In this section we report the mean absolute forecast error (MAFE) in the simulations. The MAFE is defined as

$$\text{MAFE} = E[|y_{T+1} - \hat{\mathbf{w}}' \mathbf{f}_{T+1}|] - \sigma_y \sqrt{2/\pi},$$

where the unpredictable component, $\sigma_y \int_{-\infty}^{\infty} |x| \frac{1}{\sqrt{2\pi}} \exp(-x^2/2) dx = \sigma_y \sqrt{2/\pi}$, is subtracted. In simulations we know σ_y but it is unknown in empirical applications.

The results are collected in Table 8. Results by oracle, Lasso, Ridge, and the ℓ_2 -relaxation tend to decrease with T , yet results by SA and PC may diverge as T increases. Similar to the MSFE results in the main text, the ℓ_2 -relaxation is always the best feasible estimator in all cases.

Table 8: MAFE for Simulations

T	N	K	Oracle	SA	LASSO	Ridge	PC			ℓ_2 -relax
							$q = 5$	$q = 10$	$q = 20$	
<i>Panel A: DGP1 with low SNR</i>										
50	100	2	0.336	0.721	0.457	0.472	0.604	0.614	0.609	0.346
100	200	4	0.224	1.315	0.273	0.327	0.726	0.799	0.798	0.245
200	300	6	0.177	1.595	0.209	0.215	0.809	0.837	0.849	0.195
<i>Panel B: DGP1 with high SNR</i>										
50	100	2	0.063	0.356	0.233	0.422	0.278	0.278	0.305	0.085
100	200	4	0.076	0.845	0.153	0.351	0.377	0.398	0.448	0.099
200	300	6	0.017	1.067	0.051	0.104	0.432	0.403	0.404	0.034
<i>Panel C: DGP2 with low SNR</i>										
50	100	2	0.340	0.753	0.479	0.486	0.612	0.650	0.643	0.343
100	200	4	0.236	1.313	0.285	0.336	0.740	0.780	0.809	0.259
200	300	6	0.177	1.608	0.210	0.216	0.791	0.834	0.862	0.195
<i>Panel D: DGP2 with high SNR</i>										
50	100	2	0.105	0.378	0.246	0.388	0.290	0.286	0.325	0.111
100	200	4	0.081	0.777	0.173	0.389	0.356	0.407	0.375	0.089
200	300	6	0.053	1.053	0.094	0.158	0.428	0.447	0.428	0.069
<i>Panel E: DGP3 with low SNR</i>										
50	100	2	0.435	0.799	0.578	0.599	0.720	0.709	0.705	0.441
100	200	4	0.364	1.402	0.426	0.488	0.961	1.018	1.032	0.383
200	300	6	0.350	1.609	0.372	0.390	1.083	1.120	1.126	0.364
<i>Panel F: DGP3 with high SNR</i>										
50	100	2	0.163	0.406	0.387	0.550	0.356	0.339	0.360	0.197
100	200	4	0.073	0.856	0.191	0.426	0.448	0.463	0.506	0.090
200	300	6	0.047	1.085	0.176	0.290	0.586	0.625	0.603	0.054

B.5 MAFE for the Empirical Applications

Tables 9–11 report the relative (to the benchmark) MAFEs in the three empirical applications, respectively. The results and patterns are robust in comparison with those based on MSFE.

Table 9: Relative MAFE of Box Office Forecast

n_{ev}	PMA	Lasso	Ridge	ℓ_2 -relax
10	1.000	0.999	1.078	0.985
20	1.000	1.003	1.072	0.984
30	1.000	1.003	1.064	0.976
40	1.000	1.008	1.074	0.965

Note: The MAFE of PMA is normalized as 1.

Table 10: Relative MAFE of Inflation Forecast

horizon	SA	Lasso	Ridge	ℓ_2 -relax
One-year-ahead	1.000	0.971	0.996	0.940
Two-year-ahead	1.000	0.857	0.994	0.709

Note: The MAFE of SA is normalized as 1.

Table 11: Relative MAFE of Stock Volatility Forecast

horizon h	HAR	OLS	Lasso	Ridge	ℓ_2 -relax
1	1.0000	1.2700	0.9510	0.9389	0.9318
5	1.0000	1.1610	0.8914	0.7726	0.8154
10	1.0000	1.2781	0.9501	0.8446	0.8884
22	1.0000	1.3716	1.0596	1.0470	0.9397

Note: The MAFE of HAR is normalized as 1.