

Augment Large Covariance Matrix Estimation With Auxiliary Network Information

Shuyi Ge*, Shaoran Li,[†] Oliver Linton,[‡] and Weiguang Liu[§]

Faculty of Economics, University of Cambridge

June 14, 2021

Abstract

This paper aims to incorporate auxiliary information about the location of significant correlations into the estimation of high-dimensional covariance matrices. With the development of machine learning techniques such as textual analysis, granular linkage information among firms that used to be notoriously hard to get are now becoming available to researchers. Our proposed method provides an avenue for combining those auxiliary network information with traditional economic datasets to improve the estimation of a large covariance matrix. Simulation results show that the proposed adaptive correlation thresholding method generally performs better in the estimation of covariance matrices than previous methods, especially when the true covariance matrix is sparse and the auxiliary network contains genuine information. As a preliminary application, we apply the method to the estimation of the covariance matrix of asset returns. There are several extensions and improvements that we are considering.

* Author email: sg751@cam.ac.uk

[†] Author email: sl736@cam.ac.uk

[‡] Author email: obl20@cam.ac.uk

[§] Author email: wl342@cam.ac.uk

1 Introduction and Literature Review

Covariance matrix estimation is an important area of research in both finance and statistics. Suppose we have independent observations $X_t = (X_{1t}, \dots, X_{Nt})^T$, $t = 1, \dots, T$ of a N -dimensional random vector \mathbf{X} that has mean μ and variance $\Sigma_X = E((\mathbf{X} - \mu)(\mathbf{X} - \mu)')$. The most straightforward estimator is the sample covariance estimator, which is defined as follows:

$$\hat{\Sigma}_X = \frac{1}{T} (X - \bar{X})(X - \bar{X})^\top = [\hat{\sigma}_{ij}]_{N \times N}, \quad (1)$$

where X is the $N \times T$ matrix of observations; $\bar{X} = \frac{1}{T} X 1_T$ is the sample time series average, with 1_T being a $T \times 1$ vectors of 1. However, in the high-dimensional case, where the dimension N is not negligible comparing to sample size T , the sample covariance matrix is ill-conditioned and inconsistent.

One of the structures often imposed in the high-dimensional settings is sparsity, which assumes that Σ_X is sparse (i.e., has lots of zeros or small elements) or conditionally sparse (i.e., has lots of zeros or small elements once we condition on some variables like the common risk factors).

Given the sparsity structure, several estimation strategies have been proposed in the literature such as banding, tapering, shrinkage, thresholding, etc.

Banding and tapering are applicable when X_{it} are indexed and we have a way to define a “distance” $d(i, j)$ between X_i and X_j . For example, in time series, a natural distance is $|i - j|$ and we have reason to believe that large distance between X_i and X_j will imply lower correlation. Such structure is appropriate for applications where there are natural orderings of variables, such as time series, climatology and spectroscopy. Banding methods (Bickel and Levina, 2008b) keeps only elements within a k -neighbourhood of each individual X_i , that is

$$\hat{\Sigma}_{B,k} = [\hat{\sigma}_{ij} \mathbf{1}_{d(i,j) \leq k}]$$

whereas the tapering estimator does a elementwise multiplication of $\hat{\Sigma}_X$ with a positive definite matrix \mathbf{T} that has smaller element \mathbf{T}_{ij} decreases with $d(i, j)$.

When we don't have the information about the distance among i, j , one of the viable approaches is thresholding. Let $s_\lambda(\cdot)$ be a generalized thresholding operator satisfying

1. $|s_\lambda(z)| \leq c|y|$ for all $|z - y| \leq \lambda$
2. $s_\lambda(z) = 0$ for $|z| \leq \lambda$
3. $|s_\lambda(z) - z| \leq \lambda$.

which

incorporates commonly used thresholding operators such as hard thresholding, soft thresholding, SCAD, etc. Then a thresholding estimator is

$$\hat{\Sigma}_\lambda = [\tilde{\sigma}_{ij}] \quad \text{and} \quad \tilde{\sigma}_{ij} = \begin{cases} \hat{\sigma}_{ii} & i = j \\ s_\lambda(\hat{\sigma}_{ij}) & i \neq j \end{cases}$$

Bickel and Levina, 2008a develop the theory for universal thresholding which assumes that the diagonal of Σ is uniformly bounded. Cai and W. Liu, 2011 proposes an adaptive thresholding estimator. They relax the uniform boundedness assumption and accounts for the variance of estimator of each $\hat{\sigma}_{ij}$ the establishes entry-adaptive threshold to $\hat{\Sigma}_X$. Fan, Liao, and Mincheva, 2013 argue that common factors should be extract first before applying thresholding selection when there are "extremely spiked" eigenvalues in $\hat{\Sigma}_X$ (i.e., the covariance matrix is conditionally sparse). Shu and Nan, 2019 obtains the convergence rate allowing for temporal dependence. Bickel and Levina, 2008a also compares the convergence rates of banding estimator and thresholding estimator and by utilizing the location information, banding estimator shows a superior convergence rate. See Fan, Liao, and H. Liu, 2015 for a review on estimation of large dimensional covariance matrices.

Apart from elementwise regularization methods, Ledoit and Wolf, 2004a and Ledoit and Wolf, 2012 have proposed linear and nonlinear shrinkage estimators that apply shrinkage to the eigenvalues of the sample covariance. The linear shrinkage does that by finding the linear combination of sample covariance and a well-conditioned matrix such as the identity matrix and nonlinear shrinkage estimator corrects the eigenvalues using the asymptotic Marcenko–Pastur distribution. Shrinkage estimators have been successfully applied the estimators in portfolio construction (Ledoit and Wolf, 2004b, Ledoit and Wolf, 2017).

Although in general we won't have an ordering or know distance between (i, j) , we do have some idea about who might be connected with whom using augment information apart from the observations of X . To proxy for pairwise connectivity among entities, apart from purely statistical correlation, several methods have been proposed. Hoberg and Phillips, 2016 use textual analysis to identify peers. Kaustia and Rantala, 2013 identify peers analyst co-coverage, and Ge and O. B. Linton, 2021 identify peers using business news co-mentioning. Those network information gathered from other

sources can help us to identify the locations of non-zeros apart from the information from X itself. Similar location-based thresholding ideas have applied in Fan, Furger, and Xiu, 2016b and Brownlees, Gumundsson, and Lugosi, 2020. Fan, Furger, and Xiu, 2016b apply a hard thresholding method in a way that $\sigma_{ij} = 0$ when i and j are from different sector/industry. Brownlees, Gumundsson, and Lugosi, 2020 first detect community structure using a spectral clustering-based procedure, and then apply a block-by-block thresholding to the off-diagonal elements of $\hat{\Sigma}_X$. In particular, they do not apply and thresholding to $\hat{\sigma}_{ij}$ if i and j are from the same community.

We argue in this paper that we can incorporate such auxiliary information in the estimation of the covariance matrix Σ when they help reveal the locations of the larger elements (or nonzero elements in the strictly sparse case). The augmented thresholding estimator will relax the condition put on the sparsity of the covariance Σ and show superior performance.

In this paper, we utilize granular network information gathered from other sources that could imply the locations of non-zeros in the de-factored residuals covariance matrix. The first candidate is the new-implied network. It has been documented that common news coverage reveals information about linkages among companies, which are related to many economically important relationships like business alliances, partnerships, banking and financing, customer-supplier, and production similarity (Scherbina and Schlusche, 2015, Schwenkler and Zheng, 2019). Ge and O. B. Linton, 2021 document that stocks linked by news co-mentioning exhibit additional co-movement beyond what can be explained by common risk factors. Same as Ge and O. B. Linton, 2021, we use news data from RavenPack Equity files Dow Jones Edition for the period between the beginning of 2004 to the end of 2015. This comprehensive news dataset combines relevant content from multiple sources, including Dow Jones Newswires, Wall Street Journal, and Barron’s MarketWatch, which produce the most actively monitored streams of news articles in the financial system. We identify linkages among firms by news co-mentioning. The second candidate network is IBES analyst co-coverage network. Israelsen, 2016 documents that stocks linked by analysts exhibit excess co-movement. To construct the analyst co-coverage-based adjacency matrix, we use the Institutional Brokers Estimate System (IBES) detail history files. (after getting the network, our procedure)?

Although here we are applying augment network information to the estimation of

large static covariance matrix, similar idea can be extended to the estimation of large dynamic covariance matrix. For example, dynamic network information could be well incorporated into the conditioning information set in Chen, D. Li, and O. Linton, 2019.

2 Literature Review

There has been extensive research on high-dimensional covariance estimation. Some important lines of thinking include element-wise banding and thresholding method, shrinkage method, factor models, etc. For a book-length review see Pourahmadi, 2013.

Bickel and Levina, 2008a considers banding or tapering the sample covariance matrix. Bickel and Levina, 2008a considers covariance regularization by hard thresholding. They also compare the results between banding when there is a natural ordering (for example, time series autocorrelation) and thresholding where we need to pay a log p price in the convergence rate to learn the locations. Cai and W. Liu, 2011 considers adaptive thresholding where threshold takes the form:

$$\hat{\sigma}_{ij}^* = s_{t_{ij}}(\hat{\sigma}_{ij}) \quad (2)$$

where 1. $|s_{\lambda}(z)| \leq c|y|$ for all $|z - y| \leq \lambda$ 2. $s_{\lambda}(z) = 0$ for $|z| \leq \lambda$ 3. $|s_{\lambda}(z) - z| \leq \lambda$. The convergence rate is the same, although here the uniformity class is larger. Fan, Liao, and H. Liu, 2015 proposes thresholding on the correlation matrix. The choice of thresholding functions can be found Rothman, Levina, and Zhu, 2009, Fan and R. Li, 2001, etc.

As an application of thresholding method, Fan, Furger, and Xiu, 2016a use hard thresholding method in a high-frequency setting based on the sector/industry classifier. $s_{ij}(\sigma_{ij}) = \sigma_{ij}$ if ij are in the same industry. The network they use is a block-diagonal matrix and our results accommodate more general and flexible network information.

Ledoit and Wolf, 2004b develops an estimation strategy based on linear shrinkage, where the target is identity matrix. This shrinkage guarantees that the estimated covariance matrix is well-conditioned. This approach can be thought of as decreasing

variance at the expense of increasing bias a little. There are articles discuss multiple targets, for example, Schäfer and Strimmer, 2005, Lancewicki and Aladjem, 2014 and Gray et al., 2018, but their targets are either fixed or data-driven, so different from our guided method where we bring in new information from auxiliary network information. Ledoit and Wolf, 2012 and Ledoit and Wolf, 2017 propose nonlinear shrinkage where the eigenvalues are pulled towards the “correct level” solving a nonrandom limit loss function. The shrinkage method has been shown to have really good performance in estimating large-dimensional covariance matrix, however they are a global method whereas our method is designed to emphasize “economically meaning” links. There is also a vast literature on factor models in high-dimensional models and applications in empirical finance. We refer to Connor, Hagmann, and O. Linton, 2012, Fan, Liao, and H. Liu, 2015 and Fan, Liao, and Wang, 2016 and literature review therein.

3 Estimator and Convergence Rate

Assume we have observations $X = (X_1, \dots, X_T)$, where X_t are independent drawn from a N -dimensional distribution F with mean μ and variance Σ . In Bickel and Levina, 2008a, they consider the following uniformity class of covariance matrices:

$$\mathcal{U}_\tau(q, c_0(p), M) = \left\{ \Sigma : \sigma_{i,i} \leq M, \sum_{j=1}^p |\sigma_{ij}|^q \leq c_0(p), \quad \text{for all } i \right\}$$

And the convergence rate will depend on $c_0(p)$ and q . Notice that in order to bound the i -th row sparsity index $\mathcal{S}_i(q) := \sum_{j=1}^p |\sigma_{ij}|^q$ by $c_0(p)$, the coefficient q is important. Suppose $q = 0$, then $\mathcal{S}_i = \#\{\sigma_{ij} \neq 0\}$. If the i -th row Σ_i has a lot of nonzero but small elements, then to bound $\mathcal{S}_i(0)$ would require a higher $c_0(p)$. On the other hand when $q \rightarrow 1$, the large elements of Σ_i will dominate. Hence if Σ is sparse in the sense that it contains a small number of relatively large elements and a large number of small elements, it's advisable to state conditions separately for these elements, and we consider the following uniformity class:

$$\mathcal{U}(q, c_0, c_1, M, L) = \left\{ \Sigma : \sigma_{ii} \leq M, \sum_j L_{ij}^1 \leq c_1(p), \sum_j L_{ij}^0 |\sigma_{ij}|^q \leq c_0(p) \quad \text{for all } i \right\}$$

where L_{ij} represents the location of the large elements and for $s \in \{0, 1\}$, $L_{ij}^s = \mathbf{1}_{L_{ij}=s}$. This uniformity class controls the number of the large elements at locations $((i, j) : L_{ij}^1 = 1)$ and the growth rate of the remaining small elements.

Of course, a priori we don't know the location of the large elements, but suppose we have observations from auxiliary dataset that allow us to form an estimator \hat{L} for L , independent of the sample X , we can design an estimator that takes into account the additional information in \hat{L} . A simple choice is to do banding based on the location information and apply thresholding on the remainder terms. Here we define a *Network Guided Estimator* to be

$$T_{L,\lambda}(\hat{\Sigma}) = [s_{L,\lambda}(\hat{\sigma}_{ij})]_{N \times N}$$

$$s_{L,\lambda}(\sigma_{ij}) = \begin{cases} \sigma_{ij} & \text{if } i = j \text{ or } L_{ij} = 1 \\ s_{\lambda}(\sigma_{ij}) & \text{otherwise} \end{cases}$$

where $s_{\lambda}(x)$ is the generalized thresholding operator and $\hat{\sigma}_{ij}$ are elements of the sample covariance matrix, then the feasible Network Guided Estimator is $T_{\hat{L},\lambda}(\hat{\Sigma})$

Assumption 1. We make the following assumptions:

1. $\max_{ij} |\hat{\sigma}_{ij} - \sigma_{ij}| = O_p(\sqrt{\log N/T})$;
2. $\max_{ij} |\hat{L}_{ij} - L_{ij}| = O_p(k_T)$ where $k_T \rightarrow 0$ as $T \rightarrow \infty$.

Remark. 1. The first assumption can be verified in various settings, for example, if F is Gaussian or sub-Gaussian (Cai and W. Liu, 2011). We can even replace the independence assumption and allow for temporal dependence (see Lemma A.2 in Shu and Nan, 2019).

2. The second assumption appears restrictive, but since L_{ij} are estimated independently from different datasets, we hope it's not too stringent to require that for each (i, j) , L_{ij} can be estimated consistently. In addition, simulation shows that as long as we don't make too much II-type error: $\hat{L}_{ij} = 1$ when $L_{ij} = 0$, using the additional information still improves the performance.

Given that the estimation of \hat{L}_{ij} is independent of the sample (X_t) , then perhaps we can find a less restrictive condition for the result to hold.

Theorem 1. Suppose F is Gaussian and for sufficiently large M ,

$$\lambda = M\sqrt{\frac{\log N}{T}}$$

and $\frac{\log N}{T} \rightarrow 0$ as $T \rightarrow \infty$, then for the operator norm $\|M\| = \max_j |\lambda_j(M)|$ where $\lambda_1, \dots, \lambda_N$ are the eigenvalues of M , we have

$$\left\| T_{\hat{L}, \lambda}(\hat{\Sigma}) - \Sigma \right\| = O_p \left(c_1(p) \sqrt{\frac{\log N}{T}}, c_0(p) \left(\frac{\log N}{T} \right)^{\frac{1-q}{2}} \right)$$

Proof. We have the following decomposition:

$$\left\| T_{\hat{L}, \lambda}(\hat{\Sigma}) - \Sigma \right\| \leq \left\| T_{\hat{L}, \lambda}(\Sigma) - \Sigma \right\| + \left\| T_{\hat{L}, \lambda}(\hat{\Sigma}) - T_{\hat{L}, \lambda}(\Sigma) \right\| = \mathbf{I} + \mathbf{II}$$

The first term can be bounded by

$$\begin{aligned} \mathbf{I} &\leq \max_i \sum_j \left| s_{\hat{L}, \lambda}(\sigma_{ij}) - \sigma_{ij} \right| \\ &= \max_i \sum_j \hat{L}_{ij}^0 |s_{\lambda}(\sigma_{ij}) - \sigma_{ij}| \\ &= \max_i \sum_j \left[L_{ij}^0 |s_{\lambda}(\sigma_{ij}) - \sigma_{ij}| + (\hat{L}_{ij}^0 - L_{ij}^0) |s_{\lambda}(\sigma_{ij}) - \sigma_{ij}| \right] \\ &\leq (1 + o_p(1)) \max_i \sum_j \left[L_{ij}^0 |s_{\lambda}(\sigma_{ij}) - \sigma_{ij}| \right] \\ &\leq (1 + o_p(1)) \max_i \sum_j \left[L_{ij}^0 |\sigma_{ij}| \mathbf{1}\{\sigma_{ij} \leq \lambda\} + (s_{\lambda}(\sigma_{ij}) - \sigma_{ij}) \mathbf{1}\{\sigma_{ij} > \lambda\} \right] \\ &\leq (1 + o_p(1)) \max_i \sum_j \left[L_{ij}^0 |\sigma_{ij}|^q \lambda^{1-q} \right] \\ &\leq (1 + o_p(1)) c_0(p) \lambda^{1-q} \end{aligned}$$

And the second term can be bounded similar to Rothman, Levina, and Zhu, 2009,

$$\begin{aligned}
\Pi &\leq \max_i \sum_j [\hat{L}_{ij}^1 |\hat{\sigma}_{ij} - \sigma_{ij}| + \hat{L}_{ij}^0 |s_\lambda(\hat{\sigma}_{ij}) - s_\lambda(\sigma)_{ij}|] \\
&\leq (1 + o_p(1)) c_1(p) \max_{ij} |\hat{\sigma}_{ij} - \sigma_{ij}| + (1 + o_p(1)) \max_i \sum_j L_{ij}^0 |s_\lambda(\hat{\sigma}_{ij}) - s_\lambda(\sigma)_{ij}| \\
&= O_p \left(c_1(p) \sqrt{\frac{\log N}{T}} + c_0(p) \left(\lambda^{1-q} + \lambda^{-q} \sqrt{\frac{\log N}{T}} \right) \right)
\end{aligned}$$

Hence we have

$$\begin{aligned}
\|T_{\hat{L}, \lambda}(\hat{\Sigma}) - \Sigma\| &= O_p \left(c_1(p) \sqrt{\frac{\log N}{T}} + c_0(p) \left(\lambda^{1-q} + \lambda^{-q} \sqrt{\frac{\log N}{T}} \right) \right) \\
&= O_p \left(c_1(p) \sqrt{\frac{\log N}{T}} + c_0(p) \left(\frac{\log N}{T} \right)^{\frac{1-q}{2}} \right)
\end{aligned}$$

□

4 Simulations

We demonstrate the network guided estimator and examine its small-sample performance using the following simple simulations. First, we consider the case where the true covariance Σ comes from an AR(1) model. So for $\{(i, j) : i = 1, \dots, N, j = 1, \dots, N\}$, $\sigma_{ij}^2 = \sigma_i \sigma_j \rho_{ij}$ and $\rho_{ij} = \rho^{|i-j|}$. We take $N = 400$.

$$S_{ij} = 3 * \rho^{|i-j|}$$

and assume we observe a matrix $G(l)$ indicating the location of highly correlated pairs $L_{ij}(l) = \mathbf{1}\{\rho_{ij} \geq l\}$. Conditional on $L_{ij} = 1$, we observe $G_{ij} = 1$ with probability p and conditional on $L_{ij} = 0$, $G_{ij} = 1$ with probability q . Hence p, q reflect the probability of missing important locations and including false important locations respectively.

We then generate $T = 100$ independent drws of observations $X \sim N(0, \Sigma)$ and estimate Σ using 1. Sample covariance ; 2. Linear Shrinkage estimator; 3. Nonlinear

Shrinkage estimator; 4. Universal thresholding on the correlation; 5. and Network Guided thresholding estimator. We now compare their performance. It's worth collecting here the parameters that we will adjust in the experiments:

Parameter	Description
ρ	Determines how strong the correlation is and the sparsity of the covariance matrix Σ
l	Observation level, determines how we classify a pair (i, j) as important, i.e., $L_{ij} = 1$.
p	Conditional on $L_{ij} = 1$, the probability of actually observing $G_{ij} = 1$.
q	Conditional on $L_{ij} = 0$, the probability of observing $G_{ij} = 1$
τ	The Threshold level when we apply generalized thresholding operator on σ_{ij} where $G_{ij} = 0$.

Table 1: Description of varying parameters.

Table 2: The estimation error of various estimators in terms of the Frobenius Norm

		Sample Cov	Linear Shrinkage	Nonlinear Shrinkage	Universal Threshold	Network Guided
ρ	Threshold Level					
0.70	0.0	59.43	42.23	41.61	59.43	59.43
	0.1	59.43	42.23	41.61	49.91	49.85
	0.2	59.43	42.23	41.61	41.94	41.67
	0.3	59.43	42.23	41.61	35.66	34.97
	0.4	59.43	42.23	41.61	31.20	29.81
	0.5	59.43	42.23	41.61	28.56	26.18
	0.6	59.43	42.23	41.61	27.54	23.93
	0.7	59.43	42.23	41.61	27.78	22.81
	0.8	59.43	42.23	41.61	28.83	22.48
	0.9	59.43	42.23	41.61	30.33	22.59
	1.0	59.43	42.23	41.61	32.08	22.98

Continued on next page

Table 2: The estimation error of various estimators in terms of the Frobenius Norm

ρ	Threshold Level	Sample	Linear	Nonlinear	Universal	Network
		Cov	Shrinkage	Shrinkage	Threshold	Guided
0.80	0.0	62.54	47.59	46.59	62.54	62.54
	0.1	62.54	47.59	46.59	52.60	52.80
	0.2	62.54	47.59	46.59	44.30	44.52
	0.3	62.54	47.59	46.59	37.89	37.85
	0.4	62.54	47.59	46.59	33.55	32.85
	0.5	62.54	47.59	46.59	31.32	29.48
	0.6	62.54	47.59	46.59	30.95	27.58
	0.7	62.54	47.59	46.59	31.93	26.78
	0.8	62.54	47.59	46.59	33.75	26.71
	0.9	62.54	47.59	46.59	36.02	27.05
	1.0	62.54	47.59	46.59	38.49	27.59
0.90	0.0	63.06	53.18	52.83	63.06	63.06
	0.1	63.06	53.18	52.83	53.98	54.60
	0.2	63.06	53.18	52.83	46.91	47.84
	0.3	63.06	53.18	52.83	42.16	42.89
	0.4	63.06	53.18	52.83	39.78	39.67
	0.5	63.06	53.18	52.83	39.57	37.96
	0.6	63.06	53.18	52.83	41.08	37.43
	0.7	63.06	53.18	52.83	43.74	37.69
	0.8	63.06	53.18	52.83	47.08	38.41
	0.9	63.06	53.18	52.83	50.74	39.33
	1.0	63.06	53.18	52.83	54.57	40.32
0.95	0.0	57.97	52.58	51.77	57.97	57.97
	0.1	57.97	52.58	51.77	51.42	52.21
	0.2	57.97	52.58	51.77	47.65	48.39

Continued on next page

Table 2: The estimation error of various estimators in terms of the Frobenius Norm

ρ	Threshold Level	Sample Cov	Linear Shrinkage	Nonlinear Shrinkage	Universal Threshold	Network Guided
0.99	0.3	57.97	52.58	51.77	46.74	46.40
	0.4	57.97	52.58	51.77	48.35	45.96
	0.5	57.97	52.58	51.77	51.82	46.66
	0.6	57.97	52.58	51.77	56.46	48.04
	0.7	57.97	52.58	51.77	61.73	49.73
	0.8	57.97	52.58	51.77	67.33	51.51
	0.9	57.97	52.58	51.77	73.05	53.24
	1.0	57.97	52.58	51.77	78.77	54.84
	0.0	104.67	115.10	106.32	104.67	104.67
	0.1	104.67	115.10	106.32	114.91	106.31
	0.2	104.67	115.10	106.32	125.35	108.13
	0.3	104.67	115.10	106.32	135.88	110.08
	0.4	104.67	115.10	106.32	146.37	112.10
	0.5	104.67	115.10	106.32	156.67	114.08
	0.6	104.67	115.10	106.32	166.73	115.97
	0.7	104.67	115.10	106.32	176.56	117.77
	0.8	104.67	115.10	106.32	186.19	119.50
	0.9	104.67	115.10	106.32	195.60	121.16
	1.0	104.67	115.10	106.32	204.74	122.77

Table 3: The estimation error of various estimators in terms of the Matrix-1 Norm

ρ	Threshold Level	Sample Cov	Linear Shrinkage	Nonlinear Shrinkage	Universal Threshold	Network Guided
0.70	0.0	63.77	33.42	33.94	63.77	63.77
	0.1	63.77	33.42	33.94	49.17	49.42
	0.2	63.77	33.42	33.94	37.39	38.23
	0.3	63.77	33.42	33.94	28.37	29.20
	0.4	63.77	33.42	33.94	22.21	22.31
	0.5	63.77	33.42	33.94	18.53	17.86
	0.6	63.77	33.42	33.94	16.18	14.59
	0.7	63.77	33.42	33.94	14.59	12.45
	0.8	63.77	33.42	33.94	13.41	11.03
	0.9	63.77	33.42	33.94	12.65	10.22
	1.0	63.77	33.42	33.94	12.56	9.85
0.80	0.0	73.98	44.02	42.83	73.98	73.98
	0.1	73.98	44.02	42.83	58.63	59.50
	0.2	73.98	44.02	42.83	45.86	47.60
	0.3	73.98	44.02	42.83	35.31	37.91
	0.4	73.98	44.02	42.83	26.69	30.04
	0.5	73.98	44.02	42.83	23.68	24.71
	0.6	73.98	44.02	42.83	21.47	20.60
	0.7	73.98	44.02	42.83	20.10	17.40
	0.8	73.98	44.02	42.83	19.66	15.48
	0.9	73.98	44.02	42.83	19.47	14.92
	1.0	73.98	44.02	42.83	19.72	14.85
0.90	0.0	69.65	60.41	60.57	69.65	69.65
	0.1	69.65	60.41	60.57	57.72	58.98
	0.2	69.65	60.41	60.57	50.38	49.66

Continued on next page

Table 3: The estimation error of various estimators in terms of the Matrix-1 Norm

ρ	Threshold Level	Sample Cov	Linear Shrinkage	Nonlinear Shrinkage	Universal Threshold	Network Guided
0.95	0.3	69.65	60.41	60.57	46.30	42.43
	0.4	69.65	60.41	60.57	43.35	38.31
	0.5	69.65	60.41	60.57	41.31	35.10
	0.6	69.65	60.41	60.57	39.91	32.66
	0.7	69.65	60.41	60.57	39.22	31.48
	0.8	69.65	60.41	60.57	39.73	30.93
	0.9	69.65	60.41	60.57	40.38	30.56
	1.0	69.65	60.41	60.57	41.17	30.41
	0.0	95.18	94.47	92.88	95.18	95.18
	0.1	95.18	94.47	92.88	88.39	85.83
0.99	0.2	95.18	94.47	92.88	81.86	76.60
	0.3	95.18	94.47	92.88	76.00	68.02
	0.4	95.18	94.47	92.88	71.17	60.34
	0.5	95.18	94.47	92.88	67.30	57.41
	0.6	95.18	94.47	92.88	69.93	56.36
	0.7	95.18	94.47	92.88	73.19	56.04
	0.8	95.18	94.47	92.88	76.37	55.70
	0.9	95.18	94.47	92.88	79.61	55.75
	1.0	95.18	94.47	92.88	82.59	56.93
	0.0	49.73	43.91	50.83	49.73	49.73
	0.1	49.73	43.91	50.83	46.98	48.49
	0.2	49.73	43.91	50.83	58.00	47.60
	0.3	49.73	43.91	50.83	70.41	50.18
	0.4	49.73	43.91	50.83	82.87	54.76
	0.5	49.73	43.91	50.83	94.81	58.76

Continued on next page

Table 3: The estimation error of various estimators in terms of the Matrix-1 Norm

ρ	Threshold Level	Sample Cov	Linear Shrinkage	Nonlinear Shrinkage	Universal Threshold	Network Guided
	0.6	49.73	43.91	50.83	105.96	61.89
	0.7	49.73	43.91	50.83	115.92	63.92
	0.8	49.73	43.91	50.83	128.73	66.98
	0.9	49.73	43.91	50.83	141.60	69.59
	1.0	49.73	43.91	50.83	154.46	71.30

In Table 2, we show the general performance of these estimators when we simulate using different ρ and thresholding level τ . Here we have taken the thresholding operator to be soft thresholding. It can be seen that generally speaking, when the covariance matrix becomes denser, linear, nonlinear shrinkage estimators and the sample covariance estimator become superior to

Then we consider simulations with varying observation levels l . In Figure 1 when we set observation level equal to 0, the network guided estimator will be the same as the sample covariance estimator, on the other extreme, when observation level is set to 1, the network guided estimator is equivalent to universal thresholding. In between these cases, when we have information about the locations of the important pairs, we have a range where the estimation error is lowered.

Table 4: The estimation error of the Network Guided Estimator with varying probabilities p, q that determine how G is generated.

q	0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
p											
0.0	47.09	48.44	49.91	51.19	52.43	53.68	54.89	56.14	57.37	58.50	59.61
0.1	46.69	48.00	49.51	50.80	52.11	53.36	54.62	55.92	57.09	58.09	59.32
0.2	46.35	47.81	49.19	50.38	51.73	53.10	54.43	55.47	56.76	57.80	58.99

Continued on next page

Table 4: The estimation error of the Network Guided Estimator with varying probabilities p, q that determine how G is generated.

q	0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
p											
0.3	45.91	47.40	48.82	50.12	51.50	52.80	53.88	55.10	56.39	57.47	58.79
0.4	45.53	47.08	48.40	49.61	51.01	52.37	53.61	54.72	56.00	57.14	58.45
0.5	45.15	46.74	48.00	49.46	50.71	52.16	53.33	54.60	55.77	57.00	57.99
0.6	44.78	46.14	47.79	48.98	50.38	51.47	52.93	54.10	55.32	56.59	57.85
0.7	44.37	45.94	47.10	48.96	50.10	51.29	52.52	53.96	55.03	56.23	57.44
0.8	44.06	45.58	46.85	48.36	49.66	50.87	52.15	53.62	54.81	55.90	57.22
0.9	43.51	45.06	46.28	48.02	49.33	50.56	51.99	53.27	54.46	55.69	56.82
1.0	43.14	44.71	46.08	47.52	48.87	50.15	51.69	52.80	54.13	55.29	56.53

In Table 4, we have when $p = q = 0$ the estimation error of the universal thresholding estimator, and $p = q = 1$ the sample covariance estimation error. As we can see, as long as q is not large, the estimation error will be smaller when we have a higher probability p of observing the true large elements. It should be noted that q in fact cannot be very large, given that the whole matrix is sparse.

5 Empirical Studies

5.1 Global Minimum Variance Portfolio

We apply the Network Guided Estimator to a portfolio management similar to Ledoit and Wolf, 2004b. We collect daily return data on SP500 stock from 2004 to 2019 from CRSP, together with daily data on Fama-French 3 factors and the risk free rate.

Assume that the excess returns follow the following factor model

$$Y_{it} = B_i' F_t + \varepsilon_{it}$$

and we assume that $\Sigma = [E\varepsilon_i \varepsilon_j]_{1 \leq i, j \leq N}$ is sparse.

We do a rolling window analysis, each window consists of an estimation period

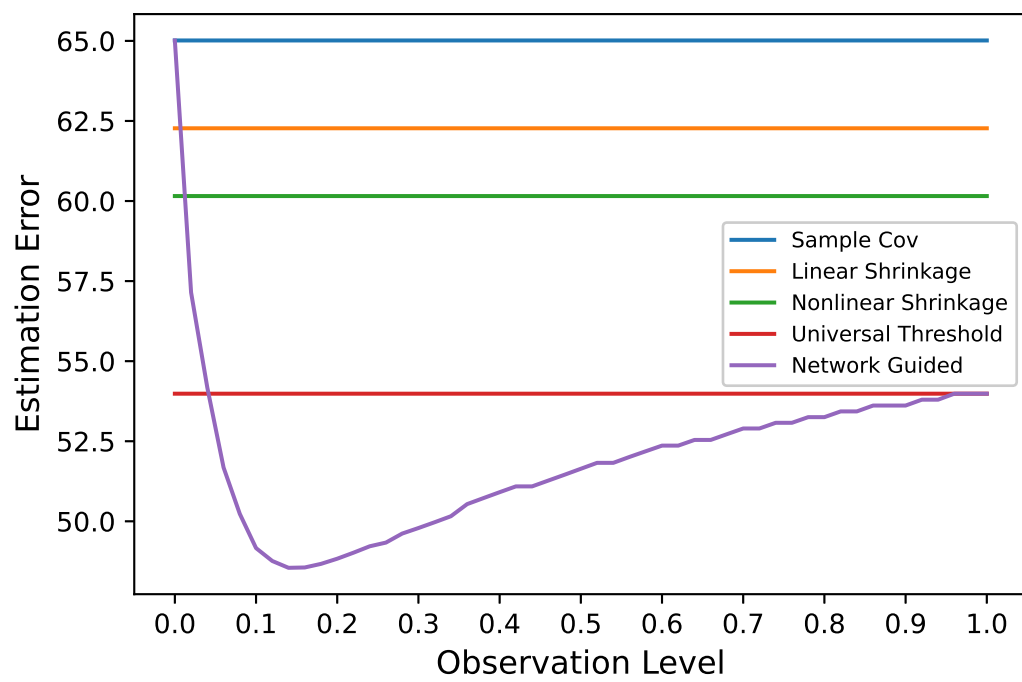


Figure 1: The estimation error against the observation level

of 252 days and a testing period of 21 days. In the estimation period, we estimate the factor loadings by linear time series regression of excess return Y_{it} on F_t , hence allowing the betas to vary over time, and find the de-factored excess return by

$$\hat{\varepsilon}_{it} = Y_{it} - \hat{B}_i' F_t$$

and in order to estimate the covariance matrix of $Y = (Y_1, \dots, Y_N)'$, we have, under the assumption that ε 's are independent of F_t ,

$$\Sigma_Y = B \Sigma_F B' + \Sigma_\varepsilon$$

We replace the factor covariance component by $\hat{B} \hat{\Sigma}_F \hat{B}'$, where $\hat{\Sigma}_F$ is the sample covariance of factors in that period, and we estimate Σ_ε by the Network Guided Estimator applied to $\hat{\Sigma}_\varepsilon = \frac{1}{T} \sum_t \hat{\varepsilon}_t \hat{\varepsilon}_t'$.

In order to apply the Network Guided Estimator, we consider two G matrix that comes from analysts co-coverage: 1. IBES: 2. Dow Jones:

In order to keep G sparse and mitigate the noisy observations, we set for IBES, $G_{ij} = 1$ if firms (i, j) are mentioned more than 18 times, and for Dow Jones data, $G_{ij} = 1$ if firms (i, j) are co-mentioned more than 100 times. With this choice, the total number of links is around 1% of the whole network matrix. We select the thresholding parameter using cross-validation with the constraint that the resulting estimate is positive definite. When the thresholding level becomes higher, the resulting estimate becomes more sparse, in the limit, it'll be a diagonal matrix and p.d., so we keep the thresholding level above the minimum level for the estimate to be p.d.

After using the data from estimation period to estimate \hat{B} , $\hat{\Sigma}_F$, $\hat{\Sigma}_\varepsilon$, we construct

$$\hat{\Sigma}_Y = \hat{B} \hat{\Sigma}_F \hat{B}' + \hat{\Sigma}_\varepsilon$$

and construct the *global minimum variance* portfolio where weights are given by

$$w = \frac{\hat{\Sigma}_Y \mathbf{1}}{\mathbf{1}' \hat{\Sigma}_Y \mathbf{1}}$$

where $\mathbf{1}$ is a conforming vector of ones. We collect the portfolio return over the next 21-day testing period. This concludes one of the rolling windows. Then we move forward

21 days and repeat this exercise. Using 2004- 2014 daily data, we can construct a daily portfolio return from 2005 to 2015, where the portfolio is rebalanced every 21 days. We compute the holding period return of this portfolio and its standard deviation. In [autoreft:4](#) we show the result together with mean and Sharpe ratio and compare it with global minimum variance portfolio constructed using linear shrinkage and universal thresholding. It's worth mentioning that given we are comparing global minimum variance portfolio, the standard deviation is the relevant indicator of performance.

6 Data

We consider daily returns of *S&P* 500 stocks for our application. All the stock market related data are from the Center for Research in Security Prices (CRSP). Daily factor returns are obtained from Kenneth French's website.

6.1 News Implied Network

The news data are obtained from RavenPack Equity files Dow Jones Edition for the period January 2004 to December 2015. This comprehensive news dataset combines relevant content from multiple sources, including Dow Jones Newswires, Wall Street Journal, and Barron's MarketWatch, which produce the most actively monitored streams of news articles in the financial system. Each unique news story (identified by a unique story ID) tags the companies mentioned in the news by their unique and permanent entity identifier codes (RP_ENTITY_ID), by which we link to stock identifier TICKER and PERMNO.

As as Ge and O. B. Linton, [2021](#), we identify links by news co-mentioning. That is, if a piece of business news reports two companies together, they share a link. We do not consider news that co-mention more than two companies since although news they may carry potential information about links, they provide noisier information. We also remove news with topics including analyst recommendations, rating changes, and index movements as these types of news might stack multiple companies together when they actually do not have real links. ?? provides descriptive statistics for RavenPack Equity files Dow Jones Edition dataset during the sample period. Since our comprehensive news dataset combines several sources, given a similar length of sample period, the

number of unique news stories is more than ten times larger than that from Scherbina and Schlusche, 2015 and more than eight hundred times than that from Schwenkler and Zheng, 2019. For link identification purposes, we only use sample news (1) are not about topics mentioned above (2) tag S&P 500 companies and (3) mention exactly two companies, which is a subsample of 1,637,256 unique news stories.

6.2 IBES Analyst Coverage Network

We use the Institutional Brokers Estimate System (IBES) detail history files to construct the analyst co-coverage-based adjacency matrix. For each year in the sample, we consider a stock is covered by an analyst if the analyst issues at least one FY1 or FY2 earnings forecast for the stock during the year. And we consider two stocks as linked if there are common analysts during the year, weighted by the number of common analysts.

7 Conclusion and Further Works

This paper considers the problem of incorporating ever-increasing auxiliary data from machine learning techniques such as textual analysis into the estimation of large covariance matrices. This current version is preliminary with ongoing research on the following applications.

Firstly, we are applying the covariance estimation technique on portfolio construction, following the problem considered in Ledoit and Wolf, 2004b and Ledoit and Wolf, 2017, where the estimation of the sparse covariance matrices are vital for constructing the minimum-variance portfolio.

Secondly, the method can be applied to study spatial-APT under large N case. kou2018asset finds that common risk factors are insufficient to capture all the significant inter-dependencies in asset returns, and local interactions are also important. Spatial-APT and spatial CAPM type of models have not been popular in large N case since the measure of contiguity is challenging. Our method can uncover contemporaneously correlated entities by combining market-based information and auxiliary network information, thus providing a natural contiguity measure. Relying solely on either statistical methods or external network information is not as desirable as the

links identified by the former are hard to interpret and the external network may miss some important links.

Thirdly, we are expanding the set of auxiliary networks beyond the Hoberg’s network as well as applying the technique on larger datasets. We have collected IBES analysts cocoverage network and are constructing new network based on firms’ characteristics. The flexibility of the methods allows us many potential improvements.

References

- Bickel, Peter J. and Elizaveta Levina (2008a). “Covariance Regularization by Thresholding”. In: *The Annals of Statistics* 6, pp. 2577–2604 (cit. on pp. 3, 5, 6).
- (2008b). “Regularized Estimation of Large Covariance Matrices”. In: *The Annals of Statistics* 1, pp. 199–227 (cit. on p. 2).
- Brownlees, Christian, Gumundur Stefán Gumundsson, and Gábor Lugosi (2020). “Community detection in partial correlation network models”. In: *Journal of Business & Economic Statistics*, pp. 1–11 (cit. on p. 4).
- Cai, Tony and Weidong Liu (2011). “Adaptive Thresholding for Sparse Covariance Matrix Estimation”. In: *Journal of the American Statistical Association* 494, pp. 672–684 (cit. on pp. 3, 5, 7).
- Chen, Jia, Degui Li, and Oliver Linton (2019). “A new semiparametric estimation approach for large dynamic covariance matrices with multiple conditioning variables”. In: *Journal of Econometrics* 212.1, pp. 155–176 (cit. on p. 5).
- Connor, Gregory, Matthias Hagmann, and Oliver Linton (2012). “Efficient Semiparametric Estimation of the Fama–French Model and Extensions”. In: *Econometrica* 2, pp. 713–754. ISSN: 1468-0262. DOI: [10.3982/ECTA7432](https://doi.org/10.3982/ECTA7432) (cit. on p. 6).
- Fan, Jianqing, Alex Furger, and Dacheng Xiu (Sept. 15, 2016a). “Incorporating Global Industrial Classification Standard Into Portfolio Allocation: A Simple Factor-Based Large Covariance Matrix Estimator With High-Frequency Data”. In: *Journal of Business & Economic Statistics*. ISSN: 0735-0015. URL: <https://www.tandfonline.com/doi/pdf/10.1080/07350015.2015.1052458?casatoken=FtSi9R7i-SkAAAAA:VJxD-t09DDtdi7QcU4Adch4gVpy6SXPZ0N9zWAL0yXC8aAm7H9HT7NieTYfbdn210pcJP5boOA> (cit. on p. 5).

- Fan, Jianqing, Alex Furger, and Dacheng Xiu (2016b). "Incorporating global industrial classification standard into portfolio allocation: A simple factor-based large covariance matrix estimator with high-frequency data". In: *Journal of Business & Economic Statistics* 34.4, pp. 489–503 (cit. on p. 4).
- Fan, Jianqing and Runze Li (2001). "Variable Selection via Nonconcave Penalized Likelihood and Its Oracle Properties". In: *Journal of the American statistical Association* 456, pp. 1348–1360 (cit. on p. 5).
- Fan, Jianqing, Yuan Liao, and Han Liu (Apr. 12, 2015). *An Overview on the Estimation of Large Covariance and Precision Matrices*. arXiv: 1504.02995 [stat]. URL: <http://arxiv.org/abs/1504.02995> (cit. on pp. 3, 5, 6).
- Fan, Jianqing, Yuan Liao, and Martina Mincheva (2013). "Large covariance estimation by thresholding principal orthogonal complements". In: *Journal of the Royal Statistical Society. Series B, Statistical methodology* 75.4 (cit. on p. 3).
- Fan, Jianqing, Yuan Liao, and Weichen Wang (Feb. 2016). "Projected Principal Component Analysis in Factor Models". In: *Annals of Statistics* 1, pp. 219–254. ISSN: 0090-5364, 2168-8966. DOI: 10.1214/15-AOS1364 (cit. on p. 6).
- Ge, Shuyi and Oliver B Linton (2021). "News-Implied Linkages and Local Dependency in the Equity Market". In: *Available at SSRN 3827902* (cit. on pp. 3, 4, 19).
- Gray, Harry et al. (Sept. 21, 2018). *Shrinkage Estimation of Large Covariance Matrices Using Multiple Shrinkage Targets*. arXiv: 1809.08024 [stat]. URL: <http://arxiv.org/abs/1809.08024> (cit. on p. 6).
- Hoberg, Gerard and Gordon Phillips (2016). "Text-based network industries and endogenous product differentiation". In: *Journal of Political Economy* 124.5, pp. 1423–1465 (cit. on p. 3).
- Israelsen, Ryan D (2016). "Does common analyst coverage explain excess comovement?" In: *Journal of Financial and Quantitative Analysis*, pp. 1193–1229 (cit. on p. 4).
- Kaustia, Markku and Ville Rantala (2013). "Common analyst-based method for defining peer firms". In: *Available at SSRN* (cit. on p. 3).
- Lancewicki, Tomer and Mayer Aladjem (Dec. 2014). "Multi-Target Shrinkage Estimation for Covariance Matrices". In: *IEEE Transactions on Signal Processing* 24, pp. 6380–6390. ISSN: 1941-0476. DOI: 10.1109/TSP.2014.2364784 (cit. on p. 6).

- Ledoit, Olivier and Michael Wolf (Feb. 1, 2004a). “A Well-Conditioned Estimator for Large-Dimensional Covariance Matrices”. In: *Journal of Multivariate Analysis* 2, pp. 365–411. ISSN: 0047-259X. DOI: [10.1016/S0047-259X\(03\)00096-4](https://doi.org/10.1016/S0047-259X(03)00096-4) (cit. on p. 3).
- (July 31, 2004b). “Honey, I Shrunk the Sample Covariance Matrix”. In: *The Journal of Portfolio Management* 4, pp. 110–119. ISSN: 0095-4918, 2168-8656. DOI: [10.3905/jpm.2004.110](https://doi.org/10.3905/jpm.2004.110) (cit. on pp. 3, 5, 16, 20).
 - (Apr. 2012). “Nonlinear Shrinkage Estimation of Large-Dimensional Covariance Matrices”. In: *Annals of Statistics* 2, pp. 1024–1060. ISSN: 0090-5364, 2168-8966. DOI: [10.1214/12-AOS989](https://doi.org/10.1214/12-AOS989) (cit. on pp. 3, 6).
 - (Dec. 1, 2017). “Nonlinear Shrinkage of the Covariance Matrix for Portfolio Selection: Markowitz Meets Goldilocks”. In: *The Review of Financial Studies* 12, pp. 4349–4388. ISSN: 0893-9454, 1465-7368. DOI: [10.1093/rfs/hhx052](https://doi.org/10.1093/rfs/hhx052) (cit. on pp. 3, 6, 20).
- Pourahmadi, Mohsen (2013). *High-Dimensional Covariance Estimation: With High-Dimensional Data*. John Wiley & Sons (cit. on p. 5).
- Rothman, Adam J., Elizaveta Levina, and Ji Zhu (2009). “Generalized Thresholding of Large Covariance Matrices”. In: *Journal of the American Statistical Association* 485, pp. 177–186 (cit. on pp. 5, 9).
- Schäfer, Juliane and Korbinian Strimmer (2005). “A Shrinkage Approach to Large-Scale Covariance Matrix Estimation and Implications for Functional Genomics”. In: p. 32 (cit. on p. 6).
- Scherbina, Anna and Bernd Schlusche (2015). “Economic linkages inferred from news stories and the predictability of stock returns”. In: *Available at SSRN 2363436* (cit. on pp. 4, 20).
- Schwenkler, Gustavo and Hannan Zheng (2019). “The network of firms implied by the news”. In: *Available at SSRN 3320859* (cit. on pp. 4, 20).
- Shu, Hai and Bin Nan (June 1, 2019). “Estimation of Large Covariance and Precision Matrices from Temporally Dependent Observations”. In: *The Annals of Statistics* 3. ISSN: 0090-5364. DOI: [10.1214/18-AOS1716](https://doi.org/10.1214/18-AOS1716). arXiv: [1412.5059](https://arxiv.org/abs/1412.5059) (cit. on pp. 3, 7).