

Minimax Estimation of Bandable Precision Matrices

Authored by:

Sahand Negahban
Addison Hu

Abstract

The inverse covariance matrix provides considerable insight for understanding statistical models in the multivariate setting. In particular, when the distribution over variables is assumed to be multivariate normal, the sparsity pattern in the inverse covariance matrix, commonly referred to as the precision matrix, corresponds to the adjacency matrix representation of the Gauss-Markov graph, which encodes conditional independence statements between variables. Minimax results under the spectral norm have previously been established for covariance matrices, both sparse and banded, and for sparse precision matrices. We establish minimax estimation bounds for estimating banded precision matrices under the spectral norm. Our results greatly improve upon the existing bounds; in particular, we find that the minimax rate for estimating banded precision matrices matches that of estimating banded covariance matrices. The key insight in our analysis is that we are able to obtain barely-noisy estimates of k times k subblocks of the precision matrix by inverting slightly wider blocks of the empirical covariance matrix along the diagonal. Our theoretical results are complemented by experiments demonstrating the sharpness of our bounds.

1 Paper Body

Imposing structure is crucial to performing statistical estimation in the high-dimensional regime where the number of observations can be much smaller than the number of parameters. In estimating graphical models, a long line of work has focused on understanding how to impose sparsity on the underlying graph structure. Sparse edge recovery is generally not easy for an arbitrary distribution. However, for Gaussian graphical models, it is well-known that the graphical structure is encoded in the inverse of the covariance matrix $\Sigma^{-1} = \Omega$, commonly referred to as the precision matrix [12, 14, 3]. Therefore, accurate recovery of the precision matrix is paramount to understanding the structure of the graphical model. As a consequence, a great deal of work has focused on

sparse recovery of precision matrices under the multivariate normal assumption [8, 4, 5, 17, 16]. Beyond revealing the graph structure, the precision matrix also turns out to be highly useful in a variety of applications, including portfolio optimization, speech recognition, and genomics [12, 23, 18]. Although there has been a rich literature exploring the sparse precision matrix setting for Gaussian graphical models, less work has emphasized understanding the estimation of precision matrices under additional structural assumptions, with some exceptions for block structured sparsity [10] or bandability [1]. One would hope that extra structure should allow us to obtain more statistically efficient solutions. In this work, we focus on the case of bandable precision matrices, which capture ? Addison graduated from Yale in May 2017. Up-to-date contact information may be found at <http://huisaddison.com/>.

31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

a sense of locality between variables. Bandable matrices arise in a number of time-series contexts and have applications in climatology, spectroscopy, fMRI analysis, and astronomy [9, 20, 15]. For example, in the time-series setting, we may assume that edges between variables X_i , X_j are more likely when i is temporally close to j , as is the case in an auto-regressive process. The precision and covariance matrices corresponding to distributions with this property are referred to as bandable, or tapering. We will discuss the details of this model in the sequel. Past work: Previous work has explored the estimation of both bandable covariance and precision matrices [6, 15]. Closely related work includes the estimation of sparse precision and covariance matrices [3, 17, 4]. Asymptotically-normal entrywise precision estimates as well as minimax rates for operator norm recovery of sparse precision matrices have also been established [16]. A line of work developed concurrently to our own establishes a matching minimax lower bound [13]. When considering an estimation technique, a powerful criterion for evaluating whether the technique performs optimally in terms of convergence rate is minimaxity. Past work has established minimax rates of convergence for sparse covariance matrices, bandable covariance matrices, and sparse precision matrices [7, 6, 4, 17]. The technique for estimating bandable covariance matrices proposed in [6] is shown to achieve the optimal rate of convergence. However, no such theoretical guarantees have been shown for the bandable precision estimator proposed in recent work for estimating sparse and smooth precision matrices that arise from cosmological data [15]. Of note is the fact that the minimax rate of convergence for estimating sparse covariance matrices matches the minimax rate of convergence of estimating sparse precision matrices. In this paper, we introduce an adaptive estimator and show that it achieves the optimal rate of convergence when estimating bandable precision matrices from the banded parameter space (3). We find, satisfyingly, that analogous to the sparse case, in which the minimax rate of convergence enjoys the same rate for both precision and covariance matrices, the minimax rate of convergence for estimating bandable precision matrices matches the minimax rate of convergence for estimating bandable covariance matrices that has been established in the literature [6]. Our contributions: Our goal is to estimate a banded precision matrix based on

n i.i.d. observations. We consider a parameter space of precision matrices Σ with a power law decay structure nearly identical to the bandable covariance matrices considered for covariance matrix estimation [6]. We present a simple-to-implement algorithm for estimating the precision matrix. Furthermore, we show that the algorithm is minimax optimal with respect to the spectral norm. The upper and lower bounds given in Section 3 together imply the following optimal rate of convergence for estimating bandable precision matrices under the spectral norm. Informally, our results show the following bound for recovering a banded precision matrix with bandwidth k . Theorem 1.1 (Informal). The minimax risk for estimating the precision matrix Σ over the class \mathcal{P}_k given in (3) satisfies:

$$\frac{2}{n} (k + \log p)$$

where this bound is achieved by the tapering estimator $\hat{\Sigma}_k$ as defined in Equation (7). An important point to note, which is shown more precisely in the sequel, is that the rate of convergence as compared to sparse precision matrix recovery is improved by a factor of $\min(k \log(p), k^2)$. We establish a minimax upper bound by detailing an algorithm for obtaining an estimator given observations x_1, \dots, x_n and a pre-specified bandwidth k , and studying the resultant estimator's risk properties under the spectral norm. We show that an estimator using our algorithm with the optimal choice of bandwidth attains the minimax rate of convergence with high probability. To establish the optimality of our estimation routine, we derive a minimax lower bound to show that the rate of convergence cannot be improved beyond that of our estimator. The lower bound is established by constructing subparameter spaces of (3) and applying testing arguments through Le Cam's method and Assouad's lemma [22, 6]. To supplement our analysis, we conduct numerical experiments to explore the performance of our estimator in the finite sample setting. The numerical experiments confirm that even in the finite sample case, our proposed estimator exhibits the minimax rate of convergence. \square

The remainder of the paper is organized as follows. In Section 2, we detail the exact model setting and introduce a blockwise inversion technique for precision matrix estimation. In Section 3, theorems establishing the minimaxity of our estimator under the spectral norm are presented. An upper bound on the estimator's risk is given in high probability with the help of a result from set packing. The minimax lower bound is derived by way of a testing argument. Both bounds are accompanied by their proofs. Finally, in Section 4, our estimator is subjected to numerical experiments. Formal proofs of the theorems may be found in the longer version of the paper [11]. Notation: We will now collect notation that will be used throughout the remaining sections. Vectors will be denoted as lower-case x while matrices are upper-case A . The spectral or operator norm of a matrix is defined to be $\|A\| = \sup_{\|x\|=1} \|Ax\|$ while the matrix '1' norm of a symmetric matrix $A \in \mathbb{R}^{p \times p}$ is defined to be $\|A\|_1 = \sum_{i=1}^p \lambda_i(A)$.

2

Background and problem set-up

In this section we present details of our model and the estimation procedure. If one considers observations of the form $x_1, \dots, x_n \in \mathbb{R}^p$ drawn from a distribution with precision matrix Σ^{-1} and zero mean, the goal then is to estimate the unknown matrix Σ^{-1} based on the observations $\{x_i\}_{i=1}^n$. Given a random sample of p -variate observations x_1, \dots, x_n drawn from a multivariate distribution with population covariance $\Sigma = \Sigma^{-1}^{-1}$, our procedure is based on a tapering estimator derived from blockwise estimates for estimating the precision matrix $\Sigma^{-1} = \Sigma^{-1}$. The maximum likelihood estimator of Σ is $n^{-1} \sum_{i=1}^n x_i x_i^T$.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

\bar{x} is the empirical mean of the vectors x_i . We will construct estimators of the precision matrix where x_i along the diagonal, and averaging over the resultant subblocks. $\Sigma^{-1} = \Sigma^{-1}$ by inverting blocks of Σ . Throughout this paper we adhere to the convention that Σ_{ij} refers to the ij th element in a matrix Σ . Consider the parameter space \mathcal{F} , with associated probability measure P , given by: $(\Sigma) \in \mathcal{F} \iff \Sigma = \Sigma^{-1} (M_0, M) = \Sigma : \max_{i,j} \{|\Sigma_{ij}| : -i \leq j \leq k\} \leq M$ for all k , $\Sigma_{ii} \in [M_0, M_0]$ $\forall i$.

(3) where $\lambda_i(\Sigma)$ denotes the i th eigenvalue of Σ , with $\lambda_i \geq \lambda_j$ for all $i \leq j$. We also constrain $\lambda_i \geq 0$, $M \geq 0$, $M_0 \geq 0$. Observe that this parameter space is nearly identical to that given in Equation (3) of [6]. We take on an additional assumption on the minimum eigenvalue of $\Sigma \in \mathcal{F}$, which is used in the technical arguments where the risk of estimating Σ under the spectral norm is bounded in terms of the error of estimating $\Sigma = \Sigma^{-1}$. Observe that the parameter space intuitively dictates that the magnitude of the entries of Σ decays in power law as we move away from the diagonal. As with the parameter space for bandable covariance matrices given in [6], we may understand Σ in (3) as a rate of decay for the precision entries Σ_{ij} as they move away from the diagonal; it can also be understood in terms of the smoothness parameter in nonparametric estimation [19]. As will be discussed in Section 3, the optimal choice of k depends on both n and the decay rate λ .

Estimation procedure

We now detail the algorithm for obtaining minimax estimates for bandable Σ , which is also given as pseudo-code² in Algorithm 1. The algorithm is inspired by the tapering procedure introduced by Cai, Zhang, and Zhou [6] in the case of covariance matrices, with modifications in order to estimate the precision matrix. Estimating

In the pseudo-code, we adhere to the NumPy convention (1) that arrays are zero-indexed, (2) that slicing an array `arr` with the operation `arr[a:b]` includes the element indexed at `a` and excludes the element indexed at `b`, and (3) that if `b` is greater than the length of the array, only elements up to the terminal element are included, with no errors.

3

the precision matrix introduces new difficulties as we do not have direct access to the estimates of elements of the precision matrix. For a given integer

$k, 1 \leq k \leq p$, we construct a tapering estimator as follows. First, we calculate the maximum likelihood estimator for the covariance, as given in Equation (2). Then, for all integers $1 \leq m \leq l \leq p$ and $m \geq 1$, we define the matrices with square blocks of size at most $3m$ along the diagonal: $\Sigma(3m) = (\Sigma_{ij} 1\{l - m \leq i - j \leq l + 2m, l - m \leq j - i \leq l + 2m\})_{p \times p}$

(4)

$\Sigma(3m)$, we replace the nonzero block with its inverse to obtain $\Sigma^*(3m)$. For a given l , we For each $1 \leq m \leq l$ refer to the individual entries of this intermediate matrix as follows: $\Sigma^*(3m) = (\Sigma^*_{ij} 1\{l - m \leq i - j \leq l + 2m, l - m \leq j - i \leq l + 2m\})_{p \times p}$ (5)

(3m)

For each l , we then keep only the central $m \times m$ subblock of $\Sigma^*(3m)$ to obtain the blockwise estimate $(\Sigma^*)_{l,l} : \Sigma^*(m) = (\Sigma^*_{ij} 1\{l - m \leq i - j \leq l + m, l - m \leq j - i \leq l + m\})_{p \times p}$ (6)

ij

Note that this notation allows for $l \leq 0$ and $l + m \leq p$; in each case, this out-of-bounds indexing allows us to cleanly handle corner cases where the subblocks are smaller than $m \times m$. For a given bandwidth k (assume k is divisible by 2), we calculate these blockwise estimates for both $m = k$ and $m = k/2$. Finally, we construct our estimator by averaging over the block matrices: $\hat{\Sigma} = \frac{1}{2} \sum_{l=1}^{p-k/2} (\Sigma^*(k) + \Sigma^*(k/2))$ (7)

$l=1, \dots, p-k/2$

$k/2$

We note that within entries of the diagonal, each entry is effectively the sum of $k/2$ estimates, and as we move from $k/2$ to k from the diagonal, each entry is progressively the sum of one fewer entry. Therefore, within $k/2$ of the diagonal, the entries are not tapered; and from $k/2$ to k of the diagonal, the entries are linearly tapered to zero. The analysis of this estimator makes careful use of this tapering schedule and the fact that our estimator is constructed through the average of block matrices of size at most $k \times k$. 2.2

Implementation details

The naive algorithm performs $O(p + k)$ inversions of square matrices with size at most $3k$. This method can be sped up considerably through an application of the Woodbury matrix identity and the Schur complement relation [21, 2]. Doing so reduces the computational complexity of the algorithm from $O(pk^3)$ to $O(pk^2)$. We discuss the details of modified algorithm and its computational complexity below. $\Sigma(3m)$ and are interested in obtaining $\Sigma^*(3m)$. We observe that the nonzero block Suppose we have $\Sigma^*(l, m+1) \Sigma(3m)$ corresponds to the inverse of the nonzero block of $\Sigma(3m)$, which only differs by one of $\Sigma^*(l, m+1)$

$\Sigma^*(l, m+1)$

$\Sigma(3m)$, the matrix for which the inverse of the nonzero block corresponds row and one column from $\Sigma^*(l, m) \Sigma(3m)$, $\Sigma^*(l, m)$ to Σ^* , which we have already computed. We may understand the movement from $\Sigma^*(l, m)$

$\Sigma^*(l, m)$

$\Sigma^*(l, m)$

(3m)

$\mathbf{I} \in \mathbb{R}^{l \times m+1}$ (to which we already have direct access) and $\mathbf{I} \in \mathbb{R}^{l \times m+1}$ as two rank-1 updates. Let us view (3m) (3m) $\mathbf{I} \in \mathbb{R}^{l \times m+1}$ the nonzero blocks of $\mathbf{I} \in \mathbb{R}^{l \times m+1}$ as the block matrices:

$\mathbf{A} \in \mathbb{R}^{l \times 1}$ $\mathbf{B} \in \mathbb{R}^{1 \times (3m+1)}$ $\mathbf{3m} \times \mathbf{3m}$ $\mathbf{NonZero}(\mathbf{I} \in \mathbb{R}^{l \times m+1}) = \mathbf{B} \in \mathbb{R}^{1 \times (3m+1)}$ $\mathbf{C} \in \mathbb{R}^{(3m+1) \times (3m+1)}$

$\mathbf{I} \in \mathbb{R}^{1 \times (3m+1)}$ $\mathbf{A} \in \mathbb{R}^{l \times 1}$ $\mathbf{B} \in \mathbb{R}^{1 \times (3m+1)}$ $\mathbf{NonZero}(\mathbf{I} \in \mathbb{R}^{l \times m+1}) = \mathbf{B} \in \mathbb{R}^{1 \times (3m+1)}$ $\mathbf{C} \in \mathbb{R}^{(3m+1) \times (3m+1)}$ $\mathbf{3m} \times \mathbf{3m}$, $\mathbf{I} \in \mathbb{R}^{l \times m+1}$, we may trivially compute $\mathbf{C} \in \mathbb{R}^{(3m+1) \times (3m+1)}$ as The Schur complement relation tells us that given $\mathbf{I} \in \mathbb{R}^{l \times m+1}$ follows:

$\mathbf{I} \in \mathbb{R}^{l \times m+1}$ $\mathbf{B} \in \mathbb{R}^{1 \times (3m+1)}$ $\mathbf{C} \in \mathbb{R}^{(3m+1) \times (3m+1)}$ $\mathbf{A} \in \mathbb{R}^{l \times 1}$ $\mathbf{B} = \mathbf{C} \in \mathbb{R}^{(3m+1) \times (3m+1)}$ (8) $\mathbf{I} \in \mathbb{R}^{l \times m+1}$ $\mathbf{A} \in \mathbb{R}^{l \times 1}$ $\mathbf{B} \in \mathbb{R}^{1 \times (3m+1)}$ $\mathbf{C} \in \mathbb{R}^{(3m+1) \times (3m+1)}$

Algorithm 1 Blockwise Inversion Technique $\mathbf{I} \in \mathbb{R}^{l \times m+1}$ function $\mathbf{F} \in \mathbb{R}^{l \times m+1}$ $\mathbf{LOCK-WISE}(\mathbf{I} \in \mathbb{R}^{l \times m+1}, \mathbf{p} \in \mathbb{Z})$ for $\mathbf{l} \in [1, \mathbf{k}]$ do $\mathbf{I} \in \mathbb{R}^{l \times m+1} + \mathbf{B} \in \mathbb{R}^{1 \times (3m+1)}$ $\mathbf{INVERSE}(\mathbf{I} \in \mathbb{R}^{l \times m+1}, \mathbf{k})$ end for for $\mathbf{l} \in [1, \mathbf{bk}/2\mathbf{c}, \mathbf{p}]$ do $\mathbf{I} \in \mathbb{R}^{l \times m+1} + \mathbf{B} \in \mathbb{R}^{1 \times (3m+1)}$ $\mathbf{INVERSE}(\mathbf{I} \in \mathbb{R}^{l \times m+1}, \mathbf{bk}/2\mathbf{c}, \mathbf{l})$ end for return $\mathbf{I} \in \mathbb{R}^{l \times m+1}$ end function $\mathbf{I} \in \mathbb{R}^{l \times m+1}$ function $\mathbf{B} \in \mathbb{R}^{1 \times (3m+1)}$ $\mathbf{INVERSE}(\mathbf{I} \in \mathbb{R}^{l \times m+1}, \mathbf{p})$ Obtain $\mathbf{3m} \times \mathbf{3m}$ block inverse. $\mathbf{s} \in \max\{\mathbf{l} \in \mathbb{Z}, \mathbf{0}\}$ $\mathbf{f} \in \min\{\mathbf{p}, \mathbf{l} + 2\mathbf{m}\}$

$\mathbf{I} \in \mathbb{R}^{l \times m+1}$ $\mathbf{M} \in \mathbb{R}^{[s:\mathbf{f}], [s:\mathbf{f}]}$. Preserve central $\mathbf{m} \times \mathbf{m}$ block of inverse. $\mathbf{s} \in \mathbf{m} + \min\{\mathbf{l} \in \mathbb{Z}, \mathbf{0}\}$ $\mathbf{N} \in \mathbf{M}[\mathbf{s}:\mathbf{s}+\mathbf{m}, \mathbf{s}:\mathbf{s}+\mathbf{m}]$. Restore block inverse to appropriate indices. $\mathbf{s} \in \max\{\mathbf{l} \in \mathbb{Z}, \mathbf{0}\}$ $\mathbf{f} \in \min\{\mathbf{l} + \mathbf{m}, \mathbf{p}\}$ $\mathbf{P}[\mathbf{s}:\mathbf{f}, \mathbf{s}:\mathbf{f}] = \mathbf{N}$ return \mathbf{P} end function

by the Woodbury matrix identity, which gives an efficient algorithm for computing the inverse of a matrix subject to a low-rank (in this case, rank-1) perturbation. This allows us to move from the inverse of a matrix in $\mathbb{R}^{3m \times 3m}$ to the inverse of a matrix in $\mathbb{R}^{(3m+1) \times (3m+1)}$ where a row and column have been removed. A nearly identical argument allows us to move from the $\mathbb{R}^{(3m+1) \times (3m+1)}$ matrix to an $\mathbb{R}^{3m \times 3m}$ matrix where a row and column have been appended, which gives us the desired block $\mathbf{I} \in \mathbb{R}^{l \times m+1}$. of $\mathbf{I} \in \mathbb{R}^{l \times m+1}$ With this modification to the algorithm, we need only compute the inverse of a square matrix of width $2\mathbf{m}$ at the beginning of the routine; thereafter, every subsequent block inverse may be computed through simple rank one matrix updates. 2.3

Complexity details

We now detail the factor of \mathbf{k} improvement in computational complexity provided through the application of the Woodbury matrix identity and the Schur complement relation introduced in Section 2.2. Recall that the naive implementation of Algorithm 1 involves $O(\mathbf{p} + \mathbf{k})$ inversions of square matrices of size at most $3\mathbf{k}$, each of which cost $O(\mathbf{k}^3)$. Therefore, the overall complexity of the naive algorithm is $O(\mathbf{pk}^3)$, as $\mathbf{k} \leq \mathbf{p}$. Now, consider the Woodbury-Schur-improved algorithm. The initial single inversion of a $2\mathbf{k} \times 2\mathbf{k}$ matrix costs $O(\mathbf{k}^3)$. Thereafter, we perform $O(\mathbf{p} + \mathbf{k})$ updates of the form given in Equation (8). These updates simply require vector matrix operations. Therefore, the update complexity on each iteration is $O(\mathbf{k}^2)$. It follows that the overall complexity of the amended algorithm is $O(\mathbf{pk}^2)$.

3

Rate optimality under the spectral norm

Here we present the results that establish the rate optimality of the above estimator under the spectral norm. For symmetric matrices \mathbf{A} , the spectral norm, which corresponds to the largest singular value of \mathbf{A} , coincides with the ‘2 -operator norm. We establish optimality by first deriving an upper bound 5

in high probability using the blockwise inversion estimator defined in Section 2.1. We then give a matching lower bound in expectation by carefully constructing two sets of multivariate normal distributions and then applying Assouad's lemma and Le Cam's method. 3.1

Upper bound under the spectral norm

In this section we derive a risk upper bound for the tapering estimator defined in (7) under the operator norm. We assume the distribution of the ξ is subgaussian; that is, there exists $\gamma \geq 0$ such that:

$\mathbb{E} \xi^2 \leq P - \gamma \mathbb{E} \xi^2 - \gamma t^2 \leq e^{-2} (9)$ for all $t \geq 0$ and $\|v\|_2 = 1$. Let $\mathcal{P} = \mathcal{P}(M_0, M, \gamma)$ denote the set of distributions of ξ that satisfy (3) and (9). $\gamma \geq k$, defined in (7), of the precision matrix Σ_p with $p \geq$ Theorem 3.1. The tapering estimator $\hat{\Sigma}_1$

$n^{-2\gamma+1}$ satisfies:

2

$k + \log p$

?

$\mathbb{E} \|\hat{\Sigma}_1 - \Sigma\|_F^2 \leq C + Ck = O(p^{-15})$

$k \leq n P$

(10)

with $k = o(n)$, $\log p = o(n)$, and a universal constant $C \geq 0$. 1 $\leq \gamma \leq k$ with $k = n^{-2\gamma+1}$ In particular, the estimator $\hat{\Sigma}_1$ satisfies:

2

$2\gamma \log p$

?

$\mathbb{E} \|\hat{\Sigma}_1 - \Sigma\|_F^2 \leq Cn + C = O(p^{-15})$

$k \leq n P$

(11)

1

Given the result in Equation (10), it is easy to show that setting $k = n^{-2\gamma+1}$ yields the optimal rate by balancing the size of the inside-taper and outside-taper terms, which gives Equation (11). The proof of this theorem, which is given in the supplementary material, relies on the fact that when we invert a $3k \times 3k$ block, the difference between the central $k \times k$ block and the corresponding $k \times k$ block which would have been obtained by inverting the full matrix has a negligible contribution to the risk. As a result, we are able to take concentration bounds on the operator norm of subgaussian matrices, customarily used for bounding the norm of the difference of covariance matrices, and apply them instead to differences of precision matrices to obtain our result. The key insight is that we can relate the spectral norm of a $k \times k$ subblock produced by our estimator to the spectral norm of the corresponding $k \times k$ subblock of the covariance matrix, which allows us to apply concentration bounds from classical random matrix theory. Moreover, it turns out that if we apply the tapering schedule induced by the construction of our estimator to the population parameter $\Sigma = F$, we may express the tapered population Σ as a sum of block matrices in exactly the same way that our estimator is expressed as a sum of block matrices. In particular, the tapering schedule is presented next. Suppose

a population precision matrix $\Sigma \in \mathbb{R}^{p \times p}$. Then, we denote the tapered version of Σ by Σ_A , and construct: $\Sigma_A = (\Sigma_{ij} + v_{ij})_{i,j \in [p]}$ $\Sigma_B = (\Sigma_{ij} + (1 - v_{ij}))_{i,j \in [p]}$ where the tapering coefficients are given by: $v_{ij} = \frac{1}{2} \min\left(\frac{|i-j|}{k}, 1\right)$

for $|i-j| \leq k$ for $|i-j| > k$ for $|i-j| \leq k$

We then handle the risk of estimating the inside-taper Σ_A and the risk of estimating the outside-taper Σ_B separately. Because our estimator and the population parameter are both averages over $k \times k$ block matrices along the diagonal, we may then take a union bound over the high probability bounds on the spectral norm deviation for the $k \times k$ subblocks to obtain a high probability bound on the risk of our estimator. We refer the reader to the longer version of the paper for further details [11].

3.2

Lower bound under the spectral norm

In Section 3.1, we established Theorem 3.1, which states that our estimator achieves the rate of $n^{-2\gamma+1}$ convergence under the spectral norm by using the optimal choice of $k = n^{2\gamma+1}$. Next we demonstrate a matching lower bound, which implies that the upper bound established in Equation (11) is tight up to constant factors. Specifically, for the estimation of precision matrices in the parameter space given by Equation (3), the following minimax lower bound holds. Theorem 3.2. The minimax risk for estimating the precision matrix Σ over \mathcal{P} under the operator norm satisfies:

$$n^{-2\gamma} \log p$$

?

$\inf_{\Sigma \in \mathcal{P}} \sup_{\Sigma \in \mathcal{P}} \mathbb{E} \|\hat{\Sigma} - \Sigma\|_{\text{op}}^2 \geq c n^{-2\gamma+1} + c$ (12) $\mathcal{P} \subset \mathcal{P}$ As in many information theoretic lower bounds, we first identify a subset of our parameter space that captures most of the complexity of the full space. We then establish an information theoretic limit on estimating parameters from this subspace, which yields a valid minimax lower bound over the original set. Specifically, for our particular parameter space \mathcal{P} , we identify two subparameter spaces, \mathcal{F}_{11} , \mathcal{F}_{12} . The first, \mathcal{F}_{11} , is a collection of $2k$ matrices with varying levels of density. To this collection, we apply Assouad's lemma obtain a lower bound with rate $n^{-2\gamma+1}$. The second, \mathcal{F}_{12} , is a collection of diagonal matrices, to which we apply Le Cam's method to derive a lower bound with rate $\log n$. The rate given in Theorem 3.2 is therefore a lower bound on minimax rate for estimating the union $(\mathcal{F}_{11} \cup \mathcal{F}_{12}) = \mathcal{F}_1 \subset \mathcal{P}$. The full details of the subparameter space construction and derivation of lower bounds may be found in the full-length version of the paper [11].

4

Experimental results

We implemented the blockwise inversion technique in NumPy and ran simulations on synthetic datasets. Our experiments confirm that even in the finite sample case, the blockwise inversion technique achieves the theoretical rates. In the experiments, we draw observations from a multivariate normal distribution with precision parameter $\Sigma \in \mathbb{R}^{p \times p}$, as defined in (3). Following [6], for given constants γ, ϵ, p , we consider precision matrices $\Sigma = (\Sigma_{ij})_{i,j \in [p]}$ of the form:

1 for $1 \leq i = j \leq p$ $\Sigma_{ij} = (13) \Sigma_{-i \neq j} = 1$ for $1 \leq i \neq j \leq p$ Though the precision matrices considered in our experiments are Toeplitz, our estimator does not take advantage of this knowledge. We choose $\alpha = 0.6$ to ensure that the matrices generated are non-negative definite. 1

In applying the tapering estimator as defined in (7), we choose the bandwidth to be $k = \lfloor \ln n \rfloor + 1$, which gives the optimal rate of convergence, as established in Theorem 3.1. In our experiments, we varied α , n , and p . For our first set of experiments, we allowed α to take on values in $\{0.2, 0.3, 0.4, 0.5\}$, n to take values in $\{250, 500, 750, 1000\}$, and p to take values in $\{100, 200, 300, 400\}$. Each setting was run for five trials, and the averages are plotted with error bars to show variability between experiments. We observe in Figure 1a that the spectral norm error increases linearly as $\log p$ increases, confirming the $\log n$ term in the rate of convergence. Building upon the experimental results from the first set of simulations, we provide an additional sets of trials for the $\alpha = 0.2$, $p = 400$ case, with $n \in \{11000, 3162, 1670\}$. These sample sizes were chosen so that in Figure 1b, there is overlap between the error plots for $\alpha = 0.2$ and the other α regimes. As with Figure 1a, Figure 1b confirms the minimax rate of convergence given in Theorem 2.3.1. Namely, we see that plotting the error with respect to $n^{\frac{1}{2} + \alpha}$ results in linear plots with almost 3

For the $\alpha = 0.2$, $p = 400$ case, we omit the settings where $n \in \{250, 500, 750\}$ from Figure 1b to improve the clarity of the plot.

7

Setting: $n = 1000$

8

0.025

$\alpha = 0.2 \quad \alpha = 0.3 \quad \alpha = 0.4 \quad \alpha = 0.5$

Spectral Norm Error

6

0.020

Mean Spectral Norm

7

$\alpha = 0.2 \quad \alpha = 0.3 \quad \alpha = 0.4 \quad \alpha = 0.5$

Setting: $p = 400$

0.015

5 4

0.010

3 2

0.005

1 0 4.6

4.8

5.0

5.2 5.4 $\log(p)$

5.6

5.8

0.0000.02

6.0

0.04
0.06
0.08?
 $n^{-\frac{1}{2} + \frac{1}{2}}$
0.10
0.12
0.14
 $2^{-\frac{1}{2}}$

- (b) Mean spectral norm error as $n^{-\frac{1}{2} + \frac{1}{2}}$ changes.
(a) Spectral norm error as $\log p$ changes.

Figure 1: Experimental results. Note that the plotted error grows linearly as a function of $\log p$ and $n^{-\frac{1}{2} + \frac{1}{2}}$, respectively, matching the theoretical results; however, the linear relationship is less clear in the $\frac{1}{2} = 0.2$ case, due to the subtle interplay of the error terms. identical slopes. We note that in both plots, there is a small difference in the behavior for the case $\frac{1}{2} = 0.2$. This observation can be attributed to the fact that for such a slow decay of the precision matrix bandwidth, we have a more subtle interplay between the bias and variance terms presented in the theorems above.

5

Discussion

In this paper we have presented minimax upper and lower bounds for estimating banded precision matrices after observing n samples drawn from a p -dimensional subgaussian distribution. Furthermore, we have provided a computationally efficient algorithm that achieves the optimal rate of convergence for estimating a banded precision matrix under the operator norm. Theorems 3.1 and 3.2 together establish that the minimax rate of convergence for estimating precision matrices over the parameter $2^{-\frac{1}{2}}$ space $F^{\frac{1}{2}}$ given in Equation (3) is $n^{-\frac{1}{2} + \frac{1}{2}} + \log n$, where $\frac{1}{2}$ dictates the bandwidth of the precision matrix. The rate achieved in this setting parallels the results established for estimating a banded covariance matrix [6]. As in that result, we observe that different regimes dictate which term dominates in the $1 - 2^{-\frac{1}{2}}$ rate of convergence. In the setting where $\log p$ is of a lower order than $n^{-\frac{1}{2} + \frac{1}{2}}$, the $n^{-\frac{1}{2} + \frac{1}{2}}$ term dominates, and the rate of convergence is determined by the smoothness parameter $\frac{1}{2}$. However, when $1 - \log p$ is much larger than $n^{-\frac{1}{2} + \frac{1}{2}}$, p has a much greater influence on the minimax rate of convergence. Overall, we have shown the performance gains that may be obtained through added structural constraints. An interesting line of future work will be to explore algorithms that uniformly exhibit a smooth transition between fully banded models and sparse models on the precision matrix. Such methods could adapt to the structure and allow for mixtures between banded and sparse precision matrices. Another interesting direction would be in understanding how dependencies between the n observations will influence the error rate of the estimator. Finally, the results presented here apply to the case of subgaussian random variables. Unfortunately, moving away from the Gaussian setting in general breaks the connection between precision matrices and graph structure. Hence, a fruitful line of work will be to also develop methods that can be applied to estimating the banded graphical model

structure with general exponential family observations. Acknowledgements We would like to thank Harry Zhou for stimulating discussions regarding matrix estimation problems. SN acknowledges funding from NSF Grant DMS 1723128.

8

2 References

- [1] P. J. Bickel and Y. R. Gel. Banded regularization of autocovariance matrices in application to parameter estimation and forecasting of time series. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 73(5):711?728, 2011.
- [2] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge University Press, Cambridge, UK, 2004.
- [3] T. T. Cai, W. Liu, and X. Luo. A Constrained L1 Minimization Approach to Sparse Precision Matrix Estimation. *arXiv:1102.2233 [stat]*, February 2011. *arXiv: 1102.2233*.
- [4] T. T. Cai, W. Liu, and H. H. Zhou. Estimating sparse precision matrix: Optimal rates of convergence and adaptive estimation. *Ann. Statist.*, 44(2):455?488, 04 2016.
- [5] T. T. Cai, Z. Ren, H. H. Zhou, et al. Estimating structured high-dimensional covariance and precision matrices: Optimal rates and adaptive estimation. *Electronic Journal of Statistics*, 10(1):1?59, 2016.
- [6] T. T. Cai, C.-H. Zhang, and H. H. Zhou. Optimal rates of convergence for covariance matrix estimation. *The Annals of Statistics*, 38(4):2118?2144, August 2010.
- [7] T. T. Cai and H. H. Zhou. Optimal rates of convergence for sparse covariance matrix estimation. *Ann. Statist.*, 40(5):2389?2420, 10 2012.
- [8] J. Friedman, T. Hastie, and R. Tibshirani. Sparse inverse covariance estimation with the graphical Lasso. *Biostatistics*, 2007.
- [9] K. J. Friston, P. Jezzard, and R. Turner. Analysis of functional mri time-series. *Human brain mapping*, 1(2):153?171, 1994.
- [10] M. J. Hosseini and S.-I. Lee. Learning sparse gaussian graphical models with overlapping blocks. In *Advances in Neural Information Processing Systems*, pages 3808?3816, 2016.
- [11] A. J. Hu and S. N. Negahban. Minimax Estimation of Bandable Precision Matrices. *arXiv*, 2017. *arXiv: 1710.07006v1*.
- [12] S. L. Lauritzen. *Graphical Models*. Oxford Statistical Science Series. Clarendon Press, Oxford, 1996.
- [13] K. Lee and J. Lee. Estimating Large Precision Matrices via Modified Cholesky Decomposition. *arXiv:1707.01143 [stat]*, July 2017. *arXiv: 1707.01143*.
- [14] N. Meinshausen and P. B?hlmann. High-dimensional graphs and variable selection with the Lasso. *Annals of Statistics*, 34:1436?1462, 2006.
- [15] N. Padmanabhan, M. White, H. H. Zhou, and R. O?Connell. Estimating sparse precision matrices. *Monthly Notices of the Royal Astronomical Society*, 460(2):1567?1576, 2016.
- [16] Z. Ren, T. Sun, C.-H. Zhang, and H. H. Zhou. Asymptotic normality and optimalities in estimation of large Gaussian graphical models. *The Annals of Statistics*, 43(3):991?1026, June 2015.
- [17] A. J. Rothman, P. J. Bickel, E. Levina, and J. Zhu. Sparse permutation invariant covariance estimation. *Electronic Journal of Statistics*, 2:494?515, 2008.
- [18] G. Saon and J. T. Chien. Bayesian sensing hidden markov models for speech recognition. In *2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5056?5059, May 2011.
- [19] A. B. Tsybakov.

Introduction to Nonparametric Estimation. Springer Publishing Company, Incorporated, 1st edition, 2008. [20] H. Visser and J. Molenaar. Trend estimation and regression analysis in climatological time series: an application of structural time series models and the kalman filter. *Journal of Climate*, 8(5):969–979, 1995. [21] M. A. Woodbury. Inverting modified matrices. Statistical Research Group, Memo. Rep. no. 42. Princeton University, Princeton, N. J., 1950. [22] B. Yu. Assouad, Fano and Le Cam. In *Festschrift for Lucien Le Cam*, pages 423–435. Springer-Verlag, Berlin, 1997. [23] M. Yuan and Y. Lin. Model selection and estimation in the Gaussian graphical model. *Biometrika*, 94(1):19–35, 2007.