
antGLasso: An Efficient Tensor Graphical Lasso Algorithm

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The class of bigraphical lasso algorithms (and, more broadly, ‘tensor’-graphical
2 lasso algorithms) has been used to estimate dependency structures within matrix
3 and tensor data. However, all current methods to do so take prohibitively long
4 on modestly sized datasets. We present a novel tensor-graphical lasso algorithm
5 that directly estimates the dependency structure, unlike its iterative predecessors.
6 This provides a speedup of multiple orders of magnitude, allowing this class of
7 algorithms to be used on large, real-world datasets.

8 1 Introduction

9 We often make independence assumptions to simplify the modelling of complex data. The strength
10 of an independence assumption exists on a spectrum. One could assume complete independence,
11 independence except on a sparse subset, or independence of samples but not features, among others.

12 For conditional dependencies, there is a natural way to model them in Gaussian data. Two datapoints
13 x, y are *conditionally independent* (with respect to a dataset \mathcal{D}) if knowing one provides no infor-
14 mation about the other that is not already contained in the rest of the dataset: $\mathbb{P}(x|y, \mathcal{D}) = \mathbb{P}(x|\mathcal{D})$.
15 For normally distributed data, conditional dependencies are encoded in the inverse of the covariance
16 matrix (the ‘precision’ matrix). Two datapoints are conditionally dependent on each other if and only
17 if their corresponding element in the precision matrix is zero. If our dataset \mathbf{d} were in the form of a
18 vector, we could then model it as $\mathbf{d} \sim \mathcal{N}(\mathbf{0}, \Psi^{-1})$ for precision matrix Ψ .

19 However, datasets are often not simple vectors - for example, single-cell RNA sequencing data
20 (scRNA-seq) comes in the form of a matrix of gene expression counts whose rows are cells and
21 columns are genes. Video data naturally requires a third-order tensor of pixels to represent it - rows,
22 columns, and frames. To address this, we ‘vectorize’ matrices by stacking the columns vertically.
23 Vectorization of tensors proceeds in much the same way.

24 Our matrix-variate dataset \mathcal{D} can then be represented as $\text{vec}[\mathcal{D}] \sim \mathcal{N}(\mathbf{0}, \Psi^{-1})$. Suppose our matrix
25 were $n \times p$ dimensional - then $\text{vec}[\mathcal{D}]$ has np elements, and its precision matrix has n^2p^2 elements.
26 This poses a problem: the estimation of our precision matrix requires substantially more parameters
27 than we have in our dataset. Since the columns of the matrix have interdependencies and the
28 rows of the matrix have interdependencies, to reduce space we can compute the interdependencies
29 of the elements of the matrix via some deterministic function of the row-wise and column-wise
30 dependencies. That is, for some function ζ , $\text{vec}[\mathcal{D}] \sim \mathcal{N}(\mathbf{0}, \zeta(\Psi_{\text{row}}, \Psi_{\text{col}})^{-1})$.

31 2 Background

32 The Kronecker sum bigraphical lasso (BiGLasso) model was first considered by Kalaitzis et al. [7].
33 BiGLasso is the multi-axis analog to graphical lasso methods [3, 1], which are used to estimate

34 covariance matrices of data of a multivariate Gaussian distribution. The Kronecker sum of two
 35 matrices, $\mathbf{A} \oplus \mathbf{B}$, can be expressed in terms of Kronecker products: $\mathbf{A} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{B}$. When the matrices
 36 \mathbf{A}, \mathbf{B} are adjacency matrices of graphs, the Kronecker sum has the interpretation as the Cartesian
 37 product of those graphs[16]. This sum is one choice of function ζ to combine the per-axis precision
 38 matrices into the precision matrix of the vectorized dataset, $\text{vec}[\mathbf{D}] \sim \mathcal{N}(\mathbf{0}, (\Psi_{\text{row}} \oplus \Psi_{\text{col}})^{-1})$.

39 Other choices for ζ have been considered, such as using the Kronecker product[18, 23], or the square
 40 of the Kronecker sum[20]. Each method has its strengths; the benefits of a Kronecker sum structure
 41 are its interpretability as a graph product, stronger sparsity, and its allowance of inter-task transfer[7].

42 The original BiGLasso model was very slow to converge to a solution, in large part due to its non-
 43 optimal space complexity of $O(n^2 p^2)$. This prohibited its use on large datasets (measuring in a
 44 couple hundred samples and/or features). Numerous modifications have been made to the algorithm
 45 to improve its speed and achieve an optimal space complexity of $O(n^2 + p^2)$, such as scBiGLasso[10],
 46 TeraLasso[4], and EiGLasso[21]. Of these, TeraLasso is notable in that it generalizes to an arbitrary
 47 amount of axes, i.e. $\zeta(\Psi_1, \dots, \Psi_k) = \Psi_1 \oplus \dots \oplus \Psi_k$. However, all of these algorithms are still
 48 iterative. The convergence time can vary greatly depending on the precise values of the data, and
 49 in general they struggle to handle datasets with a couple thousand samples and/or features. Our
 50 proposed variant, antGLasso (**analytic tensor-graphical lasso**), does not suffer from these issues,
 51 preserves optimal memory complexity, and like TeraLasso can work on datasets with an arbitrary
 52 amount of axes.

53 All algorithms listed rely on a normality assumption. One could manually transform non-normal
 54 data to a Gaussian, if they knew the right transform to use. An alternative method is to use the
 55 Nonparanormal Skeptic[11], an algorithm which directly estimates the Gram matrices of data that
 56 has been transformed to a Gaussian by some unspecified function[10]. Since all previous BiGLasso
 57 algorithms mentioned can be framed in terms of accepting these Gram matrices as input, they can be
 58 made to circumvent the normality assumption.

59 3 Methodology and Results

60 In the following we introduce the main theorems and results of our study. For our implementation,
 61 we used Python Version 3.9.12, Numpy Version 1.22.2[5], Scipy Version 1.7.3[19], Sklearn Version
 62 1.0.2[14], Matplotlib Version 3.5.1[6], and Pandas Version 1.4.2[17, 12]. All tests were run on a
 63 MacBook Pro (13-inch, M1, 2020) with 8GB of RAM, using Jupyter Notebooks[15, 8].

64 3.1 Theorems

65 The general structure of the algorithm is as follows: we first directly estimate the eigenvectors of
 66 the solution. We then use the eigenvectors to diagonalize the covariance matrix of the tensor-variate
 67 distribution, such that we can directly estimate the eigenvalues from the variances of the data. We
 68 follow the notation of Kolda and Bader [9] for tensor operations, and Greenewald, Zhou, and Hero
 69 [4] for BiGLasso variables. \mathbf{K}_b^a is shorthand for the matrix with a on the diagonals and b elsewhere.
 70 Our proofs of the following theorems are given in the supplementary material, and rely on a lemma
 71 by Dahl et al. [2].

72 **Theorem (Eigenvectors of Precision Matrices).** \mathbf{V}_ℓ , the eigenvectors of Ψ_ℓ , are the eigenvectors of
 73 $\mathbf{S}_\ell \circ \mathbf{K}_{m_\ell}^{2m_\ell-1}$, where \mathbf{K}_b^a is the matrix with diagonal elements a and off-diagonal elements b .

74 **Theorem (Eigenvalues of Precision Matrices).** Let \mathbf{a} be the vector of length $\prod_{\ell=1}^K d_\ell$ whose
 75 $(\sum_{\ell=1}^K i_\ell \prod_{j=1}^{\ell-1} d_j)$ th element is $\frac{1}{\text{var}[(\mathcal{Y} \times_1 \mathbf{V}_1^T \times_2 \dots \times_K \mathbf{V}_K^T)_{i_1 \dots i_K}]}$. Then there is a system of linear equa-
 76 tions relating \mathbf{a} to the eigenvalues of (Ψ_1, \dots, Ψ_K) .

77 Unfortunately, the calculation of the eigenvalues requires the whole dataset \mathcal{Y} , whereas the use of the
 78 Nonparanormal Skeptic requires that the algorithm be parametrizable in terms of the Gram matri-
 79 ces. We can circumvent this by introducing a heuristic, $\text{var}[(\mathcal{Y} \times_1 \mathbf{V}_1^T \times_2 \dots \times_K \mathbf{V}_K^T)_{i_1 \dots i_K}] \approx$
 80 $\Delta_{i_\ell} \mathbf{V}_\ell^T \mathbf{S}_\ell \mathbf{V}_\ell \Delta_{i_\ell}^T$, where Δ_{i_ℓ} is a row vector of zeros except for a one at position i_ℓ . We call this
 81 variant of the algorithm ‘Heuristic antGLasso’; further details are available in the supplementary
 82 material.

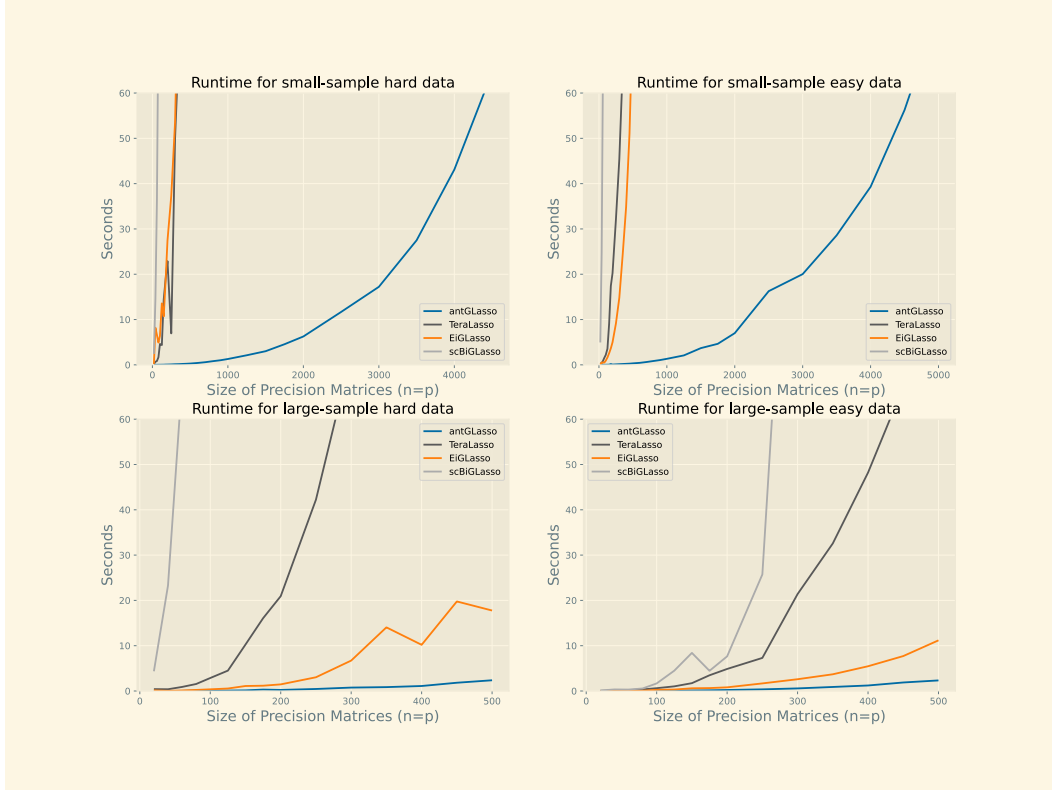


Figure 1: A comparison of runtimes of four BiGLasso algorithms on simulated data. We consider small samples ($s = 1$) and large samples ($s = 100$), as well as data drawn from distributions that tend to cause algorithms to converge slowly (‘hard’) and quickly (‘easy’).

83 3.2 Results on Simulated Data

84 We first tested our algorithm on simulated data. The precision matrices were drawn from an inverse
85 Wishart distribution. The iterative algorithms tended to converge slowly when the degrees of freedom
86 parameter of the distribution was small, and quickly when it was large. Additionally, they converged
87 much slower on small sample data (Figure 1). We can see that antGLasso is by far the fastest algorithm,
88 giving up to two orders of magnitude improvement over the next fastest algorithm (EiGLasso).

89 Not pictured are the results for tensor data, which are in the supplementary material. For three-axis
90 tensors, antGLasso is still markedly faster than TeraLasso in the small sample case - we did not
91 have enough RAM to test the large sample case. For four-axis and larger, the difference between the
92 algorithms becomes negligible - this is because the major runtime constraint becomes the calculation
93 of the Gram matrices, whose complexity grows exponentially with the amount of axes and thus
94 dominates.

95 antGLasso does not perform as well as other algorithms (Figure 2), although this gap drops as the
96 number of samples increase. Interestingly, Heuristic antGLasso performs comparably to the others
97 while preserving the speed of antGLasso. We believe this is because antGLasso directly estimates
98 the variances of the data, which will be inaccurate for small samples. The heuristic also estimates
99 variances, but by using info on the variance across other axes to in essence artificially increase the
100 sample count. The gap between the two variants drops as the number of samples increases.

101 3.3 Temporal Recovery

102 We have shown that on simulated data Heuristic antGLasso performs comparably but much quicker
103 than existing algorithms. This enables the application of BiGLasso methods to datasets in the size of
104 thousands. For real world experiments, we can consider the case of video data. Given the graph of

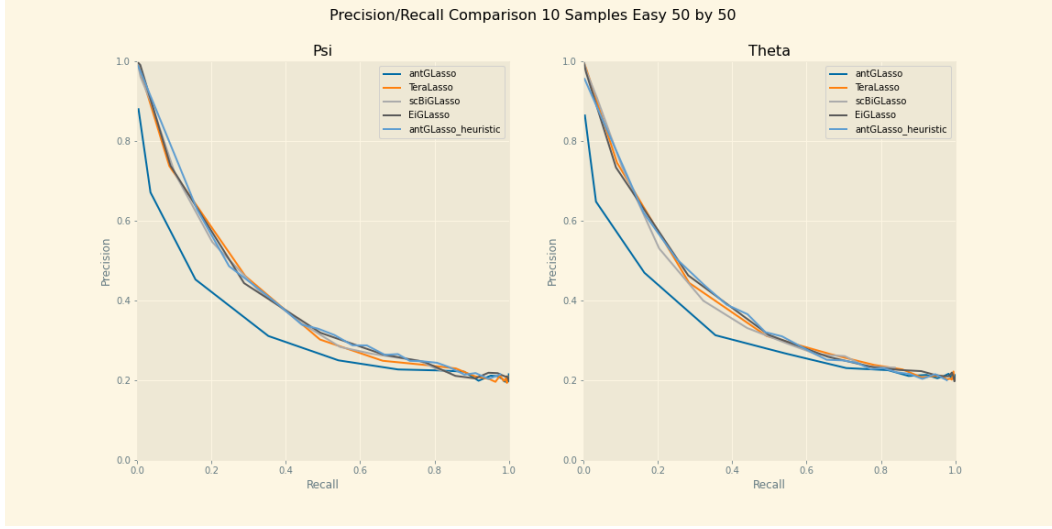


Figure 2: Precision-recall curves for two-dimensional simulated data from a distribution that converged quickly (‘easy’), for 50x50 input matrices.

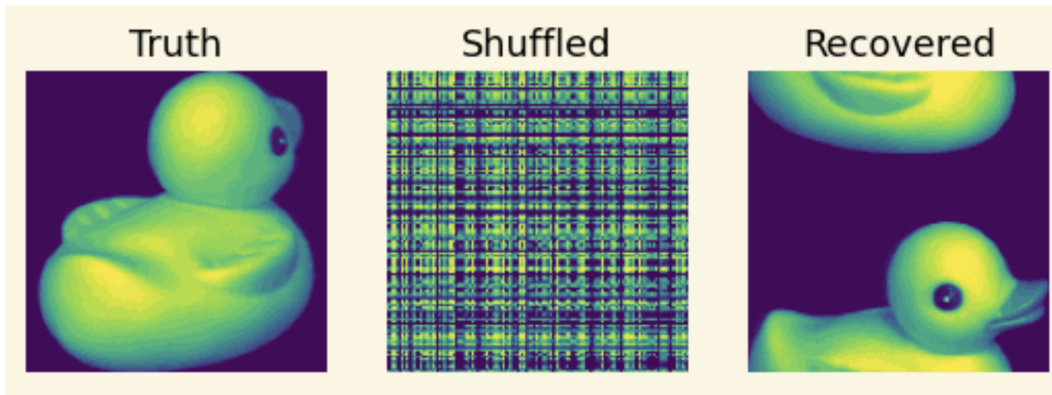


Figure 3: Recovery of a video of a spinning rubber duck. (Left) A frame from the original video, (Middle) a frame from the shuffled video, (Right) frame from the reconstructed video. Not represented in the image is the recovery of the frames of the video, which was perfect.

105 conditional dependencies, we might expect to be able to order the frames, as a frame should only be
 106 conditionally dependant on its nearest neighbors.

107 For example, consider the video of a rubber duck from the COIL-20 dataset[13]. If we scramble
 108 the rows, frames, and columns of the video, we can recover the original video from the outputs of
 109 antGLasso, as seen in Figure 3. The reconstruction is perfect, except for the duck being cut in half.
 110 This is due to the difficulty in figuring out which row to start on.

111 The strong results for video reconstruction are encouraging. However, when we modify the video by
 112 duplicating frames and adding a bit of noise our reconstruction performs much worse. Already by
 113 duplicating each frame 5 times we dip below 10% reconstruction accuracy. This is arguably because
 114 duplicating frames interferes with the conditional dependency graph - each frame has more similar
 115 frames in the rest of the dataset, which then get conditioned out, reducing informativeness. This
 116 suggests that temporal recovery can only be done in highly structured cases with few redundancies.
 117 We verified that this behavior is seen in all BiGLasso algorithms considered in the report, including
 118 antGLasso.

119 In Figure 4, we can see how choice of algorithm as well as duplicate frames and noise affect temporal
 120 recovery. When frames are not duplicated, using the Nonparanormal Skeptic performs the best.

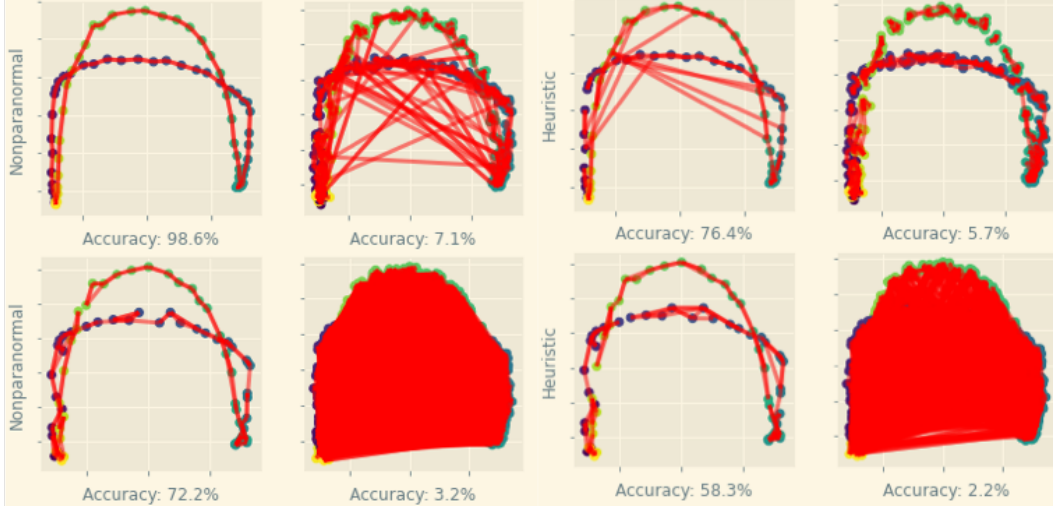


Figure 4: A comparison of temporal recovery on the COIL duck video visualized with the first two principal components. Unlike elsewhere in the paper, we do not treat the video as a three dimensional tensor, but rather we combine the two pixel dimensions so that we can make use of PCA. Lines linking nodes mean that one of the nodes has the other as one of its top two strongest connections according to the precision matrix. The left half of the graph was constructed with antGLasso+Nonparanormal Skeptic, and the right half used Heuristic antGLasso. In each half, the top-left graph shows performance on the unaltered COIL video. The bottom-left shows how performance degrades when we add noise, and the top-right shows what happens when we duplicate each frame 5 times. The bottom-right has each frame duplicated 25 times.

121 However, once frames are duplicated it seems that using Heuristic antGLasso is better. Even though
 122 the heuristic results in a lower accuracy, visually we can see that when it is wrong we can still trace a
 123 sensible progression along the graph. Once the Nonparanormal Skeptic starts going wrong, it starts
 124 connecting completely disparate parts of the graph rather than failing gracefully.

125 3.4 Clustering

126 To test antGLasso’s performance on clustering realistic data, we generated a dataset using Splatter[22].
 127 Splatter generates scRNA-seq data, which does not follow a normal distribution - hence using the
 128 Nonparanormal Skeptic is justified.

129 There are three different interpretations of the output graphs that we can use for clustering - we can
 130 use the outputs directly (the precision matrices), we can invert the outputs (covariance matrices), or
 131 we can consider only the negative values of our precision matrices. Inversion is relatively cheap as
 132 our algorithm estimates the eigenvalues and eigenvectors separately. Values in the precision matrix of
 133 a Gaussian Graphical Model can be interpreted as the ‘negative stiffness’ of a spring connecting two
 134 elements together, motivating the consideration of the negative values of our outputs.

135 In addition to the three interpretations of our outputs, we have three different models to use: Vanilla
 136 antGLasso, Heuristic antGLasso, and Heuristic antGLasso with the Nonparanormal Skeptic. The
 137 results can be seen in Figure 5. Heuristic antGLasso is the best to use, whereas the Nonparanormal
 138 Skeptic does very poorly. We also tested EiGLasso’s performance on the same dataset, and found
 139 similar results.

140 3.5 Regularization

141 The algorithm is based on an analytic solution to the maximum of the unregularized likelihood
 142 function. We do not have an analytic solution for the L1-penalized likelihood function. Instead of
 143 regularizing during the calculation of precision matrices, we first find the unregularized solution, and
 144 then apply a regularization step. By regularizing the solution at the end, rather than simultaneously, it
 145 allows us to efficiently examine the solution across many different penalty strengths. However, it

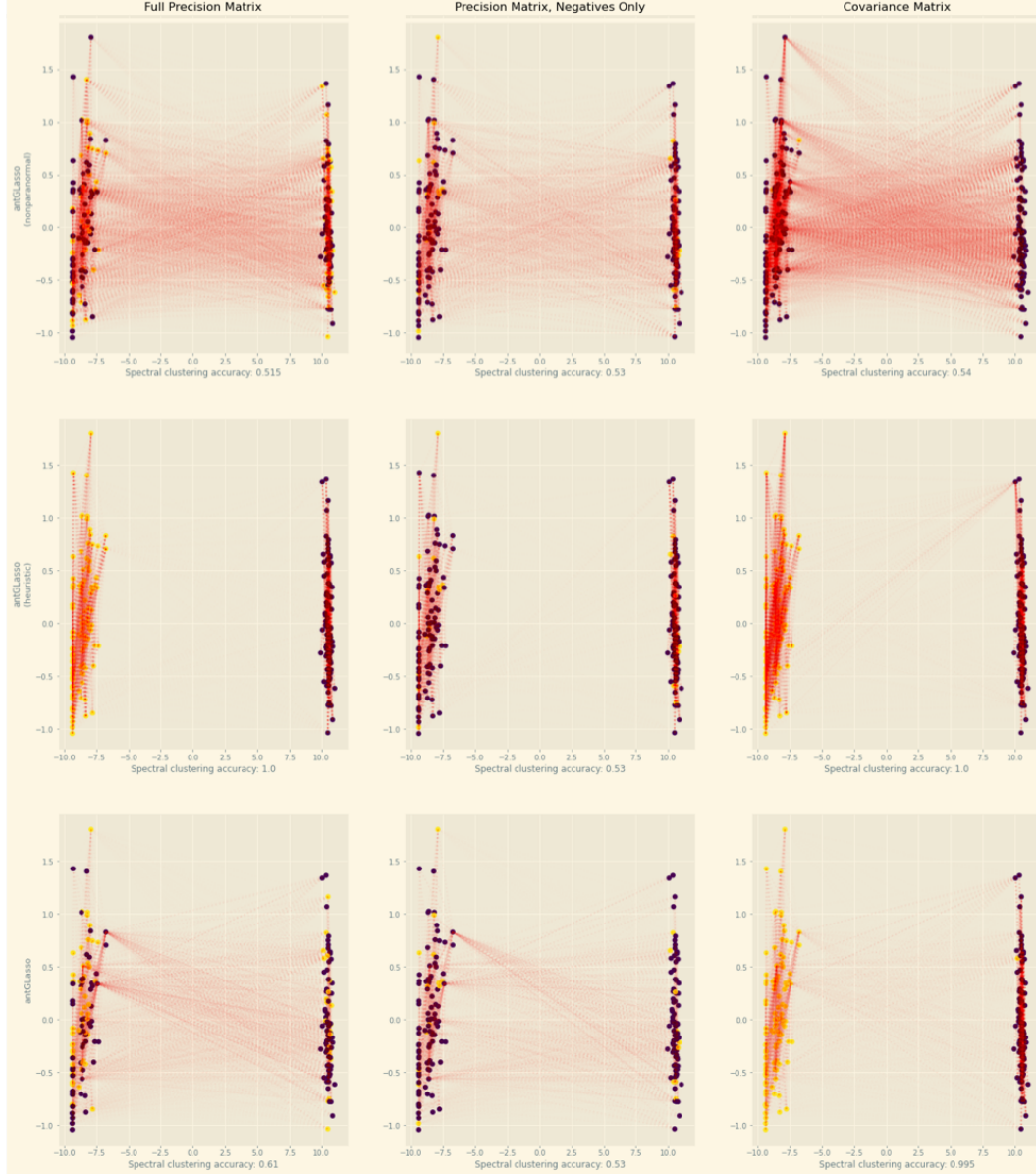


Figure 5: A comparison of different versions of the algorithm (rows) and ways to interpret the results (columns). Only Heuristic antGLasso clusters robustly in all interpretations, as evidenced by the lack of cross-connections. We report the accuracy of performing spectral clustering on these matrices as an numerical measure of clustering acumen.

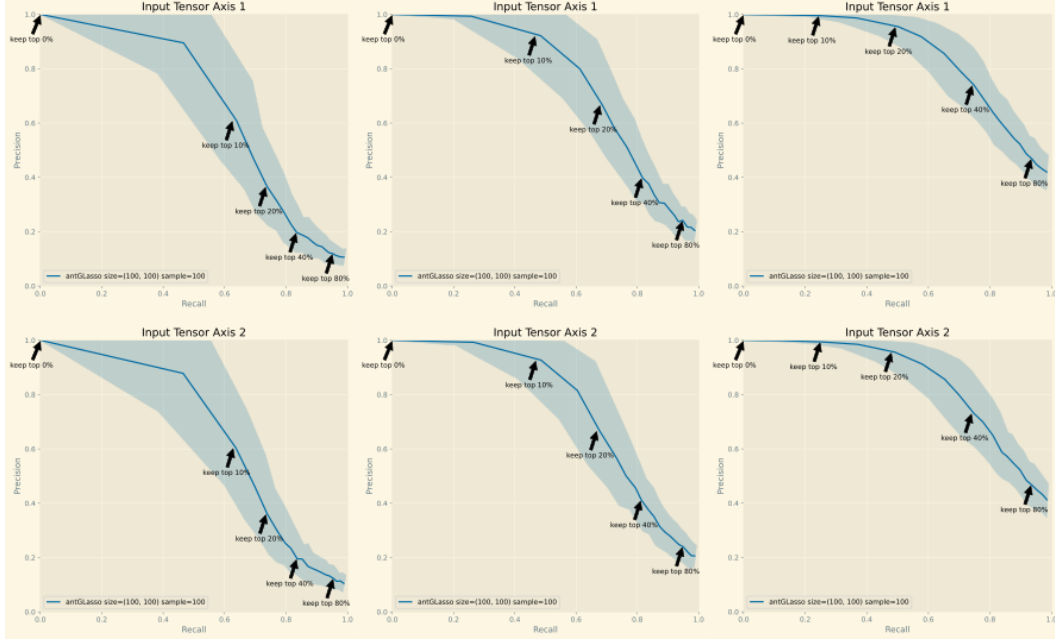


Figure 6: Precision recall curves with hyperparameter values labeled. (Left column) Simulated data with 10% true edges. (Middle column) Simulated data with 20% true edges. (Right column) Simulated data with 40% true edges.

146 does not necessarily converge to the solution to the L1-penalized likelihood function, unlike previous
 147 algorithms. As seen in Figure 2, Heuristic antGLasso still achieves state-of-the-art recovery in spite
 148 of this.

149 If we were to apply Lasso to our unregularized solution, in a similar manner to how BiGLasso and
 150 scBiGLasso apply Lasso at the end of each iteration, it would be equivalent to just thresholding. This
 151 is convenient: not only is thresholding a very fast operation, but it can also be framed in terms of
 152 the percentage of edges to keep. As seen in Figure 6, if we simulate data to have a given sparsity
 153 level, and then use the known ground truth sparsity as the hyperparameter, we achieve a balance
 154 between precision and recall. Thus, if we have an idea of how densely connected our data should be
 155 and the extent to which we wish to prioritize precision over recall, we will know in advance a good
 156 regularization strength.

157 Despite the benefits of this sparsity-based regularizer, we often chose to use problem-specific
 158 regularization regimes. For temporal recovery, a simple technique is to keep only the two highest
 159 values in each row (for the immediate past and immediate future). For clustering, it can make sense
 160 not to regularize, as the additional edges may contain useful information for the clusterer despite not
 161 actually being present in the real data.

162 4 Conclusion

163 We have created a new BiGLasso algorithm, antGLasso, which is much faster than current competitors,
 164 while preserving (in the heuristic case) their accuracy. It can be applied to tensors of arbitrary
 165 dimensions, like TeraLasso, and performs well on simulated data as well as highly-structured real-
 166 world data. It falters when applied to more complicated data, which is a major limitation. However,
 167 this is a systemic problem with BiGLasso algorithms rather than our implementation specifically.

168 For future work, we would like to understand better when the Nonparanormal Skeptic should be used.
 169 We know that the Nonparanormal Skeptic can be useful for clustering, as Li et al. [10] used it to
 170 find interpretable results on real scRNA-seq data. It seems that the Nonparanormal Skeptic performs
 171 poorly when the data is multimodal - in the case of duplicating frames in Figure 4, the duplication

has the effect of creating clear clusters of duplicates), and we believe this is what caused it to perform so poorly just as it did in the clustering tests of Figure 5. Further investigation is needed.

We would also like to investigate how to circumvent the known issues with applying this method to less structured datasets, such as scRNA data, for temporal recovery and other tasks. This may involve changing aspects of the model to better fit the idiosyncrasies of specific types of datasets.

Finally, we believe that, despite our model’s speed, it is still important for further improvements to be made. For datasets of size in the 1000s, the runtime is not significant. However, by extrapolating the (cubic) runtime curves to the largest dataset whose outputs could fit on our computer’s RAM (8GB), we can estimate that antGLasso would take 2 days to run. We hope that the algorithm can reach a point where all problems that fit on a personal computer can run over a coffee break.

Improving speed would likely require a fundamentally different approach, as the majority of the runtime is taken up by a single eigendecomposition (per axis). We considered using other decompositions, but could not find any that used orthonormal matrices, had a sparse core, and behaved nicely under the block-wise trace of the inverse of the Kronecker product of cores. If such a decomposition does exist, and is faster than eigendecomposition, then further speed gains can be made using a structurally similar algorithm to antGLasso.

Despite this, the constraint limiting BiGLasso-style algorithms’ applicability to even larger datasets is likely to be memory. Unfortunately, there is not much room for improving memory usage – modifications could likely be made to remove some intermediate products, but even in the ideal case we would be reducing memory usage by a small constant factor, since Heuristic antGLasso is memory-optimal.

References

- [1] Onureena Banerjee, Laurent El Ghaoui, and Alexandre d’Aspremont. “Model Selection Through Sparse Maximum Likelihood Estimation for Multivariate Gaussian or Binary Data”. In: *Journal of Machine Learning Research* 9.15 (2008), pp. 485–516. URL: <http://jmlr.org/papers/v9/banerjee08a.html>.
- [2] Andy Dahl et al. *Network inference in matrix-variate Gaussian models with non-independent noise*. 2013. DOI: 10.48550/ARXIV.1312.1622. URL: <https://arxiv.org/abs/1312.1622>.
- [3] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *Sparse inverse covariance estimation with the lasso*. 2007. DOI: 10.48550/ARXIV.0708.3517. URL: <https://arxiv.org/abs/0708.3517>.
- [4] Kristjan Greenewald, Shuheng Zhou, and Alfred Hero. *Tensor Graphical Lasso (TeraLasso)*. 2017. DOI: 10.48550/ARXIV.1705.03983. URL: <https://arxiv.org/abs/1705.03983>.
- [5] Charles R. Harris et al. “Array programming with NumPy”. In: *Nature* 585.7825 (Sept. 2020), pp. 357–362. DOI: 10.1038/s41586-020-2649-2. URL: <https://doi.org/10.1038/s41586-020-2649-2>.
- [6] J. D. Hunter. “Matplotlib: A 2D graphics environment”. In: *Computing in Science & Engineering* 9.3 (2007), pp. 90–95. DOI: 10.1109/MCSE.2007.55.
- [7] Alfredo Kalaitzis et al. “The Bigraphical Lasso”. In: *Proceedings of the 30th International Conference on Machine Learning*. Ed. by Sanjoy Dasgupta and David McAllester. Vol. 28. Proceedings of Machine Learning Research 3. Atlanta, Georgia, USA: PMLR, 17–19 Jun 2013, pp. 1229–1237. URL: <https://proceedings.mlr.press/v28/kalaitzis13.html>.
- [8] Thomas Kluyver et al. “Jupyter Notebooks - a publishing format for reproducible computational workflows”. In: *Positioning and Power in Academic Publishing: Players, Agents and Agendas*. Ed. by Fernando Loizides and Birgit Schmidt. Netherlands: IOS Press, 2016, pp. 87–90. URL: <https://eprints.soton.ac.uk/403913/>.
- [9] Tamara G. Kolda and Brett W. Bader. “Tensor Decompositions and Applications”. In: *SIAM Review* 51.3 (Sept. 2009), pp. 455–500. DOI: 10.1137/07070111X.
- [10] Sijia Li et al. “Scalable Bigraphical Lasso: Two-way Sparse Network Inference for Count Data”. In: (Mar. 2022).
- [11] Han Liu et al. *The Nonparanormal SKEPTIC*. 2012. DOI: 10.48550/ARXIV.1206.6488. URL: <https://arxiv.org/abs/1206.6488>.

- [12] Wes McKinney. “Data Structures for Statistical Computing in Python”. In: *Proceedings of the 9th Python in Science Conference*. Ed. by Stéfan van der Walt and Jarrod Millman. 2010, pp. 56–61. DOI: 10.25080/Majora-92bf1922-00a.
- [13] Sameer A. Nene, Shree K. Nayar, and Hiroshi Murase. *Columbia Object Image Library (COIL-20)*. Tech. rep. 1996.
- [14] F. Pedregosa et al. “Scikit-learn: Machine Learning in Python”. In: *Journal of Machine Learning Research* 12 (2011), pp. 2825–2830.
- [15] Fernando Pérez and Brian E. Granger. “IPython: a System for Interactive Scientific Computing”. In: *Computing in Science and Engineering* 9.3 (May 2007), pp. 21–29. ISSN: 1521-9615. DOI: 10.1109/MCSE.2007.53. URL: <https://ipython.org>.
- [16] Gert Sabidussi. “Graph Multiplication.” In: *Mathematische Zeitschrift* 72 (1959/60), pp. 446–457. URL: <http://eudml.org/doc/183624>.
- [17] The pandas development team. *pandas-dev/pandas: Pandas*. Version latest. Feb. 2020. DOI: 10.5281/zenodo.3509134. URL: <https://doi.org/10.5281/zenodo.3509134>.
- [18] Theodoros Tsiligkaridis and Alfred O. Hero. “Covariance Estimation in High Dimensions Via Kronecker Product Expansions”. In: *IEEE Transactions on Signal Processing* 61.21 (Nov. 2013), pp. 5347–5360. DOI: 10.1109/tsp.2013.2279355. URL: <https://doi.org/10.1109/2Ftsp.2013.2279355>.
- [19] Pauli Virtanen et al. “SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python”. In: *Nature Methods* 17 (2020), pp. 261–272. DOI: 10.1038/s41592-019-0686-2.
- [20] Yu Wang, Byoungwook Jang, and Alfred Hero. “The Sylvester Graphical Lasso (SyGlasso)”. In: *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*. Ed. by Silvia Chiappa and Roberto Calandra. Vol. 108. Proceedings of Machine Learning Research. PMLR, 26–28 Aug 2020, pp. 1943–1953. URL: <https://proceedings.mlr.press/v108/wang20d.html>.
- [21] Jun Ho Yoon and Seyoung Kim. “EiGLasso: Scalable Estimation of Cartesian Product of Sparse Inverse Covariance Matrices”. In: *Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence (UAI)*. Ed. by Jonas Peters and David Sontag. Vol. 124. Proceedings of Machine Learning Research. PMLR, Mar. 2020, pp. 1248–1257. URL: <https://proceedings.mlr.press/v124/ho-yoon20a.html>.
- [22] L. Zappia, B. Phipson, and A. Oshlack. “Splatter: simulation of single-cell RNA sequencing data”. In: *Genome Biol* 18.1 (Sept. 2017), p. 174.
- [23] Shuheng Zhou. “Gemini: Graph estimation with matrix variate normal instances”. In: *The Annals of Statistics* 42.2 (Apr. 2014). DOI: 10.1214/13-aos1187. URL: <https://doi.org/10.1214/2F13-aos1187>.

261 Checklist

262 The checklist follows the references. Please read the checklist guidelines carefully for information on
 263 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
 264 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
 265 the appropriate section of your paper or providing a brief inline description. For example:

- 266 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 267 • Did you include the license to the code and datasets? **[No]** The code and the data are
 268 proprietary.
- 269 • Did you include the license to the code and datasets? **[N/A]**

270 Please do not modify the questions and only use the provided macros for your answers. Note that the
 271 Checklist section does not count towards the page limit. In your paper, please delete this instructions
 272 block and only keep the Checklist section heading above along with the questions/answers below.

273 1. For all authors...

- 274 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 275 contributions and scope? **[Yes]** We claim that our algorithm is faster than competitors,
 276 and give evidence to this in Figure 1.

- 277 (b) Did you describe the limitations of your work? [Yes] As mentioned in the Results
 278 section, non-heuristic variant performs poorly on small sample data and BiGLasso
 279 algorithms in general perform poorly when data gets more complicated, such as
 280 duplicated video frames
- 281 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 282 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 283 them? [Yes]
- 284 2. If you are including theoretical results...
- 285 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 286 (b) Did you include complete proofs of all theoretical results? [Yes] In the supplementary
 287 material.
- 288 3. If you ran experiments...
- 289 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 290 mental results (either in the supplemental material or as a URL)? [Yes] Code that runs
 291 experiments discussed in this paper and the supplementary material is available as part
 292 of the supplementary material.
- 293 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 294 were chosen)? [Yes] In the supplementary material we talk about hyperparameters, and
 295 we talk specifically about the choice of regularization parameters in the regularization
 296 section.
- 297 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 298 ments multiple times)? [Yes] antGLasso's performance curves in the Regularization
 299 section of the supplementary material highlight the variance.
- 300 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 301 of GPUs, internal cluster, or cloud provider)? [Yes] This is Figure 1.
- 302 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 303 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 304 (b) Did you mention the license of the assets? [N/A]
- 305 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 306
- 307 (d) Did you discuss whether and how consent was obtained from people whose data you're
 308 using/curating? [N/A]
- 309 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 310 information or offensive content? [N/A]
- 311 5. If you used crowdsourcing or conducted research with human subjects...
- 312 (a) Did you include the full text of instructions given to participants and screenshots, if
 313 applicable? [N/A]
- 314 (b) Did you describe any potential participant risks, with links to Institutional Review
 315 Board (IRB) approvals, if applicable? [N/A]
- 316 (c) Did you include the estimated hourly wage paid to participants and the total amount
 317 spent on participant compensation? [N/A]