

Sparse Causal Discovery in Multivariate Time Series

Stefan Haufe* Guido Nolte
 Berlin Institute of Technology Fraunhofer First, Berlin
 Klaus-Robert Müller Nicole Krämer
 Berlin Institute of Technology Berlin Institute of Technology

January 15, 2009

Abstract

Our goal is to estimate causal interactions in multivariate time series. Using vector autoregressive (VAR) models, these can be defined based on non-vanishing coefficients belonging to respective time-lagged instances. As in most cases a parsimonious causality structure is assumed, a promising approach to causal discovery consists in fitting VAR models with an additional sparsity-promoting regularization. Along this line we here propose that sparsity should be enforced for the subgroups of coefficients that belong to each pair of time series, as the absence of a causal relation requires the coefficients for all time-lags to become jointly zero. Such behavior can be achieved by means of $\ell_{1,2}$ -norm regularized regression, for which an efficient active set solver has been proposed recently. Our method is shown to outperform standard methods in recovering simulated causality graphs. The results are on par with a second novel approach which uses multiple statistical testing.

Keywords Vector Autoregressive Model, Granger Causality, Group Lasso, Multiple Testing

1 Introduction

Causality is commonly defined based on the widely accepted assumption that an effect is always preceded by its cause. Granger (1969) postulates a measure of causal influence between two time series (*Granger Causality*). In a nutshell, time series z_i Granger-causes time series z_j if knowledge of past values of z_i improves the prediction of z_j (compared to only using past values of z_j). In the case of a set $F = \{z_1, \dots, z_M\}$ of time series, the pairwise analysis may lead to spurious detection of a causal relation. For this reason it is advisable to include

*haufe@cs.tu-berlin.de

the set $F \setminus \{z_i, z_j\}$ of all additional observable time series in both prediction tasks. Note that this approach resolves the problem of common hidden factors z_* if $z_* \in F$. (If common factors are not observable, Granger causality fails and we refer to Nolte et al. (2008) for a detailed discussion and a remedy in form of the Phase Slope Index.) While this approach, to which we refer as *complete* Granger Causality, is practical, a more elegant way to deal with multivariate data is to handle all potential causal relations between all time series at once. In this paper, we assume a linear dynamics of the underlying system, which leads to the vector autoregressive (VAR) model.

In many applications the true causality graph is assumed to be sparse, i.e. only a few causal interactions between time series are expected. However, both Ordinary Least Squares (OLS) and Ridge Regression, which are usually used for fitting VAR models, are known for producing dense coefficients. Only recently Valdes-Sosa et al. (2005) have proposed to enforce estimation of sparse AR coefficients using ℓ_1 -norm regularized models such as the Lasso (Tibshirani, 1996).

In this paper we propose a novel sparse approach which – unlike Lasso – accounts for the fact that the absence of a causal relation between z_i and z_j requires all AR coefficients belonging to that certain pair of time series to be jointly zero. Furthermore, we consider Ridge Regression in combination with the multiple statistical testing procedure provided by Hothorn et al. (2008). More details on the methodology are given in Section 3. These methods are evaluated and compared to standard approaches in extensive simulations.

2 Background

In this section, we briefly summarize related approaches to estimate sparse vector autoregressive models in the context of causal discovery. We roughly distinguish between sparse estimation methods and testing strategies.

Given a multivariate time series $\mathbf{z}(t) \in \mathbb{R}^M$, a linear vector autoregressive process of order P is defined as

$$\mathbf{z}(t) = \sum_{p=1}^P A^{(p)} \mathbf{z}(t-p) + \varepsilon(t), \quad (1)$$

where $A^{(p)} \in \mathbb{R}^{M \times M}$, $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$ and $t \in \mathbb{Z}$ indicates time. Hence, the signal at time t is modeled as a linear combination of its P past values and Gaussian measurement noise. Inspired by the initial assumption that the cause should always precede the effect, we suggest the following definition of causality in the case of vector autoregressive models. We say that time series z_i has a causal influence on time series z_j if for at least one $p \in \{1, \dots, P\}$, the coefficient $A_{ji}^{(p)}$ corresponding to the interaction between z_j and z_i at the p th time-lag is nonzero.

Thus, causal inference may be conducted by estimating the matrices $A^{(p)}$ from a sample $Z = (\mathbf{z}(1), \dots, \mathbf{z}(T))$. Let us introduce the following shortcuts. We denote by $A = (A^{(1)}, \dots, A^{(P)})^\top$ the matrix of all VAR coefficients and set $X = (Z_1, \dots, Z_P)$, $Y = Z_0$, $Z_p = (\mathbf{z}(P+1-p), \dots, \mathbf{z}(T-p))^\top$. Here $\text{vec}(\cdot)$ denotes the vectorization operation.

2.1 Sparsity

Probably the most straightforward way to estimate a sparse VAR is to use ℓ_1 -regularization on the set of coefficients,

$$\hat{A}^{\text{lasso}} = \arg \min_A \|\text{vec}(XA - Y)\|_2^2 + \lambda \|\text{vec}(A)\|_1, \lambda \geq 0.$$

Recently, Valdes-Sosa et al. (2005) proposed a combination of VAR-estimation and the Lasso (Tibshirani, 1996). While Valdes-Sosa et al. (2005) only consider a VAR model of order 1, there have been extensions to higher orders (e.g. Arnold et al., 2007). However, we note in the latter case, Lasso is not used on the VAR coefficients directly, but that the problem is transformed into the task of estimating partial correlation coefficients between time-lagged copies of the time series (see also Opgen-Rhein and Strimmer, 2007).

2.2 Testing

Just as in the case of sparse methods, it is often suggested to transform the regression task into the estimation of the matrix of partial correlation coefficients between time-lagged copies of the time series. While Drton and Perlman (2008) estimate the correlation matrix in an unregularized way, Opgen-Rhein and Strimmer (2007) propose a shrinkage estimator, which is superior in the case of high-dimensional data (Schäfer and Strimmer, 2005). Afterwards, significant partial correlations are detected by controlling false discovery rates. While the latter approach is only tested for $P = 1$, it is straightforward to extend it to higher order VAR's.

3 Our Approach

In the following, we provide the details regarding the groupwise sparsity and the alternative testing strategy respectively.

3.1 Ridge Regression and Multiple Testing

Under the assumption of Gaussian white noise it is natural to estimate the AR coefficients using regularized least squares, and probably the most straightforward way to do so is to use Ridge Regression,

$$\hat{A}^{\text{ridge}} = \arg \min_A \|\text{vec}(XA - Y)\|_2^2 + \lambda \|\text{vec}(A)\|_2^2 = (X^\top X + \lambda I)^{-1} X^\top Y, \lambda \geq 0. \quad (2)$$

Thanks to the Ridge penalty, Eq. (2) delivers solutions with small coefficients, which, however, are in general never exactly zero. In the strict sense of Granger, this corresponds to a fully-connected dependency graph, rendering Ridge Regression an improper candidate for sparse causal recovery. On the other side, many of the estimated coefficients are expected to be non-significant. Hence, we propose a sparsification through statistical testing. In contrast to e.g. bootstrapping, we derive p -values explicitly using the approximate distribution of the coefficients.

It is apparant from Eq. (2) that the estimation can be done independently for each column of A , and so does the testing. Let therefore α_k denote the k th column of A and let $\mathbf{y}_k = (z_k(P+1), \dots, z_k(T))^\top$. Neglecting the dependency of X and Y , the Ridge coefficients depend linearly on Y , and we can conclude that under the null-hypothesis $H_0 : \alpha_k = 0$, we have $\hat{\alpha}_k \sim \mathcal{N}(\mathbf{0}, \sigma_k^2 \Sigma)$ with

$$\Sigma = (X^\top X + \lambda I)^{-1} X^\top X (X^\top X + \lambda I)^{-1}.$$

Furthermore, setting $H = X (X^\top X + \lambda I)^{-1} X^\top$ an estimate of the model variance σ_k^2 is given by

$$\hat{\sigma}_k^2 = \frac{\|\mathbf{y}_k - H\mathbf{y}_k\|^2}{\text{trace}((I - H)(I - H^\top))}. \quad (3)$$

Using Eq. (3) we can now construct normalized test statistics $\tilde{\alpha}_{ik} = \hat{\alpha}_{ik} / \sqrt{\hat{\sigma}_k^2 \Sigma_{ii}}$ which are jointly normally distributed with $\tilde{\alpha} \sim \mathcal{N}(\mathbf{0}, R)$ and $R_{ij} := \Sigma_{ij} / \sqrt{\Sigma_{ii} \Sigma_{jj}}$. Suppose we want to test all individual hypotheses $H_{0,i} : \alpha_{ik} = 0$ simultaneously, then, according to Hothorn et al. (2008), the adjusted p -values are $p_i = 1 - g(R, |\tilde{\alpha}_{ik}|)$. We reject a hypothesis, if the p -value is below the predefined significance level γ . Here,

$$g(R, t) = P\left(\max_i |\tilde{\alpha}_{ik}| \leq t\right) = \int_{-t}^t \dots \int_{-t}^t \phi(\alpha_1, \dots, \alpha_{MP}) d\alpha_1 \dots d\alpha_{MP} \quad (4)$$

and $\phi(\alpha)$ is the density function of the multivariate normal distribution $\mathcal{N}(\mathbf{0}, R)$.

3.2 Group Lasso

Sparse causal discovery using Ridge Regression is a two-step procedure and may possibly suffer from the aggregation of assumptions that enter in each step. Direct estimation of sparse VAR coefficients (e.g. via Lasso) is therefore desirable, as this would allow omission of the multiple significance testing step. However, for higher order models, this approach is prone to selecting a different set of causal interactions for each of the P time lags. We here suggest that this behavior can be overcome by enforcing *joint sparsity* of the coefficient vectors that belong to a certain pair of time series. This corresponds to incorporating the prior belief that causal influences between time series are not restricted to only one particular time lag into the estimation. The positive effect of such modeling can be verified in Figure 1 (see Section 4 for more details).

The idea of imposing groupwise sparse coefficients leads to $\ell_{1,2}$ -norm regularized regression also known as the *Group Lasso* (Yuan and Lin, 2006; Meier et al., 2008), which has also applications in Multiple Kernel Learning (Bach et al., 2004; Sonnenburg et al., 2006) and the EEG/MEG inverse problem (e.g. Haufe et al., 2008). The term $\ell_{1,2}$ -norm corresponds to an ℓ_1 -norm of a vector of ℓ_2 -norms. Our proposed objective is given by

$$\begin{aligned} \hat{A}^{\text{glasso}} &= \arg \min_A \|\text{vec}(XA - Y)\|_2^2 \\ \text{s.t.} \quad &\left\| \left(A_{11}^{(1)}, \dots, A_{MM}^{(P)} \right) \right\|_2 + \sum_{i \neq j} \left\| \left(A_{ij}^{(1)}, \dots, A_{ij}^{(P)} \right) \right\|_2 \leq \kappa, \end{aligned} \quad (5)$$

This penalty leads to a groupwise variable selection, i.e. a whole block of coefficients is jointly zero. Note that the first term in Eq. (6) penalizes all MP coefficients describing univariate relations. In this way, those coefficients are shrunk and hence, overfitting is avoided. Furthermore, we remark that it is also conceivable to split the estimation of A into M subproblems (as suggested in Subsection 3.1), which is desirable in large-scale scenarios.

Eqs. (5) and (6) define a non-differentiable but convex optimization problem which can be solved by means of Second-order Cone Programming (SOCP). For problems with sparse expected structure, however, the optimization can be carried out much more efficiently using the results of Roth and Fischer (2008). By keeping a set of active coefficient groups, their algorithm needs to call the SOCP solver only for problem sizes far smaller than the original problem – leading to a considerable reduction of memory usage and computation time. In the experiments, we employ the active-set algorithm of Roth and Fischer (2008) in combination with a freely available SOCP solver (Sturm, 1999).

4 Simulations

We conduct a series of experiments in which the causal structure of simulated data has to be recovered. We compare the Group Lasso, standard Lasso, Ridge Regression with multiple testing and complete Granger Causality based on AR models. All four approaches are applied both with and without knowledge of the true model order. In the latter case $P = 10$ is chosen for the reconstruction. For all methods considered, it is also possible to estimate the model order P , e.g. via cross-validation.

4.1 Setup

Each simulated data set consists of a multivariate time series with parameters $M = 7$ and $T = 1000$ that is generated by a random VAR process of order $P = 5$ according to (1). The distribution of the noise component $\varepsilon(t)$ is chosen to be the standard normal distribution. The VAR coefficients for all but 10 randomly chosen pairs of time series are set to zero, yielding exactly 10 causal interactions. The non-zero coefficients are drawn randomly from $\mathcal{N}(0, 0.04I)$. Each set of

VAR coefficients is tested for the stability of its induced dynamical system by looking at the eigenvalues of the corresponding transition matrix. Only coefficients leading to stable systems (i.e. those with transition matrices with eigenvalues of at most 1) are accepted. We consider the following three types of problems, for each of which we create 10 instances: 1) no noise is added to the data generated by the VAR model 2) the data is superimposed by Gaussian noise of approximately the same strength, which is uncorrelated (white) both across time and sensors 3) the data is superimposed by mixed noise of approximately the same strength, which is generated as a random instantaneous mixture of M univariate AR processes of order 20. Note that in none of these cases the noise itself possesses a causal structure which would superimpose the true structure.

For measuring performance we consider Receiver Operating Characteristics (ROC) curves, which allow objective assessment of the performance in different regimes (e.g. very few false positives). As an additional measure of absolute performance we calculate the Area Under Curve (AUC). ROC curves and AUC values are averaged across the 10 problem instances and standard errors are computed for AUC.

Granger Causality is calculated using the Levinson-Wiggins-Robinson algorithm for fitting AR models (Marple, 1987), which is available in the open source Biosig toolbox (Schlögl, 2003). Note that for this particular method, we use Granger’s original definition of causal influence instead of our coefficient-based approach. That is, for a pair of time series z_i and z_j we calculate the logarithm of the ratio of the residuals of the two AR models 1) including interactions and 2) excluding interactions between z_i and z_j (*Granger score*). This score is divided by its standard deviations as estimated by the jackknife. To obtain a ROC-curve, the Granger score is threshold at different values, ranging from completely sparse to completely dense solutions.

For Ridge Regression, the regularization parameter λ is chosen via 10-fold cross-validation (with respect to prediction accuracy). For this value of λ , we derive the test statistics defined in Subsection 3.1. The multidimensional integrals in Eq. (4) are computed using Monte Carlo sampling according to Genz (1992). ROC-curves are constructed by varying the significance level γ .

For Lasso and Group Lasso, solutions ranging from completely sparse to completely dense are obtained through variation of the regularizing constant λ and κ respectively.

4.2 Results and Discussion

First, we illustrate the different behavior of the investigated methods in Figure 1. This example corresponds to the situation without noise and with known model order $P = 5$. The left figure shows the true underlying causal structure, with a black box indicating a causal interaction. The reconstructions for the different methods are based on a single estimate of the VAR coefficients. For Granger causality, we use a threshold of 2. For Ridge Regression, we use a significance level of $\gamma = 0.05$. For Lasso, Ridge Regression and Group Lasso, the regularizing constant is fixed by using 10-fold cross-validation (with respect

to prediction accuracy). We display the binary influence matrix in Figure 1. In this example Ridge Regression exhibits perfect reconstruction and outperforms all other methods. Group Lasso comes second. Note that, due to the strong tendency of Group Lasso to select the same influences for each time lag, its estimated causal dependency matrix is sparser than that of Lasso.

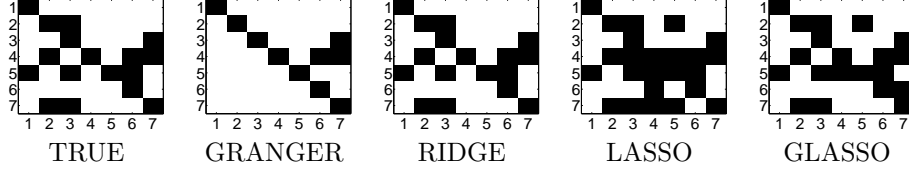


Figure 1: Simulated causal influence matrix and estimates according to Granger Causality, Ridge Regression, Lasso and Group Lasso.

Table 1 summarizes the AUC scores obtained in the experiments described above. The complementing ROC curves are shown in Figure 2. In short it can be stated that Group Lasso and Ridge Regression outperform their competitors in all scenarios, although not always significantly. While Ridge Regression performs slightly better than Group Lasso in the noiseless condition, Group Lasso has a clearly visible yet insignificant advantage over all methods in the white noise setting. Under the influence of mixed noise Ridge Regression and Group Lasso are on par. Note furthermore that the ROC curve for Lasso is below the ROC curve of Group Lasso, which shows that Lasso tends to be too dense. Interestingly, knowledge of the true model order hardly provided any significant advantage in our simulations.

		GRANGER	RIDGE	LASSO	GLASSO
$P = 5$	NO NOISE	0.991 ± 0.004	1.000 ± 0.000	0.996 ± 0.002	0.997 ± 0.002
	WHITE NOISE	0.910 ± 0.023	0.948 ± 0.020	0.941 ± 0.021	0.971 ± 0.016
	MIXED NOISE	0.896 ± 0.012	0.928 ± 0.010	0.889 ± 0.011	0.926 ± 0.012
$P = 10$	NO NOISE	0.980 ± 0.005	0.998 ± 0.002	0.996 ± 0.002	0.999 ± 0.001
	WHITE NOISE	0.885 ± 0.019	0.958 ± 0.012	0.948 ± 0.013	0.979 ± 0.005
	MIXED NOISE	0.893 ± 0.013	0.931 ± 0.015	0.861 ± 0.014	0.931 ± 0.007

Table 1: Average AUC scores and standard errors of Granger Causality, Ridge Regression, Lasso and Group Lasso in three different noise conditions and for two different model orders. Entries with significant superior score are highlighted.

5 Conclusion

We presented a novel approach for causal discovery in multivariate time series which is based on the Group Lasso. As an alternative we also discussed Ridge

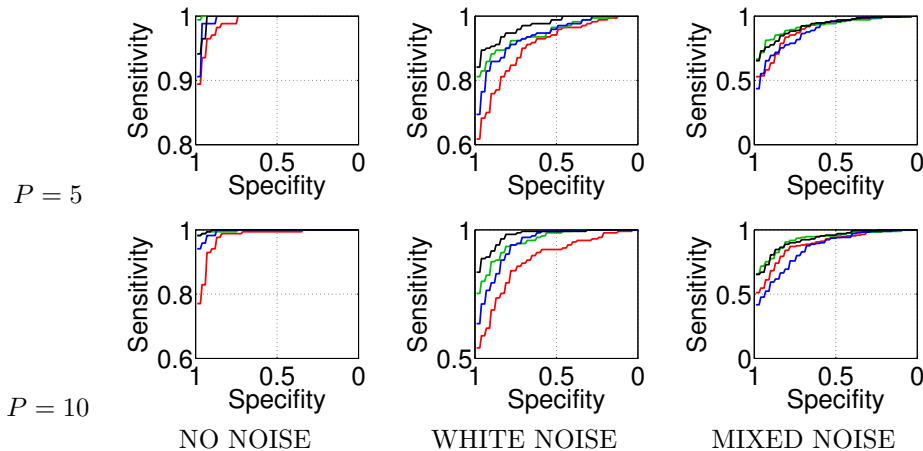


Figure 2: Average ROC curves of Granger Causality (red), Ridge Regression (green), Lasso (blue) and Group Lasso (black) in three different noise conditions and for two different model orders.

Regression with subsequent multiple testing according to Hothorn et al. (2008) which is also novel in the context of VAR modeling. Both approaches were shown to outperform standard methods in simulated scenarios. Future research will aim at applying our techniques to real-world problems. Given that the sparsity assumption is correct, our Group Lasso approach should be able to handle much larger problems than the ones that were considered here by 1) splitting the problem into M independent subproblems and 2) using the active set solver of Roth and Fischer (2008) in combination with strong regularization that ensures staying in the sparse regime. We expect that this will allow large-scale applications such as the estimation of cerebral information flow from functional Magnetic Resonance Tomography (fMRI) recordings to benefit from the improved accuracy of our approach.

Acknowledgements

This work was supported in part by the German BMBF (FKZ 01GQ0850, 01-IS07007A and 16SV2234) and the FP7-ICT Programme of the European Community under the PASCAL2 Network of Excellence, ICT-216886. We thank Thorsten Dickhaus for discussions.

References

Arnold, A., Liu, Y., and Abe, N. (2007). Temporal causal modeling with graphical granger methods. In *Proceedings of the Thirteenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 66–75.

- Bach, F., Lanckriet, G., and Jordan, M. (2004). Multiple kernel learning, conic duality and the SMO algorithm. In *Proceedings of the Twenty-first International Conference on Machine Learning*.
- Drton, M. and Perlman, M. (2008). A SInful approach to gaussian graphical model selection. *Journal of Statistical Planning and Inference*, 138(4):1179–1200.
- Genz, A. (1992). Numerical computation of multivariate normal probabilities. *Journal of Computational and Graphical Statistics*, 1:141–150.
- Granger, C. (1969). Investigating causal relations by econometric models and cross-spectral methods. *Econometrica*, 37:424–438.
- Haufe, S., Nikulin, V., Ziehe, A., Müller, K.-R., and Nolte, G. (2008). Combining sparsity and rotational invariance in EEG/MEG source reconstruction. *NeuroImage*, 42(2):726–738.
- Hothorn, T., Bretz, F., and Westfall, P. (2008). Simultaneous inference in general parametric models. *Biometrical Journal*, 3:346–363.
- Marple, S. (1987). *Digital spectral analysis with applications*. Prentice Hall, Englewood Cliffs, NJ.
- Meier, L., van de Geer, S., and Bühlmann, P. (2008). The group lasso for logistic regression. *Journal of the Royal Statistical Society: Series B*, 70(1):53–71.
- Nolte, G., Ziehe, A., Nikulin, V., Schlögl, A., Krämer, N., Brismar, T., and Müller, K. (2008). Robustly estimating the flow direction of information in complex physical systems. *Physical Review Letters*, 100(23):234101.
- Opgen-Rhein, R. and Strimmer, K. (2007). Learning causal networks from systems biology time course data: an effective model selection procedure for the vector autoregressive process. *BMC Bioinformatics*, 9.
- Roth, V. and Fischer, B. (2008). The group lasso for generalized linear models: uniqueness of solutions and efficient algorithms. In *Proceedings of the 25th International Conference on Machine Learning*, pages 848–855.
- Schäfer, J. and Strimmer, K. (2005). A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics. *Statistical Applications in Genetics and Molecular Biology*, 4:32.
- Schlögl, A. (2003). BIOSIG - an open source software library for biomedical signal processing, <http://BIOSIG.SF.NET>.
- Sonnenburg, S., Rätsch, G., Schäfer, C., and Schölkopf, B. (2006). Large scale multiple kernel learning. *The Journal of Machine Learning Research*, 7:1531–1565.

- Sturm, J. (1999). Using SeDuMi 1.02, a MATLAB toolbox for optimization over symmetric cones. *Optimization Methods and Software*, 11–12:625–653.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B*, 58:267–288.
- Valdes-Sosa, P., Sanchez-Bornot, J., Lage-Castellanos, A., Vega-Hernandez, M., Bosch-Bayard, J., Melie-Garcia, L., and Canales-Rodriguez, E. (2005). Estimating brain functional connectivity with sparse multivariate autoregression. *Philosophical Transactions of the Royal Society B*, 360:969–981.
- Yuan, M. and Lin, Y. (2006). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society Series B*, 68(1):49–67.