

Dynamic Causal Networks with Multi-scale Temporal Structure

Eric D. Kolaczyk

Dept of Mathematics and Statistics, Boston University

Joint work with Xinyu Kang & Apratim Ganguly

Supported in part by AFOSR award 12RSL042.

At 10K Feet . . .

Two simple observations:

- 1 The automated, simultaneous monitoring of each unit in a large complex system has become commonplace.
- 2 Frequently the resulting data are a high-dimensional multivariate time series.

Implication: The combination of

- ▷ *systems* and
- ▷ *time series*

perspectives suggests the use of *dynamic network modeling*, a highly active frontier in the field of network analysis.

Illustration: Neuroscience

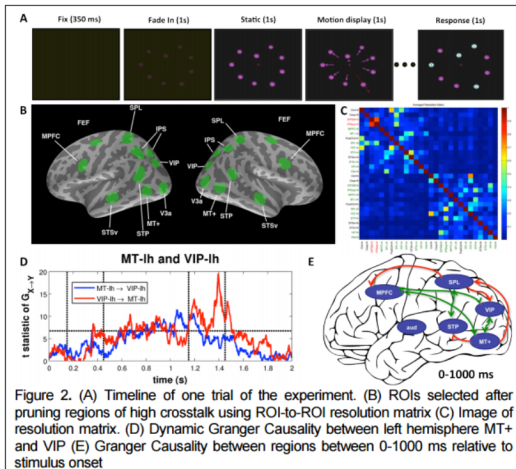
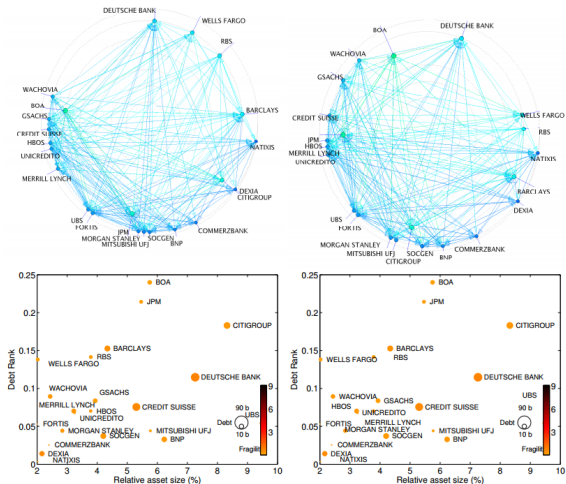


Illustration: Economics/Finance



Source: Battiston et al. (2012) *Nature Reports*

BU-Keio Workshop, August 2016

Dynamic Network Topology Inference

In these and similar application contexts, there is a basic paradigm at work:

Systems-level question of interest;

⇒ Collection of multivariate time series data;

⇒ Construction of network-based representation of system from data;

⇒ Network-centric answer to systems-level question.

The third step is known as (association) network topology inference.

Aside on Network Inference

There are a variety of methods for inference of a network¹

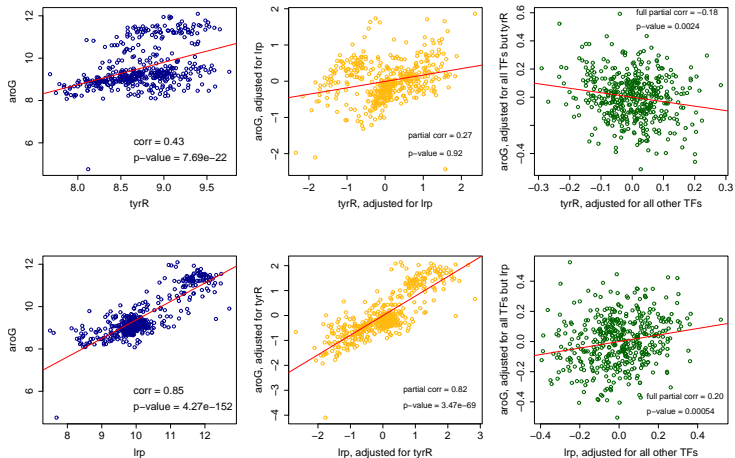
These can be largely categorized by choices in

- 1 notion of association (e.g., correlation, mutual information, etc.)
- 2 method of inference (e.g., testing, regression, etc.)
- 3 working parameters (e.g., significance level, smoothing parameter, etc.)

Note: Most formal methods seek to exploit sparseness typically encountered in empirical networks.

¹See, for example, Ch7 in either of Kolaczyk (2009). *Statistical Analysis of Network Data*, or Kolaczyk & Csardi (2014). *Statistical Analysis of Network Data with R*.

Illustration with Gene Regulation: Choices Matter!



Correlation

Partial Corr. w/ 1

Partial Corr. w/ All

An Additional Observation: Nature of Temporal Changes

In many contexts, we can expect changes in the system (local or global) across *multiple time scales*.

Examples:

- ▷ Neuroscience:
- ▷ Finance:

Suggests the need for *multi-scale analysis* . . . a concept well-established in time series analysis, but which has not yet emerged in dynamic network analysis.

Our Focus

Motivated by these various observations, we focus on the problem of **detecting dynamic connectivity changes across multiple time scales** in a network-centric representation of a system, based on high-dimensional multivariate time series observations.

Our approach combines

- ▶ Granger causal modeling, with
- ▶ partition-based multi-scale modeling

Granger Causal Models

Let \mathbf{X}_t be an N -dim time series, at time t , corresponding to a network of N nodes $v \in V$.

We will adopt a $VAR(p)$ model, i.e.,

$$\mathbf{X}_t = \sum_{\tau=1}^p \Theta_{\tau} \mathbf{X}_{t-\tau} + \epsilon_t ,$$

where Θ_{τ} collects the influence of the nodes on each other at lag τ , and the ϵ_t are independent white noise.

Say that X_v *Granger causes* X_u if $\Theta_{\tau}(u, v) \neq 0$ for some $\tau = 1, \dots, p$.

From Time Series to Network Graphs

In this setting, the notion of ‘network’ is made precise through graphs defined as a function of the underlying graphical model. That is, through *conditional independence* relations, which are coded in one-to-one correspondence with the non-zero elements of

$$\Theta = (\Theta_1, \dots, \Theta_p) .$$

Two options:

- ① $G = (V, E)$ is a directed graph with an edge from v to u iff $\|\theta(u, v)\| \neq 0$, where $\theta(u, v) = (\Theta_1(u, v), \dots, \Theta_p(u, v))^T$.
- ② $G = (V^{(p+1)}, E)$ is a directed multi-graph, with an edge from v in layer τ to u in layer 0 iff $\theta_\tau(u, v) \neq 0$.

We will focus on the first of these options.

Time-indexed Directed Graphical Model

Our interest will be in *non-stationary* multivariate time series, and the corresponding representations using dynamic networks.

We adopt a changepoint perspective, so that our model class consists of concatenations of Granger causal $VAR(p)$ models, each with its own Θ constant over a given interval of time.

The result is a time-indexed directed graphical model, from which we define a dynamic network $G_t = (V, E_t)$ in analogy to the stationary case.

Our goal is to infer (i) the changepoints distinguishing the stationary intervals, and (ii) the corresponding edge sets E_t .

Granger-causal Dynamic Network Modeling

A great deal of work has been done recently in modeling multiple time series using causal network types of models.

Extending seminal work of Meinshausen and Buhlmann², a selection of representative examples include

- ▶ Neighborhood-selection under stationarity, using group-lasso principles³
- ▶ Network Granger causality with panel data⁴
- ▶ Inference of networks defined through long-run correlation⁵
- ▶ Adaptive selection of (fixed) window size⁶.

Our work takes these ideas in a multi-scale direction.

²Meinshausen, N. and Buhlmann, P. (2006). High-dimensional graphs and variable selection with lasso. *AoS*

³Bolstad, A., Van Veen, B. D., & Nowak, R. (2011). Causal network inference via group sparse regularization. *IEEE TSP*.

⁴Basu, S., Shojaie A., & Michailidis, G. (2015). Network granger causality with inherent grouping structure. *JMLR*.

⁵Barigozzi, M., Brownlees, C. (2013). Nets: network estimation for time series. *Available at SSRN*.

⁶Long, C.J., Brown, EN., Triantafyllou, C., Aharon, I., Wald, LL., Solo, V. (2005). Nonstationary noise estimation in functional MRI. *NeuroImage*.

Multi-scale Modeling

Starting in the late 1980's (building on threads of work going back at least to the early 1900's), there was an explosion of development on methods of multi-scale modeling that were

- ▶ mathematically principled;
- ▶ computationally efficient; and, often,
- ▶ domain/problem-specific.

The quintessential example is that of methods utilizing transformations with respect to bases of *wavelets*.

Wavelets & Recursive Partitioning: A connection

While the development of wavelet-based methods for standard signal and image analysis applications proceeded apace, the development of extensions for less-traditional settings like

- ▷ non-normal noise, and
- ▷ structured data (e.g., manifolds, graphs, etc.)

proved decidedly more challenging⁷.

Helpful in many contexts was a fundamental result of Donoho⁸ relating

- 1 methods of recursive (dyadic) partitioning, and
- 2 selection of a best-orthonormal basis,

where the basis is selected from a class of (unbalanced) Haar bases.

⁷But, nevertheless, was/is being thoroughly explored as well!

⁸Donoho, D.L. (1997). CART and Best-Ortho Basis: A Connection. *AoS*.

Illustration: Multiscale GLMs

Hence, partition-based methods – which can be more amenable than wavelets to adaptation in nontraditional settings – are nevertheless effective structures upon which to build multi-scale models and methods.

A useful illustration of this principle is in the context of generalized linear models (GLMs)⁹

Given independent observations y_1, \dots, y_n , where

$$p_{\theta}(y_i | t_i) = \exp \left\{ \frac{y_i \cdot \theta(t_i) - b[\theta(t_i)]}{\tau} + c(y, \tau) \right\},$$

estimate $\theta(\cdot)$, a member of some inhomogeneous function class (e.g., Besov ball).

⁹Kolaczyk, E.D, and Nowak, R.D (2005). Multiscale generalized linear models for nonparametric function estimation. *Biometrika*.

Multiscale GLMs (cont)

Estimate $\theta = (\theta_1, \dots, \theta_n)$ as

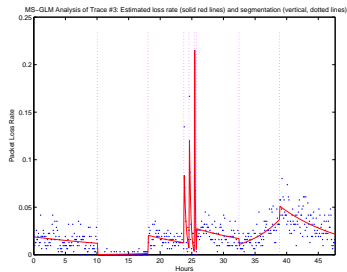
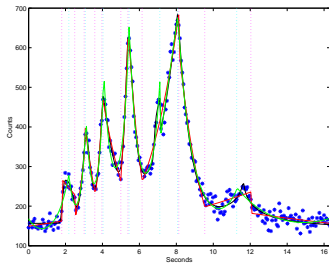
$$\hat{\theta} \equiv \arg \max_{\mathcal{P} \preceq \mathcal{P}_{Dy}^*} \max_{\theta' \in PP(\mathcal{P}; D)} \{ \ell(\theta') - 2\lambda \#(\mathcal{P}) \} ,$$

where

- ▷ \mathcal{P}_{Dy}^* is a complete recursive dyadic partition;
- ▷ $\mathcal{P} \preceq \mathcal{P}_{Dy}^*$ denotes a sub-partition;
- ▷ $PP(\mathcal{P}; D) \equiv \{ \text{Piece-wise polynomials, of order } D \text{ on } \mathcal{P} \}$; and
- ▷ $\ell(\theta)$ is the log-likelihood

Multiscale GLMs: Results Summary

- ▷ $O(n)$ alg¹⁰ and near-optimal risk behavior
- ▷ Extension to recursive partitioning, with $O(n^3)$ alg
- ▷ Application to gamma-ray bursts & Internet packet loss.



¹⁰Formally, $O(n)$ model comparison steps, where complexity of comparison step depends on models being fit

Our Contribution

We propose

- 1 a multi-scale Granger-based dynamic network model, and
- 2 a corresponding method of network topology inference

that captures the dynamics of a system in a manner sensitive to changes at multiple time scales, while encouraging sparsity of network connectivity.

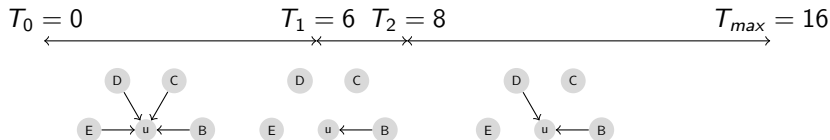
Key elements of our framework:

- ▷ We partition the non-stationary time space into blocks at various scales, with stationary assumed within each block;
- ▷ We do neighborhood selection with the group lasso, within each block.

A Cartoon Version

We estimate network topology neighborhood by neighborhood.

WLOG, consider the local neighborhood of node/series u .

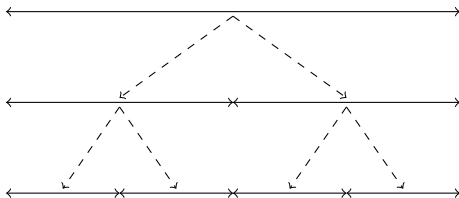


What we do:

- ▶ **Changepoint Detection:** Estimate the times T that the changes happened.
- ▶ **Neighborhood Selection:** Given the estimated change points, infer the neighborhood structure from time 0 to \hat{T}_1 and from time $\hat{T}_1 + 1$ to \hat{T}_2 and so on.

Parameterization of Changepoints

To capture potential change points, we partition time into candidate stationary intervals using RDP¹¹.



For appropriate choice of cost function, our problem can be solved in a dynamic programming fashion, in strict analogy to, e.g., MS-GLMs.

¹¹Extension to RPs straightforward, with increased computational cost.

Recovery of Neighborhood Structure

Within a candidate stationary interval, WLOG with n time points, we do neighborhood selection with group lasso, solving for

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \|\mathbf{X}_u - \sum_{v \in V \setminus \{u\}} \mathbf{X}_{(v)} \theta_{(v)}\|_2^2 + \lambda \sum_{v \in V \setminus \{u\}} \|\theta_{(v)}\|_2.$$

Here

- ▶ \mathbf{X}_u is the $n \times 1$ time series observed at node u ;
- ▶ $\mathbf{X}_{(v)}$ is the $n \times p$ matrix formed from \mathbf{X}_v at lags $\tau = 1, \dots, p$;
- ▶ $\theta_{(v)} = \theta(u, v) = (\Theta_1(u, v), \dots, \Theta_p(u, v))^T$;
- ▶ λ is a smoothing parameter (TBD).

Write the *group lasso penalized likelihood (gIPL)* as

$$g\hat{IPL} = \frac{1}{n} \left\| \mathbf{X}_u - \mathbf{X} \hat{\theta} \right\|_2^2 + \lambda \sum_{v \in V \setminus \{u\}} \left\| \hat{\theta}_{(v)} \right\|_2.$$

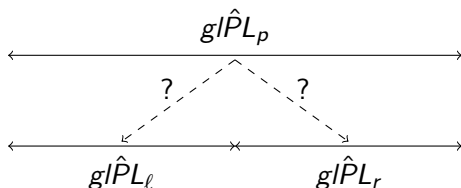
Selecting a Multiscale Dynamic Granger Causal Model

Our overall problem can be set up as solving for

$$\hat{\theta} \equiv \arg \max_{\mathcal{P} \preceq \mathcal{P}_{Dy}^*} \max_{\theta' \in PGC(\mathcal{P}; p)} \{ \ell(\theta') - \text{Pen}(\theta') \} ,$$

where $PGC(\mathcal{P}; p)$ is the collection of 'piecewise Granger Causal' models on \mathcal{P} of lag p .

Key Algorithmic Insight: This may be solved using dynamic programming, where at each potential break point, we solve the following problem:



Split if:

$$g|\hat{P}L_\ell + g|\hat{P}L_r < g|\hat{P}L_p$$

Consistency of Splitting

Define the truth to be: $\mathbf{X}_u = \begin{pmatrix} \mathbf{X}_u^{(\ell)} \\ \mathbf{X}_u^{(r)} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^{(\ell)}\theta^{(\ell)} + \epsilon^{(\ell)} \\ \mathbf{X}^{(r)}\theta^{(r)} + \epsilon^{(r)} \end{pmatrix}$, where $\mathbf{X}_u^{(\ell)}$ are the first $n/2$ observations and $\mathbf{X}_u^{(r)}$ the last $n/2$ observations.

Theorem

Under appropriate assumptions, for any sequence λ_n such that $\lambda_n \rightarrow 0$ and $\lambda_n n^{1/2} \rightarrow +\infty$,

$$\mathbb{P}_{\theta^{(\ell)}=\theta^{(r)}} \left(g\hat{I}PL_\ell + g\hat{I}PL_r \geq g\hat{I}PL_p \right) \rightarrow 1 \quad (1)$$

$$\mathbb{P}_{\theta^{(\ell)} \neq \theta^{(r)}} \left(g\hat{I}PL_\ell + g\hat{I}PL_r \leq g\hat{I}PL_p \right) \rightarrow 1 \quad (2)$$

Proof exploits consistency, within stationary intervals, established for group-lasso based neighborhood selection¹².

¹²Bolstad, A., Van Veen, B. D., & Nowak, R. (2011). Basu, S., Shojaie A., & Michailidis, G. (2015)

Finite Sample Control of Type I Error Rate

Theorem

If the penalty parameter $\lambda(\alpha)$ is chosen such that

$$\lambda(\alpha) = 2\hat{\sigma}_u \sqrt{pQ\left(1 - \frac{\alpha}{N^2}\right)}$$

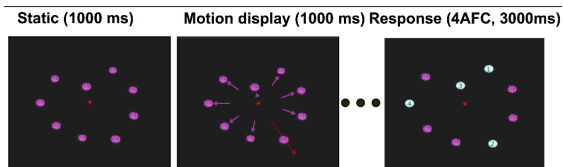
where $\hat{\sigma}_u^2 = \frac{\|\mathbf{x}_u\|_2^2}{n}$, p the order of AR lags, $N = |V|$, and $Q(\cdot)$ the quantile function of $\chi^2(p)$, then

$$\mathbb{P}(\exists u \in V : \hat{C}_u^\lambda \not\subseteq C_u) \leq \alpha.$$

Here C_u is the connected component in G to which u belongs, and \hat{C}_u^λ , its estimate.

Brain Image Data

We have trial-by-trial *visual search* data¹³.



MEG data of 8 subjects who were asked to watch moving objects displayed on a black screen.

Each subject received 160 trials and the MEG of their brain was taken accordingly.

Each resulting time series truncated to 1190 sample points, with intervention centered at time $t = 1/2$.

¹³Data courtesy of Lucia M. Vaina
BU-Keio Workshop, August 2016

1st Attempt @ Modeling: High-level Results

Optimal piecewise Granger causal model was fit (Using VAR(7); and Recursive Dyadic Partition), yielding:

- ▶ a 3-node dynamic network.
- ▶ Changepoints detected only at dyadic positions $\{1/4, 1/2, 3/4\}$.
- ▶ Causal relations detected between pairs of nodes

Aggregate Rates for Changepoints and Edges

Subject	1	2	3	4	5	6	7	8
Rate of split	(0.32)	(0.22)	(0.23)	(0.19)	(0.19)	(0.34)	(0.46)	(0.44)

	First Half				Second Half			
	IPS	FEF	V3a		IPS	FEF	V3a	
Subject 1	IPS	0	0.10	0.09	IPS	0	0.09	0.10
	FEF	0.11	0	0.08	FEF	0.1	0	0.07
	V3a	0.11	0.09	0	V3a	0.1	0.09	0
Subject 2	IPS	0	0.08	0.11	IPS	0	0.05	0.09
	FEF	0.08	0	0.03	FEF	0.04	0	0.03
	V3a	0.11	0.04	0	V3a	0.09	0.02	0
Subject 3	IPS	0	0.01	0.18	IPS	0	0.03	0.33
	FEF	0.01	0	0.01	FEF	0.03	0	0.03
	V3a	0.18	0.04	0	V3a	0.32	0.04	0
Subject 4	IPS	0	0.09	0.01	IPS	0	0.08	0.03
	FEF	0.10	0	0.04	FEF	0.09	0	0.06
	V3a	0.01	0.04	0	V3a	0.02	0.03	0
Subject 5	IPS	0	0.08	0.04	IPS	0	0.06	0.05
	FEF	0.09	0	0.01	FEF	0.08	0	0.02
	V3a	0.03	0.06	0	V3a	0.06	0.03	0
Subject 6	IPS	0	0.09	0.08	IPS	0	0.09	0.06
	FEF	0.08	0	0.12	FEF	0.09	0	0.11
	V3a	0.07	0.12	0	V3a	0.08	0.10	0
Subject 7	IPS	0	0.19	0.21	IPS	0	0.19	0.25
	FEF	0.20	0	0.08	FEF	0.19	0	0.04
	V3a	0.19	0.08	0	V3a	0.23	0.01	0
Subject 8	IPS	0	0.25	0.18	IPS	0	0.14	0.21
	FEF	0.25	0	0.09	FEF	0.15	0	0.09
	V3a	0.17	0.09	0	V3a	0.21	0.10	0

Closing Thoughts

Dynamic network analysis is arguably one of the most active frontiers in 'network science'.

Statisticians have (uncharacteristically!) jumped in early on the problem of network topology inference.

Neuroscience and economics/finance has a great deal of interest in and activity on this topic.

Still to do for piecewise Granger causal models:

- ▶ implement RP variant (more refined than RDP)
- ▶ theoretical characterization of overall model selection
- ▶ a more careful empirical study with MEG data

Thank you