

Learning Graphical Models for Stationary Time Series

Francis R. Bach and Michael I. Jordan

Abstract—Probabilistic graphical models can be extended to time series by considering probabilistic dependencies between entire time series. For stationary Gaussian time series, the graphical model semantics can be expressed naturally in the frequency domain, leading to interesting families of structured time series models that are complementary to families defined in the time domain. In this paper, we present an algorithm to learn the structure from data for directed graphical models for stationary Gaussian time series. We describe an algorithm for efficient forecasting for stationary Gaussian time series whose spectral densities factorize in a graphical model. We also explore the relationships between graphical model structure and sparsity, comparing and contrasting the notions of sparsity in the time domain and the frequency domain. Finally, we show how to make use of Mercer kernels in this setting, allowing our ideas to be extended to nonlinear models.

Index Terms—Frequency domain analysis, modeling, sparse matrices, spectral analysis, statistics, time series.

I. INTRODUCTION

TIME series arise in many problems in signal processing, bioinformatics, computer vision, and machine learning. In the statistical modeling of time series, the assumption of stationarity makes possible the use of tools from spectral analysis [1]. Much of the algorithmic research effort in this area has dealt with making such tools scalable with respect to the number of observed samples, for example, via algorithms that exploit the fast Fourier transform or forecasting algorithms such as the Durbin–Levinson algorithm [2].

Domains with large number of variables—in which Markovian or general graphical models have excelled—have not attracted the same attention in the time series field. Graphical models [3], [4] provide a general framework for defining probabilistic models over large numbers of variables by building global models out of local interaction models. Numerous special cases are familiar in signal processing, including the Kalman filter, hidden Markov models, and factor analysis. In addition, these models come equipped with standard, numerically efficient learning and inference algorithms, which are algorithms that have had applications beyond signal processing and machine learning, such as in error-control coding [5].

Manuscript received July 15, 2003; revised December 9, 2003. This work was supported by the National Science Foundation under Grant IIS-9988642 and the Multidisciplinary Research Program of the Department of Defense under Grant MURI N00014-00-1-0637. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Hans-Andrea Loeliger.

F. R. Bach is with the Computer Science Division, University of California, Berkeley, CA 94114 USA (e-mail: fbach@cs.berkeley.edu).

M. I. Jordan is with the Computer Science Division and the Department of Statistics, University of California, Berkeley, CA 94114 USA.

Digital Object Identifier 10.1109/TSP.2004.831032

Graphical models for time series are generally defined in the *time domain*. That is, they define a transition probability distribution on a set of state variables, conditioning on values of these variables at previous time steps. Such models, which are often referred to as *dynamic Bayesian networks*, have had significant applications in areas such as bioinformatics and speech processing [6], [7]. In this paper, we consider an extension of the basic notion of graphical model that considers dependencies between entire time series [8], [9]. As we show, this extension can be naturally expressed in the frequency domain, making use of the *spectral representations* for stationary Gaussian time series. In this paper, we present an algorithm to learn the structure of directed graphical models for spectral representations of time series from data. The algorithm, which is presented in Section V, is a direct extension of algorithms for learning-directed graphical models for Gaussian data [10], [11].

We also discuss the problem of forecasting, which is the problem of predicting the future given the past. For stationary Gaussian time series, algorithms for forecasting can be defined in either the time or frequency domain. In Section IV, we present a novel algorithm for efficient forecasting with graphical models. While classical algorithms such as the Durbin–Levinson algorithm incur a cubic time complexity [2], our new algorithm, by making use of the structure of the graphical model, has a quadratic complexity (in the number of variables). This new algorithm can also be used for factor analysis models for time series [12].

While our basic focus is on structured linear time series models, we also present an extension to nonlinear models in Section VI. In particular, we make use of our previous work in learning graphical models for hybrid data [13], enabling the fitting of augmented models that include nonlinearities or discrete variables. The principle of the algorithm is very simple: Although variables may not be Gaussian, by mapping them into a high-dimensional feature space, they can be considered as Gaussian for the purpose of model selection. This mapping is made implicit and computationally efficient through the use of kernel methods [14].

The graphical models that we describe in this paper have several possible applications, paralleling the many applications of graphical models for nontemporal data, such as feature selection for regression or classification, sparse and statistically sound models for forecasting, or simply a better understanding of the relationships between the different variables. In Section II, we review the necessary background on stationary Gaussian time series; in Section III, we show how the graphical framework can be naturally extended to such time series, with new numerically efficient inference procedures presented in Section IV. We then

present the structure learning algorithm in Section V for stationary Gaussian time series and extend it to nonlinear settings in Section VI. Finally, in Section VII, we report simulations on synthetic and real datasets to illustrate the validity of our algorithm and compare them with other approaches to model stationary time series.

II. STATIONARY GAUSSIAN PROCESSES

We consider a multiple time series $x(t)$, where for each $t \in \mathbb{Z}$, $x(t) = (x_1(t), \dots, x_m(t))$ has m univariate real components.¹

Throughout this paper, we assume that $x(t)$ is a zero-mean Gaussian process, that is, all finite sets of marginals are jointly Gaussian. In addition, we assume that the process $x(t)$ is *stationary*: For Gaussian processes, $x(t)$ is stationary if and only if $Ex(t+h)x(t)^\top$ does not depend on t , or equivalently, all marginals are invariant by time translation. Given a stationary process, the *autocovariance function* is defined as the following matrix-valued function $\Gamma(h)$ on \mathbb{Z} , which is defined as

$$\Gamma(h) = Ex(t+h)x(t)^\top = Ex(h)x(0)^\top.$$

For each $h \in \mathbb{Z}$, $\Gamma(h)$ is a symmetric $m \times m$ real matrix, and the function $h \mapsto \Gamma(h)$ is *non-negative*, that is, for all sets of vectors $\alpha_i \in \mathbb{R}^m$ indexed by $i \in I$, we have $\sum_{i,j \in I} \alpha_i^\top \Gamma(i-j) \alpha_j \geq 0$. Essentially, the Gaussian stationarity assumption is equivalent to modeling the variables as jointly Gaussian with tied parameters, i.e., the covariance matrix of any successive variables $x(t), x(t+1), \dots, x(t+h-1)$ is *Toeplitz* by blocks:

$$T_h(\Gamma) = \begin{pmatrix} \Gamma(0) & \Gamma(-1) & \Gamma(-2) & \cdots & \Gamma(-h+1) \\ \Gamma(1) & \Gamma(0) & \Gamma(-1) & & \\ \Gamma(2) & \Gamma(1) & \Gamma(0) & & \vdots \\ \vdots & & & \ddots & \\ \Gamma(h-1) & & \cdots & & \Gamma(0) \end{pmatrix}. \quad (1)$$

A. Spectral Density Matrix

We assume that $\sum_{-\infty}^{\infty} \|\Gamma(h)\|_2 < \infty$, where $\|M\|_2$ denotes the 2-norm of the matrix M , equal to its largest singular value, so that the *spectral density matrix* $f(\omega)$ is well defined, as a matrix-valued function on \mathbb{R} :

$$f(\omega) = \frac{1}{2\pi} \sum_{h=-\infty}^{+\infty} \Gamma(h) e^{-ih\omega}.$$

For each ω , $f(\omega)$ is an $m \times m$ Hermitian matrix. In addition, the function $\omega \mapsto f(\omega)$ is 2π -periodic. Moreover, for real-valued random variables, we have $f(-\omega) = f(\omega)^\top$, which implies that we only need to consider the spectral density matrix for $\omega \in [0, \pi]$.

Since the function $\Gamma(h)$ is non-negative, the spectral density matrix is pointwise non-negative, that is, $\forall \omega \in [0, 2\pi], \forall \alpha \in \mathbb{C}^m, \alpha^* f(\omega) \alpha \in \mathbb{R}^+$. Note that knowing the spectral density

f is equivalent to knowing the autocovariance function Γ since they form a Fourier transform pair. In particular, we have

$$\Gamma(h) = \int_0^{2\pi} f(\omega) e^{ih\omega} d\omega. \quad (2)$$

The rate of decay of $\Gamma(h)$ as h grows is determined by the smoothness of the spectral density. In particular, when designing a spectral density, care must be taken regarding its smoothness so that the resulting autocovariance function has attractive decay properties (which themselves govern the complexity of prediction).

B. Spectral Representation

Any Gaussian stationary time series with an absolutely summable autocovariance function has a spectral representation of the following form [2]:

$$x(t) = \int_0^{2\pi} e^{it\omega} dZ(\omega) \quad (3)$$

where $Z(\omega)$ is a random process with orthogonal increments such that for each $\omega_1 < \omega_2$, $\text{cov}(Z(\omega_2) - Z(\omega_1)) = \int_{\omega_1}^{\omega_2} f(\omega) d\omega$. In other words, $x(t)$ is a superposition of infinitely many independent random signals at different frequencies.

C. Autoregressive Models

In this paper, we also consider stationary autoregressive (AR) models. Stationarity can be imposed either by assuming that the process is initialized according to the stationary distribution or that the process is causal and extends to negative infinity. An autoregressive model has the following formulation:

$$x(t) = \Psi_1 x(t-1) + \cdots + \Psi_p x(t-p) + z(t)$$

where the variables $z(t)$ are mutually independent, each with covariance matrix Σ , and independent from $x(u)$, $u < t$. Each Ψ_i is an $m \times m$ matrix.

The parameters Ψ_i and the order p of AR models can be efficiently estimated from data by using classical regression methods for parameter estimation and variable selection [2]. If sparsity is desired, one can attempt to find zeros in the covariance matrix Σ or its inverse, as well as in the regression weight matrices Ψ_i , in a manner analogous to subset selection for linear regression [15]. In this paper, we consider sparsity in the frequency domain, which is a notion that is complementary to sparsity in the time domain.

D. Finite Sample

We now briefly review methods for estimating the autocovariance function $\Gamma(h)$ and the spectral density matrix $f(\omega)$ from data $x(0), \dots, x(T-1)$.

1) *Sample Autocovariance and Periodogram*: The *sample autocovariance function* is defined as

$$\hat{\Gamma}(h) = \frac{1}{T} \sum_{t=0}^{T-h-1} (x(t+h) - \bar{x})(x(t) - \bar{x})^\top$$

¹Note that the theory of stationary Gaussian processes and most of the results of this paper can be naturally extended to the complex case.

for $h \in [0, T-1]$, extended by symmetry for negative h and equal to zero for $|h|$ equal or greater than T (the vector $\bar{x} = (1/T) \sum_{t=0}^{T-1} x(t)$ is the sample mean of the data). The sample autocovariance function is consistent and asymptotically normal under weak assumptions [2].

The *periodogram* is a 2π -periodic matrix-valued function that is defined at the frequencies $\omega_k = 2\pi k/T$, $\omega_k \in [0, 2\pi]$ as the Fourier transform of the sample autocovariance function. More precisely, let $d(0), \dots, d(T-1)$ be the discrete Fourier transform (DFT) of the data

$$d(k) = \frac{1}{\sqrt{T}} \sum_{t=0}^{T-1} x(t) e^{-ikt}.$$

At the frequencies $\omega_k = 2\pi k/T$, $\omega_k \in [0, 2\pi]$, the periodogram is defined as $I(\omega_k) = (1/2\pi)d(k)d(k)^*$ and can readily be computed using m fast Fourier transforms (FFT). It is then extended as a piecewise constant periodic function on \mathbb{R} .

The periodogram does not provide a consistent estimator of the spectral density and is a notoriously bad estimator of the spectral density. However, when it is appropriately smoothed, it is a good estimator.

2) *Smoothing the Periodogram*: The periodogram is smoothed by convolving it with a smoothing window $W_r(j)$, leading to the following estimator, at frequency ω_k :

$$\hat{f}^r(\omega_k) = \sum_{j=-\infty}^{\infty} W_r(j) I(\omega_{j+k}). \quad (4)$$

$W_r(j)$ is a smoothing window that is required to be symmetric and to sum to one. In our simulations, we used the window $W_r(j)$, whose discrete Fourier transform is $w_k = (1/\sqrt{T})e^{-\omega_k^2 r^2/2}$. When r is less than $T/8$, the window is approximately equal to a Gaussian window of width r , i.e., $W_r(j) = (1/r\sqrt{2\pi})e^{-j^2/(2r^2)}$. This window has favorable properties both in time (it is positive) and frequency (it is smooth). We refer to r as the bandwidth of the smoothed periodogram. Note that if the number of samples T tends to infinity with $r(T)$ tending to infinity such that $r(T)/T$ tends to zero, then (4) provides a consistent estimate of the spectral density matrix [2].

Note that by simply inverting the Fourier transform using the FFT, we can recover the autocovariance function. With the type of smoothing we chose, we know in advance an upper bound on the decay of the autocovariance as h grows. Indeed, we can simply verify that the autocovariance function $\hat{\Gamma}^r$ derived from the estimate \hat{f}^r satisfies $\|\hat{\Gamma}^r(h)\|_2 \leq e^{-(2\pi/T)^2 h^2 r^2/2} \|\hat{\Gamma}(0)\|_2$. Thus, we can choose the time horizon H to be a constant times T/r (in simulations, we took $H = 4T/(r\sqrt{2\pi})$). Limiting the time horizon is equivalent to limiting the resolution of the spectra, that is, the spectral density is represented by its values on the grid $\omega_k = 2\pi k/H$ for a given H , and all integrals involving the density are approximated using Riemannian sums computed on this grid. Note that by the Nyquist theorem, if $\Gamma(h)$ is equal to zero for $|h| \geq H/2$, then the spectral density can be exactly represented by a finite sample with even steps $2\pi/H$. The periodogram gives a sample with precision $2\pi/T$; to obtain a sample with precision $2\pi/H$, we need to subsample the

Input: data $x_i(t)$, $i \in \{1, \dots, m\}$, $t \in \{0, \dots, T-1\}$

Algorithm:

1. Compute the m DFTs: $d_i(\cdot) = \text{DFT}(x_i(\cdot))$
2. Compute the periodogram $I(k) = \frac{1}{2\pi} d(k)d(k)^*$
3. Determine optimal smoothing r and degree of freedom df by minimizing $S(r)$ in Eq. (6), for r ranging from $T^{1/5}$ to $T^{4/5}$ with a grid of step $T^{3/10}$
4. Compute the spectral density matrix $f_k = \hat{f}(\omega_k)$ at frequencies $\omega_k = \frac{2k\pi}{T}$, $k \in \{0, \dots, T-1\}$ using Eq. (4)
6. Subsample the sequence f_k T/H times, with $H = 4T/(r\sqrt{2\pi})$

Output: $df = T/(r\sqrt{2\pi})$, f_k , $k \in \{0, \dots, H-1\}$

Fig. 1. Periodogram smoothing with automatic selection of the bandwidth.

periodogram with a positive linear filter (in order to maintain its positivity).

3) *Selecting the Bandwidth*: We use the Akaike information criterion (AIC) to select the optimal bandwidth for a given dataset. The AIC criterion is the sum of the negative log-likelihood of the data under the model, plus a function of the effective number of parameters or degrees of freedom df [15].

We use the *Whittle approximation* of the likelihood [12]:

$$\ell_w = -\frac{1}{2} \sum_{k=0}^{T-1} \left(\log \left| \hat{f}(\omega_k) \right| + \text{tr} \left\{ \hat{f}(\omega_k)^{-1} I(\omega_k) \right\} \right) - \frac{Tm}{2} \log 2\pi \quad (5)$$

where $|A|$ denotes the determinant of the matrix A . The Whittle approximation relies on the fact that the discrete Fourier transform of the data is asymptotically normal with independent components and with variance the spectral density [2].

In order to determine df , we notice that the smoothing defined in (4) is a linear smoothing, that is, $\hat{f}(\omega_k)$ is obtained from $I(\omega_k)$ by applying a linear matrix V_r , which is circulant with first row equal to the smoothing window W_r . The effective number of parameters is the trace of V_r [15], which, in our case, is simply obtained as $df_r = \text{tr} V_r = T \times W_r(0) = T/(r\sqrt{2\pi})$.

Thus, in order to find the optimal bandwidth, we minimize the following AIC criterion with respect to the smoothing parameter r :

$$S(r) = -\ell_w + \frac{df_r}{2} m^2. \quad (6)$$

(m^2 is the number of real parameters necessary to encode an $m \times m$ Hermitian matrix, and the $1/2$ comes from the fact that since our signals are real, we only need the spectral density on $[0, \pi]$). We perform an exhaustive search over $O(T^{1/2})$ values of the smoothing parameter r between $T^{1/5}$ and $T^{4/5}$ (the lower and upper bounds ensure that we obtain a consistent estimate of the spectral density matrix). The resulting algorithm is presented in Fig. 1. Note that we estimate the joint spectral density matrix first and then learn the structure of the graphical model (the topic of Section V; see Fig. 2). In Section V-C, we show how adaptive smoothing of the periodogram can be achieved in a manner that depends on the structure of the learned graphical model.

E. Forecasting

Given a finite sample from time 1 to H and a model (an autocovariance function or a spectral density), forecasting is the task

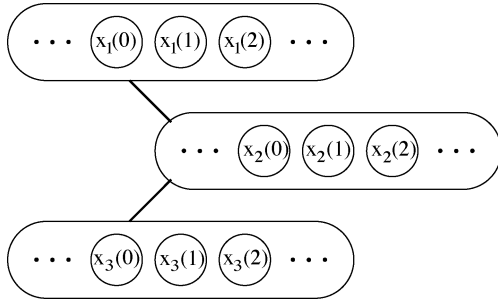


Fig. 2. Graphical model for three time series: x_1 is independent from x_3 given x_2 .

of predicting values for times $H + 1$ and later. In the Gaussian case, the best prediction is a linear approximation of $x(H + 1)$, based on $x(1), \dots, x(H)$. More precisely, we look for matrix coefficients Ψ_j such that the expected error

$$e(\Psi) = E \left\| x(H + 1) - \sum_{j=1}^H \Psi_j x(H + 1 - j) \right\|^2 \quad (7)$$

is minimal, where $\|y\| = (y^\top y)^{1/2}$ denotes the Euclidean norm of a vector y .

Minimizing $e(\Psi)$ in (7) is equivalent to solving the following system of equations [the so-called Yule–Walker equations, which are obtained by setting the derivatives to zero in (7)]:

$$\forall j \geq 1, \quad \sum_{i=1}^H \Psi_i \Gamma(i - j) = \Gamma(j). \quad (8)$$

We denote by γ_H the $H \times m$ vector such that $\gamma_H^\top = (\Gamma(1) \Gamma(2) \dots \Gamma(H))$, and Ψ as the $H \times m$ vector such that $\Psi^\top = (\Psi_1 \Psi_2 \dots \Psi_H)$. We also use the notation

$$Q_H(\Gamma) = \begin{pmatrix} \Gamma(0) & \Gamma(1) & \dots & \Gamma(H-1) \\ \Gamma(-1) & \Gamma(0) & & \\ \vdots & & \ddots & \\ \Gamma(-H+1) & \dots & & \Gamma(0) \end{pmatrix}.$$

Then, (8) can be written as

$$Q_H \Psi = \gamma_H. \quad (9)$$

In order to solve the system in (8) or (9), the Toeplitz structure of the problem can be exploited to avoid the $O(m^3 H^3)$ complexity associated with the unstructured linear system. In particular, the innovations algorithm or the Durbin–Levinson algorithm can be used to iteratively compute the prediction and the error covariance in $O(m^3 H^2)$ operations [2].

When the model is defined through the spectral density sampled at H points, then we first compute the autocovariance function using the FFT in time $O(m^2 H \log H)$ and incur a cubic time complexity for inference. In Section IV, we use a different technique to solve the Yule–Walker equations in time $O(m^2 H(d + \log H))$ for spectral densities that factorize in a graphical model with maximum fan-in d .

III. GRAPHICAL MODELS FOR TIME SERIES

The graphical model framework can be extended to multivariate time series in several ways. We follow [8] and [9] and

consider dependencies between whole time series, that is, between the entire sets $x_i \triangleq \{x_i(t), t \in \mathbb{Z}\}$, for $i = 1, \dots, m$. Thus, the time series x_i and x_j are independent if and only if the random infinite vectors $\{x_i(t), t \in \mathbb{Z}\}$ and $\{x_j(t), t \in \mathbb{Z}\}$ are independent. Similarly, time series x_i and x_j are conditionally independent, given x_k , if and only if $\{x_i(t), t \in \mathbb{Z}\}$ and $\{x_j(t), t \in \mathbb{Z}\}$ are conditionally independent, given $\{x_k(t), t \in \mathbb{Z}\}$.

For classical graphical models and Gaussian variables, marginal and conditional independence statements can be read out from zeros in the covariance matrix and its inverse [3]. It turns out that for Gaussian stationary time series, these results can be naturally extended, essentially replacing the covariance matrix by the spectral density matrix, as we now describe.

A. Gaussian Time Series and Independence

We consider a Gaussian stationary multivariate time series $x(t)$ with m components and with positive (i.e., invertible) spectral density matrix $f(\omega)$. Marginal independence and conditional independence are easily characterized in the frequency domain, as the following proposition shows [8].

Proposition 1: The time series x_i and x_j are *marginally independent* if and only if

$$\forall \omega \in [0, 2\pi], \quad f_{ij}(\omega) = 0.$$

The time series x_i and x_j are *conditionally independent* given all other time series x_k , $k \neq i, j$ if and only if

$$\forall \omega \in [0, 2\pi], \quad (f(\omega)^{-1})_{ij} = 0. \quad \blacksquare$$

Intuitively, this proposition enables us to consider each frequency component independently of the other components.

As a final step toward full equivalence between a Gaussian stationary time series and the concatenation of independent variables at each frequency, we have the following proposition, which gives a closed-form expression for the KL divergence [16], paralleling the expression for the KL divergence between zero mean Gaussian vectors [17]:

Proposition 2: The KL divergence between two zero-mean stationary vector-valued Gaussian processes U and V , with invertible spectral density matrices $f(\omega)$ and $g(\omega)$, is equal to

$$J(f(\omega) \| g(\omega)) = \frac{1}{2\pi} \int_0^{2\pi} I(f(\omega) \| g(\omega)) d\omega \quad (10)$$

where $I(F \| G) = -(1/2) \log \det[FG^{-1}] - (1/2) \text{tr}(I - FG^{-1})$ is the KL divergence between two covariance matrices. \blacksquare

As we will see in Section V, what is needed when learning a graphical model from data is the expression of the entropy of the random variable defined through a spectral density. The expression of interest is the entropy rate, which is defined as $h = \lim_{T \rightarrow \infty} (1/T) H(x(0), \dots, x(T-1))$. For Gaussian stationary time series, it can be readily computed as

$$h = \frac{1}{4\pi} \int_0^{2\pi} \log \det [4\pi^2 e f(\omega)] d\omega \quad (11)$$

which is a result known as Szëgo's theorem [18].

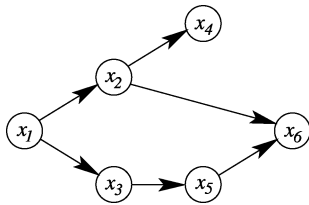


Fig. 3. Examples of graphical model and conditional independence. x_1 is independent from x_6 , given x_2, x_3 ; x_4 is independent from x_6 , given x_2 .

Note that h is also equal to the conditional entropy of $x(0)$ given the infinite past [19], that is, if G is the covariance matrix of $x(0)$ given the infinite past, we have

$$h = \lim_{T \rightarrow \infty} H(x(0)|x(-1), \dots, x(-T)) = \frac{1}{2} \log \det 2\pi e G.$$

B. Directed Graphical Models for Gaussian Time Series

In this paper, we consider directed graphical model representations for time series. This section reviews the semantics of directed graphical models [3] to make the paper self-contained. We first review the classical results for Gaussian vectors and covariance matrices and present extensions of these to Gaussian stationary time series and spectral densities.

1) *Directed Models and Conditional Independence:* Let x be a Gaussian variable with covariance matrix Σ . The variable x , or the covariance matrix Σ , is said to factorize in a directed acyclic graph G if and only if for all i , x_i is independent from its nondescendant variables, given its parents [3]. See Fig. 3 for a simple example. For a characterization in terms of the covariance matrix, see Section III-B3.

We generalize this definition to stationary time series: The time series $x(t)$, with spectral density $f(\omega)$, is said to factorize in a directed acyclic graph G , if for all $\omega \in [0, 2\pi]$, the (complex) covariance matrix $f(\omega)$ factorizes in G . This is exactly equivalent to the following property: $x(t)$ factorizes in G if and only if, for all i , x_i is independent from its nondescendant variables given its parents, where conditional independence between time series should be understood as in the previous section.

2) *Factorization of the KL Divergence:* Since the KL divergence decomposes in a directed acyclic graphical model for Gaussian variables [3], the previous propositions immediately imply that the KL divergence decomposes in directed acyclic models for Gaussian stationary time series. Thus, parameter estimation in a directed graphical model with m time series decouples into m distinct conditional density estimates for univariate time series. This makes possible the efficient learning of the structure of a directed graphical model for time series from data, as we show in Section V.

3) *Sparse Representation of the Covariance and Precision Matrices:* In this section, we show how matrix-vector products $\Sigma\alpha$ and $\Sigma^{-1}\alpha$ can be computed efficiently when a covariance matrix Σ factorizes in a directed graphical model. This is easily seen from the “regression view” of Gaussian directed graphical models [20], which implies that Σ^{-1} and Σ can be written as

$$\begin{aligned} \Sigma^{-1} &= (I - W)^* D^{-1} (I - W) \\ \Sigma &= (I - W)^{-1} D (I - W)^{-*} \end{aligned}$$

where W is a strictly lower triangular matrix (i.e. with zero diagonal), and D is diagonal. In this interpretation, W are the regression weights and D the conditional variances. In other words, a distribution that factorizes in G is defined by the parameters W and D . Note that these parameters can be easily found from the full joint covariance matrix as follows:

$$\begin{aligned} W_i &= \Sigma_{i, \pi_i} \Sigma_{\pi_i, \pi_i}^{-1} \\ d_i &= \Sigma_{i, i} - \Sigma_{i, \pi_i} \Sigma_{\pi_i, \pi_i}^{-1} \Sigma_{\pi_i, i} \end{aligned}$$

where π_i is the set of parents of node i in G , and for any two subsets A, B of $\{1, \dots, m\}$, of cardinal $|A|$ and $|B|$, $\Sigma_{A, B}$ denotes the $|A| \times |B|$ matrix extracted from Σ by keeping row indices in A and column indices in B .

If the maximal fan-in of the graph G is d , then for any vector α , we can compute $\Sigma^{-1}\alpha$ in $O(dm)$ operations because the product $W\alpha$ can be computed in $O(dm)$ operations. In order to be able to compute $\Sigma\alpha$, we just need to be able to solve the systems of the form $(I - W)x = \alpha$, which can be done in $O(dm)$ operations, since the matrix W is triangular and has at most dm nonzero elements.

IV. EFFICIENT FORECASTING WITH DIRECTED GRAPHICAL MODELS

In this section, we describe a novel algorithm for efficient forecasting when the spectral density matrix has a sparse structure. We learn the one-step-ahead predictor from the potentially infinite past, that is, we learn the predictor from the H previous time steps, where the required “horizon” H is determined by the smoothness of the spectral density. We essentially need to solve the Yule–Walker equation of order H . Direct methods such as the Durbin–Levinson algorithm do not scale well with the number of variables m . However, by solving the linear system using the conjugate gradient technique, we can make use of the structure of our spectral density.

A. Problem Setup

We assume that the spectral density is known through H samples f_k at frequencies $\omega_k = 2k\pi/H$ and that H is large enough so that the autocovariance function is equal to the discrete Fourier transform of the sequence f_k . In addition, we assume that for each k , f_k has a *sparse structure*. By sparse structure, we mean that for each vector $\alpha \in \mathbb{R}^m$, the product $f_k\alpha$ and $f_k^{-1}\alpha$ can be computed in $O(sm)$ operations, where s is a constant independent of m instead of $O(m^2)$. A first example is when f_k factorizes in a directed graphical model with maximal fan-in that is less than d , where we have $s = O(d)$, as shown in Section III-B3. A second example is where the covariance $f = f_k$ follows a factor analysis model, that is, $f = WW^T + \psi$, where W is an $m \times p$ matrix, and ψ is diagonal. In that case, it is easy to show that $s = O(p^3)$ (the argument is trivial for $f\alpha$, and it follows by the matrix inversion lemma for $f^{-1}\alpha$).

In this section, our interest is the design of an algorithm to determine the one-step-ahead predictor from the last H steps; therefore, we need to solve the following Yule–Walker equations $Q_H(\Gamma)\Psi = \gamma_H$, as defined in (9). We assume that the values $\Gamma(h)$, for $h = 0, \dots, H/2 - 1$, are obtained as

the first $H/2$ values of the FFT of order H of the sequence (f_0, \dots, f_{H-1}) , whereas $\Gamma(h) = 0$ for $h \geq H/2$.

B. Conjugate Gradient Methods and Toeplitz Systems

Conjugate gradient methods can be used to solve large linear systems of equations $Ax = b$ of size n , where the matrix A is Hermitian positive and where the evaluation of the product Ax can be performed in $O(n)$ operations [instead of $O(n^2)$]. It is an iterative method where at each iteration, one matrix-vector product has to be performed. The number of iterations is governed by the condition number of the matrix A —the ratio between its largest eigenvalue and its smallest eigenvalue. Unfortunately, in many cases, A is ill-conditioned, and the number of iterations can be very large. A common solution is to precondition the matrix A ; instead of solving $Ax = b$, the system $C^{-1}AC^{-1}y = C^{-1}b$ is solved, which yields a solution for the original system through $x = C^{-1}y$. In order to make conditioning practical, there are two requirements for the matrix C : The linear system involving C^2 must be solved cheaply, and the condition number of $C^{-1}AC^{-1}$ must be small [21], i.e., $C^{-1}AC^{-1}$ has to be as close as possible to a multiple of the identity matrix.

For Toeplitz matrices, there is a natural choice of C obtained through the approximate diagonalization of Toeplitz matrices in the Fourier basis [22]. Indeed, any block symmetric Toeplitz matrix of the form

$$T = \begin{pmatrix} T_0 & T_{-1} & \cdots & T_{-H+1} \\ T_1 & T_0 & & \\ \vdots & & \ddots & \\ T_{H-1} & \cdots & & T_0 \end{pmatrix}$$

can be approximated by a circulant matrix of the form

$$C = \begin{pmatrix} C_0 & C_{-1} & \cdots & C_{-H+1} \\ C_1 & C_0 & & \\ \vdots & & \ddots & \\ C_{H-1} & \cdots & & C_0 \end{pmatrix}$$

where $C_k = T_k$ for $|k| \leq H/2 - 1$, and $C_k = C_{n-k}$ (which defines the values for $|k| \geq H/2$).

The circulant approximation is very useful because it can be diagonalized in a fixed basis, that is

$$C = \mathcal{F}\mathcal{D}\mathcal{F}^*$$

where \mathcal{F} is defined by blocks, that is, $\mathcal{F}_{jk} = (1/\sqrt{H})\exp(2ijk\pi/H)I$, and \mathcal{D} is block diagonal with diagonal blocks $\hat{C}_k = (1/\sqrt{H})\sum_{j=0}^{H-1} C_k \exp(-2ijk\pi/H)$. Since circulant matrices can be diagonalized in the Fourier basis \mathcal{F} , they can be easily inverted as

$$C^{-1} = \mathcal{F}\mathcal{D}^{-1}\mathcal{F}^*.$$

For Toeplitz systems, it is natural to precondition the associated inverse circulant matrix [23], [24]. In our case, it is equivalent to being preconditioned by the circulant matrix generated by the sequence $(\bar{f}_k)^{-1}$. The preconditioning is efficient because products of the form $\bar{f}_k\alpha$ and $(\bar{f}_k)^{-1}$ can be computed in linear time in m . In Fig. 4, we give an outline of the resulting algorithm.

C. Computing ν -Step-Ahead Predictors

Using the concept of spectral factorization, we can compute any ν -step-ahead predictor and error covariances efficiently. The solution of the linear system in (8) is the optimal one-step-ahead filter, with transfer function $A_1(\omega) = \sum_{k=0}^{\infty} \Psi_{k+1} e^{-ik\omega}$. Let A_ν be the transfer function of the ν -step-ahead filter and G_ν its error. They can be computed using the following theorem [25], [26].

Theorem 1: The error covariance $G = G_1$ of the one-step-ahead predictor is

$$G = \int_0^{2\pi} (I - e^{-i\omega} A_1(\omega)) f(\omega) d\omega.$$

If $\psi(\omega) = I - e^{-i\omega} A_1(\omega)$ and $\phi(\omega) = \psi^{-1}(\omega) = I - \sum_{k=1}^{\infty} \phi_k e^{-ik\omega}$, the ν -step-ahead filter can be written as

$$A_\nu(\omega) = e^{i\nu\omega} \left(I - \sum_{k=0}^{\nu-1} \phi_k \phi^{-1}(\omega) \right)$$

and the error covariance is $G_\nu = \sum_{k=0}^{\nu-1} \phi_k G \phi_k^T$. ■

Thus, once the one-step-ahead filter is found, all other filters and errors can be computed. However, this requires inverting H matrices of size m . This can be avoided by noticing that $\phi(\omega)$ can be obtained by computing the one-step-ahead predictor corresponding to the spectral density $\bar{f}^{-1}(\omega)$. If m is very large, such that multiplying two full matrices of size m is prohibitive, we can obtain the ν -step-ahead predictor and error covariance by solving the Toeplitz system $Q_H \Psi_\nu = \gamma_{H,\nu}$, where $\gamma_{H,\nu} = (\Gamma(\nu), \dots, \Gamma(\nu + H - 1))$.

D. Comments

1) *Convergence Analysis:* Results from [27] show that the number of iterations is independent of the order H and depends on the smoothness of the spectral density as well as how far $f(\omega)$ is from singular.

2) *Interpretation in the Frequency Domain:* The previous algorithm has an interpretation in the frequency domain that links it to previous approaches [25], [26] that do not rely on efficient methods for linear systems.

3) *Total Complexity:* Let H be the number of samples required for the computation of the Fourier transforms, m the number of samples, and d the maximum fan-in. We thus have the following:

- each iteration: $O(m^2 H \log H)$ for the FFTs, $O(m^2 d H)$ for the multiplication by $f(\omega)$;
- overall complexity of computing the predictor: $O(m^2 H (\log H + d))$.

V. LEARNING-DIRECTED GRAPHICAL MODELS

In this section, we present our algorithm for learning-directed graphical models for stationary Gaussian time series. Not surprisingly, we make heavy use of the corresponding algorithms for the Gaussian temporally independent case. We cast the structure learning as a model selection problem where the model is defined by the directed graph G . We use the Akaike Information

Input: Spectral density values $f_k \in \mathbb{R}^{m \times m}$, at frequencies $\omega_k = 2k\pi/H$, $k \in \{0, \dots, H-1\}$, precision ε

Algorithm:

1. Computing covariances:

$$\begin{aligned} \forall h \in \{0, \dots, H-1\}, r(:, :, h) &\leftarrow \bar{f}(:, :, h) \\ \forall (i, j) \in \{1, \dots, m\}^2, r(i, j, :) &\leftarrow \text{IFFT}(r(i, j, :)) \\ r(i, j, :) &\leftarrow r(i, j, 0 : H/2-1) \end{aligned}$$

2. Initialization: $k \leftarrow 0$, $\gamma_1 \leftarrow 0$, $\rho \leftarrow \sum_{i,j,t} |r(i, j, t)|^2$

3. Compute Ψ : while $\sqrt{\rho} > \varepsilon$:

$$\begin{aligned} k &\leftarrow k + 1 \\ \forall (i, j) \in \{1, \dots, m\}^2, z(i, j, :) &\leftarrow [r(i, j, :); 0 \dots 0] \\ z(i, j, :) &\leftarrow \text{IFFT}(z(i, j, :)) \\ \forall h \in \{0, \dots, H-1\}, z(:, :, h) &\leftarrow \bar{f}(:, :, h)^{-1} z(:, :, h) \\ \forall (i, j) \in \{1, \dots, m\}^2, z(i, j, :) &\leftarrow \text{FFT}(z(i, j, :)) \\ z(i, j, :) &\leftarrow z(i, j, 0 : H/2-1) \end{aligned}$$

$$\begin{aligned} \gamma_0 &\leftarrow \gamma_1, \gamma_1 \leftarrow \sum_{i,j,t} \bar{r}(i, j, t) z(i, j, t) \\ \text{if } k = 1 \text{ then } p &= z \text{ else } \beta \leftarrow \gamma_1 / \gamma_0, p \leftarrow z + \beta p \\ \forall (i, j) \in \{1, \dots, m\}^2, w(i, j, :) &\leftarrow [p(i, j, :); 0 \dots 0] \\ w(i, j, :) &\leftarrow \text{IFFT}(w(i, j, :)) \\ \forall h \in \{0, \dots, H-1\}, w(:, :, h) &\leftarrow \bar{f}(:, :, h) w(:, :, h) \\ \forall (i, j) \in \{1, \dots, m\}^2, w(i, j, :) &\leftarrow \text{FFT}(w(i, j, :)) \\ w(i, j, :) &\leftarrow w(i, j, 0 : H/2-1) \end{aligned}$$

$$\begin{aligned} \alpha &\leftarrow \sum_{i,j,t} \bar{p}(i, j, t) w(i, j, t), \alpha \leftarrow \gamma_1 / \alpha \\ g &\leftarrow g + \alpha p, r \leftarrow r - \alpha w \\ \rho &\leftarrow \sum_{i,j,t} |r(i, j, t)|^2 \end{aligned}$$

4. Compute G and Ψ_i :

$$\begin{aligned} \forall h \in \{0, \dots, H/2-1\}, g(:, :, h) &\leftarrow g(:, :, h)^\top \\ \forall (i, j) \in \{1, \dots, m\}^2, w(i, j, :) &\leftarrow [0; g(i, j, :); 0 \dots 0] \\ w(i, j, :) &\leftarrow \text{FFT}(w(i, j, :)) \\ \forall h \in \{0, \dots, H-1\}, w(:, :, h) &\leftarrow (I - w(:, :, h)) \bar{f}(:, :, h) \\ G &\leftarrow \frac{1}{H} \sum_t w(:, :, t) \\ \forall h \in \{1, \dots, H/2\}, \Psi_h &= g(:, :, h-1) \end{aligned}$$

Output: $G, \Psi_h, h \in \{1, \dots, H/2\}$

Fig. 4. Solving the Yule–Walker equations using the spectral density.

Criterion (AIC) in order to define the score $J(G)$ to be minimized. Once the score is determined, the task of minimizing it is performed with greedy algorithms, as detailed in Section V-B.

A. AIC Score for Time Series

We are given T consecutive samples of m univariate times series $x_i(t)$, $i \in \{1, \dots, m\}$, and $t \in \{1, \dots, T\}$. We first estimate the joint spectral density matrix $\hat{f}(\omega)$ as well as the estimated degrees of freedom df , as shown in Section II-D. The spectral density is represented as a set of H samples f_k at frequencies $\omega_k = 2k\pi/H$.

For directed models, the AIC score $J(G)$ is equal to the maximum of the likelihood plus a penalty score. It is known to factorize [28], [29]

$$J(G) = \sum_{i=1}^m J_i(\pi_i(G)) \quad (12)$$

where $\pi_i = \pi_i(G)$ are the parents of node i in G , and

$$J_i(\pi_i(G)) = -T \hat{H}(x_i | x_{\pi_i}) + \#\{i, \pi_i(G)\}$$

Input: data $x_i(t)$, $i \in \{1, \dots, m\}$, $t \in \{1, \dots, T\}$

Algorithm:

1. Estimate the joint spectral density $f_k = \hat{f}(\omega_k)$ and the degree of freedom df using the algorithm in Figure 1, for $\omega_k = 2k\pi/H$
2. Minimization of the AIC score $J(G)$
 - a. Initialization: $G = \emptyset$
 - b. while no decrease $J(G)$, select best local moves: addition, deletion or reversal
3. Compute optimal smoothing r_i separately for each clique $\{i\} \cup \pi_i(G)$ using local smoothing
4. Compute the spectral density matrix f_k^G that factorizes in G using the local sufficient statistics

Output: G, f_k^G

Fig. 5. Learning graphical models for stationary Gaussian time series.

where $\#\{i, \pi_i(G)\}$ is the number of parameters required to encode the conditional probability distribution of x_i given the parents $x_{\pi_i(G)}$, and $\hat{H}(x|y)$ is conditional entropy of the random vector x , given the random vector y , computed using the estimated distribution of the joint (spectral) density.

Using the expression of the KL divergence and entropy rate in (10) and (11), we get the following expression for the local score (where we omit the terms that are independent of the choice of the graph G):

$$J_i(\pi_i) = \frac{-T}{4\pi} \int_0^{2\pi} \log \frac{|\hat{f}_{\{i\} \cup \pi_i}(\omega)|}{|\hat{f}_{\pi_i}(\omega)|} d\omega + (2|\pi_i| + 1) \frac{df}{2} \quad (13)$$

where $|\pi_i|$ is the cardinality of π_i (i.e., the number of parents), and $f_A(\omega)$ denotes the square block (A, A) of the $m \times m$ matrix $f(\omega)$. The previous local score is approximated using the samples of $\hat{f}(\omega)$ as

$$J_i(\pi_i) = \frac{-T}{2H} \sum_{k=0}^{H-1} \log \frac{|(f_k)_{\{i\} \cup \pi_i}|}{|(f_k)_{\pi_i}|} + (2|\pi_i| + 1) \frac{df}{2}. \quad (14)$$

B. Learning Algorithm

In order to learn the graph structure G , we need to minimize the score $J(G)$ defined by (12) and (13). This minimization problem is known to be NP-complete [30], and greedy algorithms are commonly used. In particular, since the score $J(G)$ factorizes as a sum of local scores, hillclimbing search using local moves—edge addition, deletion or removal—is computationally efficient (in particular through the caching of already computed local scores) and, with the possible use of random restarts, yields good local minima [28], [29]. An outline of the final algorithm is presented in Fig. 5. The exact complexity of the search procedure that uses caching of local scores is presented in Section V-D.

C. Local Smoothing and Efficiency

One of the major gains from learning a sparse structure for the spectral density matrix is that we can perform (and optimize) the smoothing of the periodogram locally in the graph, restricting

ourselves to elements that share a clique, that is, instead of applying the algorithm of Fig. 1 once with m variables, we can apply it m times with $|\pi_i(G)| + 1$ variables $i = 1, \dots, m$. In the algorithm presented in Fig. 5, a joint pilot estimate of the spectral density is used to determine the graph G ; then, the local smoothing is performed. We now explore some of the issues, both numerical and statistical, associated with such smoothing.

Numerically, learning the best smoothing parameter for the joint spectral density is an $O(m^3 T \log T)$ operation. In order to overcome the cubic time complexity in m , it is possible to learn a different smoothing parameter r_i for each local potential that is required, which makes the algorithm only $O(d^3 T \log T)$, where d is the maximal fan-in of the directed graph. Note that in the algorithm of Fig. 5, it is possible to avoid the cubic time complexity of the computation of the pilot estimate simply by considering a bandwidth that is the mean of the optimal bandwidth for each of the variables taken separately.

Statistically, learning a smoothing parameter for the joint density when there is a strong local structure would lead to oversmoothing (since the AIC penalty becomes relatively too important): Local smoothing is substantially more efficient. In the algorithm presented in Fig. 5, the local smoothing is performed once the structure is learned, thus requiring m distinct smoothing parameter searches.

Finally, alternating between learning local smoothing parameters and directed graphs is possible in a manner analogous to the procedure of [29]. In that situation, the global score to be optimized is the concatenation of (6), (12), and (13). More precisely, if r_i denotes the smoothing parameters for the clique $\{i\} \cup \pi_i(G)$ and \hat{f}^{r_i} denotes the smoothed periodogram with parameter r_i , we need to compute the Whittle likelihood and then add the number of parameters, as shown in the next theorem.

Theorem 2: The Whittle likelihood for a spectral density that factorizes in a directed graph G , where the conditional spectral densities are computed using a smoothing constant r_i and with resulting estimates $\hat{f}_{\{i\} \cup \pi_i}^{r_i}(\omega)$ for the cliques $\{i\} \cup \pi_i$, is approximately equal to

$$\ell_W(G, r) = \frac{T}{2H} \sum_{k=0}^{H-1} \sum_{i=1}^m \left(\log \frac{|(f_k^{r_i})_{\{i\} \cup \pi_i}|}{|(f_k^{r_i})_{\pi_i}|} + \text{tr} \left\{ (f_k^{r_i})_{\{i\} \cup \pi_i}^{-1} I_{\{i\} \cup \pi_i}(\omega_k) \right\} - \text{tr} \left\{ (f_k^{r_i})_{\pi_i}^{-1} I_{\pi_i}(\omega_k) \right\} \right). \quad \blacksquare$$

The overall AIC score is thus equal to

$$J(G, r) = -\ell_W(G, r) + \sum_{i=1}^m (2|\pi_i| + 1) \frac{df_i}{2}. \quad (15)$$

Note that the optimal local smoothing that results from optimizing r_i independently from G and other smoothing constants r_j , $j \neq i$ is different from the optimal smoothing on the clique $\pi_i \cup \{i\}$. In simulations on a wide variety of examples, it turns out that the difference between the two types of smoothing is very small, with a slight advantage for the conditional smoothing, which makes the density estimate more robust (see Section VII-A).

D. Running Time Complexity

Denoting by m the number of these variables, d the maximum fan-in that we impose on our networks, T the number of observations, and H the time horizon, we can compute the running time complexity of our learning algorithm: As seen in Fig. 5, the structure learning has several successive stages, whose running time complexities are the following.

- 1) Estimate (marginal) smoothing parameters r_i for each of the m variables. Cost: $O(mT^{3/2} \log T)$.
- 2) Compute/smooth the periodogram for the joint density. Cost: $O(m^2 T \log T)$.
- 3) Structure learning using greedy search: When using the efficient scheme of [31], it can be made to be $O(m^3 d^2 + m^2 d^4 H)$.
- 4) Estimate local smoothing parameters. Cost: $O(md^3 T^{3/2} \log T)$.
- 5) The overall cost is $O(md^3 T^{3/2} \log T + m^2 T \log T + m^3 d^2 + m^2 d^4 H)$.

VI. NONLINEARITY THROUGH MERCER KERNELS

In this section, we show how the methods presented thus far can be extended to nonlinear models. In particular, in Section VI-B, we show how these methods can be “kernelized”; based on this, we develop algorithms for learning in Section VI-C and prediction in Section VI-D. As already mentioned, the basic principle underlying our approach is to first map the data x into a “feature space” \mathcal{F} using a “feature map” $\Phi(x)$ and to use the algorithms developed in earlier sections on the transformed data $\Phi(x)$. What makes the approach feasible is that only dot products between data points are needed by these algorithms, as seen in Section VI-B and C. Thus, the algorithm only involves manipulation of the values $k(x, y) = \langle \Phi(x), \Phi(y) \rangle$ for x and y because they are two data points in the original space. $k(x, y)$ is usually referred to as the *kernel function*. A function $k(x, y)$ is a Mercer kernel if and if only there exists a feature space \mathcal{F} and a feature map $\Phi(x)$ such that $k(x, y) = \langle \Phi(x), \Phi(y) \rangle$.

Since the algorithm only involves dot products and, thus, kernel values, algorithms can benefit from the expressive power of the feature space, without incurring the computational cost of actually mapping the points into this feature space, which is a method usually referred to as the “kernel trick” [14]. Mercer kernels can be defined on many different input spaces, thus making the applicability of the presented technique wide.

A. Notation

We assume that the variables $x_i(t)$ are mapped to $\varphi_i(t) = \Phi_i(x_i(t)) \in \mathcal{F}_i$, where Φ_i are given feature maps from the “input space” \mathcal{X}_i to the feature space \mathcal{F}_i . Let $k_i(x_i, y_i) = \langle \Phi_i(x_i), \Phi_i(y_i) \rangle$ be the dot product between two mapped elements. We are given T samples from $t = 0$ to $t = T - 1$. We can build m “Gram matrices” K_i (one per variable), which are defined as $T \times T$ matrices composed of the pairwise kernel values $k_i(x_i(u), x_i(v))$ between two data points.

B. Kernelization of the Periodogram

We now show how we can find a basis in which the periodogram has a particularly simple expression involving products of Gram matrices. Since the model selection scores are independent of the basis, we can use these expressions in Section VI-C to construct our model selection criteria.

We assume that the data are centered,² that is, $\sum_t \varphi_i(t) = 0$. The sample covariance matrix $\hat{\Gamma}(h)$, which is defined by blocks, has its (i, j) -block equal to, for $h \geq 0$

$$\hat{\Gamma}_{ij}(h) = \frac{1}{T} \sum_{t=h}^{T-1} \varphi_i(t) \varphi_j(t-h)^\top.$$

For $s, u \in \{0, \dots, T-1\}$, we have

$$\begin{aligned} \varphi_i(s)^\top \hat{\Gamma}_{ij}(h) \varphi_j(u) &= \frac{1}{T} \sum_{t=h}^{T-1} \varphi_i(s)^\top \varphi_i(t) \varphi_j(t-h)^\top \varphi_j(u) \\ &= \frac{1}{T} \sum_{t=h}^{T-1} (K_i)_{s,t} (K_i)_{t-h,u} \\ &= \frac{1}{T} \delta_s K_i J_h K_j^\top \delta_u \end{aligned}$$

where δ_s is the vector in \mathbb{R}^T with only zeros, except for a one at index s , K_i is the Gram matrix associated with variable i , and J_h is the $T \times T$ matrix such that $(J_h)_{ab} = \delta(a-b-h)$. The only nonzero elements of J_h belong to the h th lower diagonal. Thus, in the *data basis* defined by the data points, the autocovariance function has its (i, j) th block equal to $(1/T)K_i J_h K_j$. The Fourier transform of the autocovariance has the following expression:

$$\begin{aligned} I_{ij}(\omega) &= \sum_{h=-(T-1)}^{T-1} \Gamma_{ij}(h) e^{-ih\omega} \\ &= \frac{1}{T} \sum_{h=-(T-1)}^{T-1} K_i J_h e^{-ih\omega} K_j^\top \\ &= \frac{1}{T} K_i \left(\sum_{h=-(T-1)}^{T-1} J_h e^{-ih\omega} \right) K_j^\top \\ &= \frac{1}{T} K_i f_\omega f_\omega^* K_j^\top \\ &= \frac{1}{T} K_i f_\omega (K_i f_\omega)^* \end{aligned}$$

where f_ω is the vector in \mathbb{R}^T with components $e^{-ih\omega}$, $h \in \{0, \dots, T-1\}$. The periodogram is exactly $I_{ij}(\omega_k)$ for $\omega_k = 2k\pi/T$.

Since the feature space might have infinite dimension, smoothing (in space) is usually required to limit the number of parameters that are implicitly estimated. Smoothing (in space) the periodogram by an isotropic Gaussian is equivalent to adding κI ; in the data basis, the regularized (cross)-periodogram is thus

$$I_{ij}(\omega) = \frac{1}{T} K_i f_\omega f_\omega^* K_j^\top + \kappa \delta_{ij} K_i.$$

²The data can be implicitly centered in feature space by replacing Gram matrices K by $(I - (1/T)\mathbf{1})K(I - (1/T)\mathbf{1})$, where $\mathbf{1}$ is a $T \times T$ matrix composed of ones [14].

Smoothing (in frequency) can be done as in the nonkernelized version by averaging consecutive values of the periodogram. We denote $(f_k)_{ij}$, $k = 1, \dots, H$ as the values of the estimated spectral density after smoothing. Each $(f_k)_{ij}$ is a $T \times T$ matrix.

C. Model Selection Criteria

Following [13], the effective dimension of variables i is taken to be $d_i = \text{tr}(K_i + \kappa I)^{-1} K_i$, and we use the following AIC score:

$$J(G) = \sum_{i=1}^m \left\{ \frac{-T}{2H} \sum_{k=0}^{H-1} \log \frac{|(f_k)_{\{i\} \cup \pi_i, \{i\} \cup \pi_i}|}{|(f_k)_{\pi_i, \pi_i}|} + \left(2 \sum_{j \in \pi_i} d_j + 1 \right) \frac{df}{2} \right\}.$$

The same optimization techniques used for classical graphical models can be used to optimize this score. Efficient implementation techniques based on low-rank approximations are presented in [13].

D. Prediction with Kernels

Once we have a model, we would like to perform prediction. Given our framework, we obtain a predicted value $\hat{\phi} = (\hat{\phi}_1, \dots, \hat{\phi}_m)$ in the feature space, and this predicted value might not correspond to an x_i such that $\hat{\phi}_i = \Phi_i(x_i)$. As for the covariance error, it is not defined in input space. In order to get an estimate and an error, we propose the following two-stage procedure. Once the predicting filters $\hat{\phi}_i(t) = f_i(\phi(t-1), \dots, \phi(t-H))$ are computed, do the following.

- 1) For each variable i and each index $t \in [H, T-1]$, find

$$\hat{x}_i(t) = \arg \min_{x_i} \|\Phi_i(x_i) - f_i(\phi(t-1), \dots, \phi(t-H))\|.$$

This is usually referred to as the “pre-image,” and several techniques are available [14], [32].

- 2) Compute the average errors. In particular, if all variables are continuous, we get

$$G = \frac{1}{T-H} \sum_{t \in [H, T-1]} (\hat{x}(t) - x(t)) (\hat{x}(t) - x(t))^\top.$$

VII. SIMULATIONS

A. Synthetic Examples

Our first goal is to see whether, when the data are generated from a sparse graphical model, the search procedure can find this model. We use the following procedure: We generate a random spectral density matrix, project the spectral density matrix function to a random graph with maximal fan-in equal to two, generate data from the new spectral density matrix with various numbers of samples, and learn the model from data using the algorithm described in Fig. 5. We then compute the KL divergence from the generating model as well as the number of undirected edges not found by the learning algorithm. We compare the performance in KL divergence for the following models: a fully connected graphical model, a learned graphical model without post-smoothing, and a learned graphical model

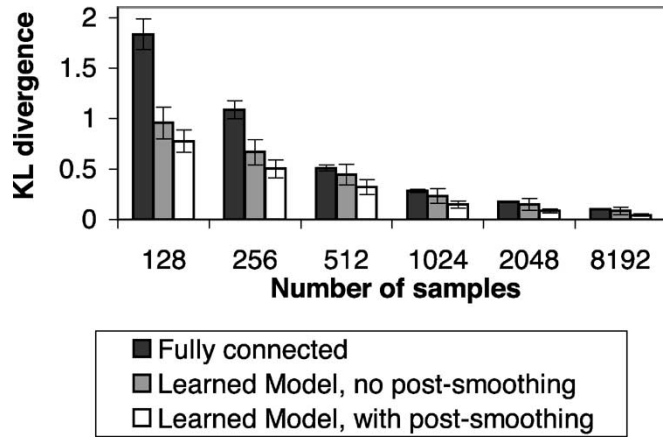


Fig. 6. Synthetic examples. KL divergence from the truth versus the number of samples.

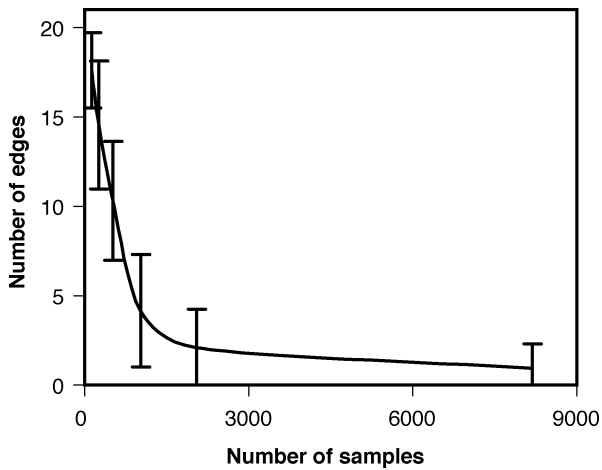


Fig. 7. Synthetic examples. Number of nonrecovered edges versus number of samples for the structure learning algorithm presented in Section V-B.

with post-smoothing (as described in Section V-C). We can see in Fig. 6 that the structure learning algorithm manages to predict significantly better than the fully connected graphical model, especially for small samples (when the sample gets large, the KL divergences for all methods go to zero; see also Fig. 7). In addition, the post-smoothing leads to a slight improvement; in particular, the variance of the error is reduced.

B. Real-Life Datasets

We compare our algorithms with a classical approach that exhibits reasonable scaling of computing time to large datasets. Note, in particular, that with large datasets such as climate data, two issues need to be addressed: a) There are a very large number of variables m , and b) there are relatively few time samples compared to the number of variables. The model we compare to is the sparse autoregressive model SAR(p) of maximal order p . These are AR(p) models that favor zeros in the weights Φ_i and are such that the error matrix Σ factorizes in a sparse graphical model G . They are learned using forward selection of edges with the AIC criterion.

We used climate datasets extracted from the National Climatic Data Center (NCDC) “summary of the day” data, which are composed of daily mean temperatures and precipitation for

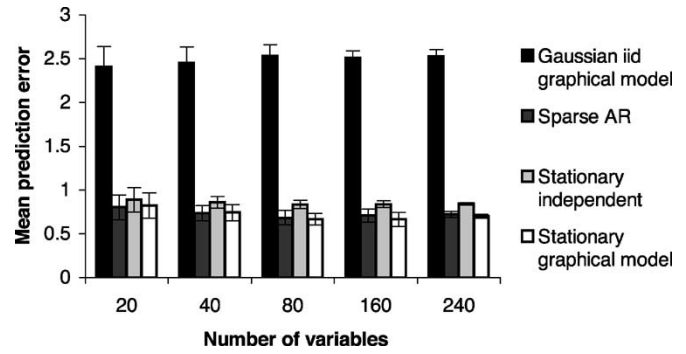


Fig. 8. Results for the climate datasets: mean prediction error.

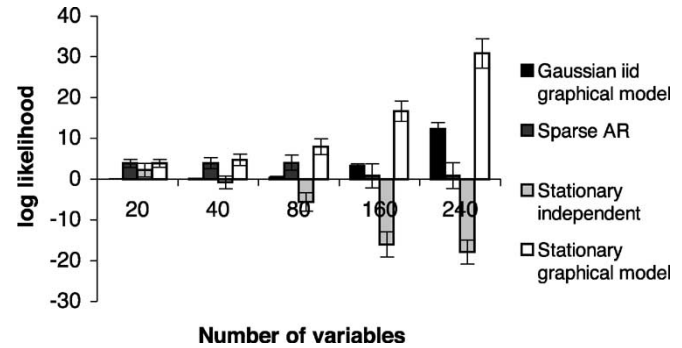


Fig. 9. Results for the climate datasets. Predictive log likelihood normalized by the log likelihood for the full Gaussian iid model (larger values are better).

more than 10 000 stations across the United States. We subsampled this dataset by choosing random locations and choosing the closest stations. The number of samples was 1024 consecutive days. Each of the datasets that was so constructed was divided in a training set (of size 1024 days) and a testing set. For each of these sets, the obvious seasonal component (year) is removed by using the first six corresponding Fourier components. We report the negative log likelihood of the one-step-ahead prediction of the test data. We normalize the result by subtracting the log likelihood of the corresponding iid Gaussian variables with a full covariance matrix. We also report the average prediction error. We report results in Figs. 8 and 9. We can see in Fig. 8 that on the mean prediction error, the sparse methods, in time or frequency, perform better than the methods that do not use sparsity; in terms of log likelihood, the sparse frequency domain approach outperforms the time domain approach.

VIII. CONCLUSIONS

Probabilistic graphical models provide an elegant general framework for the representation of complex sets of dependencies among an interacting set of variables. Standard algorithms are available for probabilistic inference, and a variety of parameter estimation and model selection algorithms have been devised. The graphical structure of these models is essential for making these algorithms computationally efficient. It has important statistical implications as well, yielding models in which the graphical structure corresponds directly to a natural notion of sparsity of the representation.

Applications of graphical models to time series analysis have generally taken the form of state space models. In this setting,

the graphical model machinery is aimed at capturing structure in the state transition matrix, and sparsity in the graph has an interpretation in the time domain—a given state variable has a probabilistic dependence on a limited number of variables in the past.

In the current paper, we have described an alternative methodology for making use of graphical models in time series analysis. In particular, we have developed a frequency domain approach in which the structure captured by a graphical model is related to sparseness in the spectral density matrix. We have described parameter estimation and structure learning algorithms that are geared for this setting. We have provided computational complexity analyses throughout, emphasizing the development of methods that are appropriate for large-scale time series models.

As seen in the experiments with climate data, methods that attempt to uncover structure in the frequency domain can lead to equivalent predictive performance to analogous methods operating in the time domain but higher likelihoods. The frequency approach is looking for conditional independence relationships among the variables and should work well when such relationships are expected.

Finally, it is worth noting that although we have restricted ourselves to regularly sampled time series in this paper, our algorithms apply immediately to irregularly sampled time series once the spectral density is estimated [33].

REFERENCES

- [1] P. Bloomfield, *Fourier Analysis of Time Series: An Introduction*. New York: Wiley, 2000.
- [2] P. J. Brockwell and R. A. Davis, *Time Series: Theory and Methods*. New York: Springer-Verlag, 1991.
- [3] S. L. Lauritzen, *Graphical Models*. London, U.K.: Clarendon, 1996.
- [4] M. I. Jordan, "Graphical models," *Stat. Sci.* (Special Issue on Bayesian Statistics), 2003, to be published.
- [5] T. J. Richardson and R. L. Urbanke, "The capacity of low-density parity-check codes under message-passing decoding," *IEEE Trans. Inform. Theory*, vol. 47, pp. 599–618, Apr. 2001.
- [6] T. Dean and K. Kanazawa, "A model for reasoning about persistence and causation," *Comput. Intel.*, vol. 5, no. 3, pp. 142–150, 1989.
- [7] K. Murphy, "Dynamic Bayesian Networks: Representation, Inference and Learning," Ph.D. dissertation, Comput. Sci. Div., Univ. Calif. Berkeley, Berkeley, CA, 2002.
- [8] D. R. Brillinger, "Remarks concerning graphical models for time series and point processes," *Rev. Econ.*, vol. 16, pp. 1–23, 1996.
- [9] R. Dahlhaus, "Graphical interaction models for multivariate time series," *Metrika*, vol. 51, pp. 157–172, 2000.
- [10] D. Geiger and D. Heckerman, "Learning Gaussian networks," in *Proc. Uncertainty Artificial Intelligence*, 1994.
- [11] W. Lam and F. Bacchus, "Learning Bayesian belief networks: an approach based on the MDL principle," *Comput. Intel.*, vol. 10, no. 4, pp. 269–293, 1994.
- [12] D. R. Brillinger, *Time Series: Data Analysis and Theory*. Philadelphia, PA: SIAM, 2001.
- [13] F. R. Bach and M. I. Jordan, "Learning graphical models with Mercer kernels," in *Advances Neural Inform. Process. Syst.*, vol. 15, 2003.
- [14] B. Schölkopf and A. J. Smola, *Learning With Kernels*. Cambridge, MA: MIT Press, 2001.
- [15] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*. New York: Springer-Verlag, 2001.
- [16] D. Kazakos and P. Papantoni-Kazakos, "Spectral distances between Gaussian processes," *IEEE Trans. Automat. Contr.*, vol. AC-25, pp. 950–959, 1980.

- [17] T. P. Speed and H. T. Kiiveri, "Gaussian Markov distributions over finite graphs," *Ann. Stat.*, vol. 14, no. 1, pp. 138–150, 1986.
- [18] E. J. Hannan, *Multiple Time Series*. New York: Wiley, 1970.
- [19] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York: Wiley, 1991.
- [20] R. Shachter and R. Kenley, "Gaussian influence diagrams," *Manage. Sci.*, vol. 35, no. 5, pp. 527–550, 1989.
- [21] G. H. Golub and C. F. Van Loan, *Matrix Computations*. Baltimore, MD: Johns Hopkins Univ. Press, 1996.
- [22] R. M. Gray, "Toeplitz and Circulant Matrices: A Review, Tech. Rep.," Inform. Syst. Lab., Stanford Univ., Stanford, CA, 2002.
- [23] R. H. Chan, "Toeplitz preconditioners for Toeplitz systems with non-negative generating functions," *IMA J. Numer. Anal.*, vol. 11, no. 3, pp. 333–345, 1991.
- [24] R. H. Chan, J. G. Nagy, and R. J. Plemmons, "FFT-based preconditioners for Toeplitz-block least squares problems," *SIAM J. Numer. Anal.*, vol. 30, no. 6, pp. 1740–1768, 1993.
- [25] N. Wiener and P. Masani, "The prediction theory of multivariate stochastic processes. I. The regularity conditions," *Acta Math.*, vol. 98, pp. 111–150, 1957.
- [26] —, "The prediction theory of multivariate stochastic processes. II. The linear predictor," *Acta Math.*, vol. 99, pp. 93–137, 1958.
- [27] R. H. Chan and M. K. Ng, "Conjugate gradient methods for Toeplitz systems," *SIAM Rev.*, vol. 38, no. 3, pp. 427–482, 1996.
- [28] D. Heckerman, D. Geiger, and D. M. Chickering, "Learning Bayesian networks: the combination of knowledge and statistical data," *Mach. Learning*, vol. 20, no. 3, pp. 197–243, 1995.
- [29] N. Friedman and M. Goldszmidt, "Learning Bayesian networks with local structure," in *Learning in Graphical Models*. Cambridge, MA: MIT Press, 1998.
- [30] D. M. Chickering, "Learning Bayesian networks is NP-complete," in *Learning from Data: Artificial Intelligence and Statistics 5*. New York: Springer-Verlag, 1996.
- [31] P. Giudici and R. Castelo, "Improving Markov chain Monte-Carlo model search for data mining," *Machine Learning*, vol. 50, pp. 127–158, 2003.
- [32] J. Weston, O. Chapelle, A. Elisseeff, B. Schölkopf, and V. Vapnik, "Kernel dependency estimation," in *Adv. NIPS*, vol. 15, 2003.
- [33] E. Parzen, *Time Series Analysis of Irregularly Observed Data*. New York: Springer-Verlag, 1983.



Francis R. Bach graduated from the Ecole Polytechnic, Palaiseau, France, in 1997 and received the Masters degree in mathematics from the Ecole Normale Supérieure, Cachan, France, in 2000. He is currently pursuing the Ph.D. degree at the Computer Science Division, University of California at Berkeley.

His research interests include machine learning, graphical models, kernel methods, and statistical signal processing.



Michael I. Jordan received the Masters degree from Arizona State University, Tempe, and the Ph.D. degree from the University of California at San Diego, La Jolla.

He is Professor with the Department of Electrical Engineering and Computer Science and the Department of Statistics, University of California at Berkeley. He was a professor at the Massachusetts Institute of Technology, Cambridge, from 1988 to 1998. In recent years, his research has focused on topics at the interface of statistics and computation,

including probabilistic graphical models, kernel machines, and applications of statistical machine learning to problems in bioinformatics, information retrieval, and signal processing.