

CONVERGENCE ANALYSIS OF THE GRAPH-TOPOLOGY-INFERENCE KERNEL LMS ALGORITHM

Mircea Moscu* Ricardo Borsoi*,† Cédric Richard*

* Université Côte d'Azur, CNRS, OCA, Nice, France

† Federal University of Santa Catarina, Florianópolis, Brazil

ABSTRACT

Identifying directed connectivity patterns from nodal measurements is an important problem in network analysis. Recent works proposed to leverage the performance and flexibility of strategies operating in reproducing kernel Hilbert spaces (RKHS) to model nonlinear interactions between network agents. Moreover, several applications require online and efficient solutions, which motivated the consideration of distributed adaptive learning strategies inspired by algorithms such as the kernel least mean square (KLMS). Despite showing good performance, a thorough theoretical understanding of the behavior of such algorithms is still missing. This makes applying them in practice challenging, especially because the set-up of adaptive algorithms involves additional parameters like the step size and a dictionary of kernel functions. In this paper, we present a convergence analysis of the graph-topology-inference KLMS algorithm. Monte Carlo simulations demonstrate the accuracy of the theoretical models.

Index Terms— convergence analysis, kernel least mean squares, topology inference, online processing

1. INTRODUCTION

Graphs have proven to be a fundamental tool in several applications, such as the analysis of socio-economical interactions [1] or brain activity [2]. In this context, information on the network structure is important and required by traditional graph signal processing algorithms [3]. It however appears that the topology of the graph is usually unknown beforehand. Recent works have addressed this issue by estimating the network topology directly from a set of measurements acquired at each node of the graph [4]. Although most topology identification algorithms assume linear dependencies between the nodal signals, nonlinear interactions between measurements at different nodes can be observed in many real-world applications, e.g., brain activity [2]. Many big data-oriented applications also require algorithms which are distributed, scalable, and can operate online with adaptive capabilities.

Although several solutions have been proposed for nonlinear online topology identification, a theoretical evaluation of the performance of such methods remain a challenging but important issue. These analyses are important both in understanding and designing the algorithm. In this paper, we present a theoretical model for the convergence behavior graph-topology-inference KLMS algorithm.

Background and prior work: Many approaches have been proposed for topology identification in the presence of linear interactions. Such works include inverse covariance estimation based on

the graphical Lasso [5], or the recovery of the topology from spectral templates [6]. More recently, online strategies have been proposed to address this problem, including diffusion-based [7] and distributed [8] adaptive algorithms, and variable-splitting based methods for stationary data [9]. However, these works do not consider nonlinear interactions that may happen.

Several strategies have been proposed to address the topology inference problem with nonlinear interactions, such as polynomial structural equation models [10, 11] and their nonlinear counterparts [12]. However, these methods require the nonlinear interaction models to be specified *a priori*. Kernel-based strategies, on the other hand, offer an efficient means of representing nonlinear functions by selecting an appropriate RKHS associated to a positive definite reproducing kernel [13]. The flexibility, elegance and mathematical tractability of kernel-based methods made their use widespread in topology identification. An in-depth overview of kernel-based topology inference can be found in [4], where several solutions are discussed. For instance, in [14], reproducing kernels are used to represent nodal measurements acquired at different time instants. An auto-regressive framework is employed to model graph connectivity over time, and a variable splitting-based batch algorithm is proposed to estimate the topology. A multi-kernel is also proposed in [15], where the graph topology is encoded via partial correlations. Despite performing well when the nonlinearities in the data are unknown, these strategies are not adaptive and their computational cost can be high.

An important aspect of kernel-based machinery is that it allows one to develop strategies that can address the online and adaptive topology identification problem in nonlinear settings. In [16], a batch strategy is modified to operate in an online fashion by considering a moving window with a fixed amount of samples. A stochastic gradient-based solution using a multi-kernel model was later proposed in [17] by considering a coherence-based online dictionary construction [18]. This method provides a competitive performance for estimating brain connectivity at a very small computational complexity. These models, however, are based on additive nonlinear interactions. Recently, a strategy was proposed to consider more general nonlinearities in the data [19] by employing a nonparametric topology identification framework.

However, to this date no theoretical evaluation of the performance of the online topology estimation approaches is available. This is an important issue since adaptive algorithms require the careful design of, e.g., a dictionary construction strategy [18] and an optimized step size in order to obtain a satisfactory performance. In this work, we provide a statistical performance analysis of the graph-topology-inference KLMS with a Gaussian kernel and a preselected dictionary. The presented model characterizes the mean and mean-squared behavior of the algorithm, both in transient and steady-state. Although previous works have analyzed the behavior of the KLMS

The work of C. Richard was funded in part by the ANR under grant ANR-19-CE48-0002, and by the 3IA Côte d'Azur Senior Chair program.

algorithm [20, 21, 22], we consider a distributed setting, with separate kernels and dictionaries for each node, and a temporal dependence in the input data, which allows the algorithm to have memory and make use of past data. This last characteristic is important in applications such as functional brain topology estimation, where there is a 10–20 ms delay in signal propagation between nodes [23]. Moreover, considering the dictionary as a fixed parameter of the algorithm allows one to analyze its influence on the performance of the model, and to adjust it so as to maximize the performance [20]. Simulation results illustrate the validity of the model. In future works, we will explore the use of the presented model to devise parameter selection strategies, providing guaranteed means of inferring the topology of a graph over time.

Definitions: A graph \mathcal{G} consists of a set \mathcal{N} of $(N + 1)$ nodes, and a set \mathcal{E} of edges such that $(m, n) \in \mathcal{E}$ if, and only if node n is linked to node m . At node level, we collect a real-valued signal $\mathbf{y}(i) \triangleq [y_1(i), \dots, y_{N+1}(i)]^\top$, where $y_n(i)$ is the sample of the signal $\mathbf{y}(i)$ at node n and time instant i . The adjacency matrix \mathbf{A} [24], is defined as an $(N + 1) \times (N + 1)$ matrix whose entries a_{nm} are zero if $(m, n) \notin \mathcal{E}$ and one otherwise. We also consider $a_{nn} = 0$, for all n .

Notations: Normal font letters denote scalars, while boldface lowercase and uppercase letters stand for column vectors and matrices, respectively. Uppercase calligraphic letters denote sets, of cardinality $|\cdot|$. Finally, $\mathbb{E}\{\cdot\}$ is the expectation operator.

2. PROBLEM FORMULATION

Let us consider measurements $\mathbf{y}(i) \in \mathbb{R}^{N+1}$ acquired over an $(N + 1)$ -node graph with adjacency matrix \mathbf{A} . We suppose the measurements $\mathbf{y}(i)$ to be acquired sequentially at all time instants $i \geq 0$. Signal $y_n(i)$ at each node $n = 1, \dots, N$ of the graph is non-linearly coupled to the signals at all nodes in its neighborhood, according to the topology described in \mathbf{A} . By considering an additive nonlinear model, the local measurements at the n -th node can be represented as:

$$y_n(i) = \sum_{m \in \mathcal{N} \setminus \{n\}} f_{nm}(\mathbf{y}_{L_m}(i)) + v_n(i), \quad (1)$$

where $\mathbf{y}_{L_m} = [y_m(i), \dots, y_m(i - L_m + 1)]^\top$, functions $f_{nm} : \mathbb{R}^{L_m} \rightarrow \mathbb{R}$ represent the interactions between the different nodes in the network, and $v_n(i)$ denotes innovation noise. By relating how each node $m \in \mathcal{N} \setminus \{n\}$ influences node n , functions f_{nm} encode the connectivity in \mathbf{A} directly as $a_{nm} = 0$ if and only if $f_{nm} \equiv 0$. For ease of notation, we assume that $n \equiv (N + 1)$, i.e., we identify n with the $(N + 1)$ -th node of the graph, which allows us to denote $\mathcal{N} \setminus \{n\} = \{1, \dots, N\}$. Thus, considering the available graph signal measurements $\mathbf{y}(i)$ for $1 \leq \ell \leq i$, the non-parametric local topology estimation problem at node n can be written as [16, 17]:

$$\min_{f_{n1}, \dots, f_{nN} \in \mathcal{H}} \frac{1}{2i} \sum_{\ell=1}^i \mathbb{E} \left\{ \left[y_n(i) - \sum_{m=1}^N f_{nm}(\mathbf{y}_{L_m}(i)) \right]^2 \right\} + \Psi(\|f_{n1}\|_{\mathcal{H}}, \dots, \|f_{nN}\|_{\mathcal{H}}) \quad (2)$$

where \mathcal{H} is a (normed) function space and $\Psi : \mathbb{R}^N \rightarrow [0, \infty[$ is a regularization functional which attempts to promote sparsity in the underlying adjacency matrix by favoring solutions in which many functions f_{nm} are identically zero.

Although several approaches can be considered for nonlinear modeling (such as, e.g., polynomial structural equation models [25]), kernel methods are particularly appealing due to their

elegance and efficiency. We consider \mathcal{H} to be an RKHS associated with a positive reproducing kernel $\kappa(\cdot, \cdot)$ [13]. Let us now assume that $f_{nm} \in \mathcal{H}$, $m \in \mathcal{N} \setminus \{n\}$, and that function Ψ can be decomposed as $\Psi(x_1, \dots, x_N) = \sum_{m=1}^N \psi_m(x_m)$, where each $\psi_m : \mathbb{R} \rightarrow [0, \infty[$ is non-decreasing. Then, since (2) employs a convex loss function, the conditions of the Representer Theorem are satisfied [26], which means that the solution to (2) admits a finite-dimensional representation of the form:

$$f_{nm}^*(\cdot) = \sum_{p=1}^i \alpha_{nmp} \kappa_m(\cdot, \mathbf{y}_{L_m}(p)), \quad m = 1, \dots, N, \quad (3)$$

where $\alpha_{nmp} \in \mathbb{R}$ are the representation coefficients.

An immediate observation concerning (3) is that the number of coefficients α_{nmp} becomes prohibitively large as i increases. This makes solving (2) online unfeasible. A solution to this problem is the use of kernel dictionaries, which represent f_{nm} as a linear combination of a small number of appropriately selected kernel functions $\kappa_m(\cdot, y_m(\omega_j))$ [18]. Such dictionaries, which we denote by $\mathcal{D}_m = \{\kappa_m(\cdot, y_m(\omega_j)) : \omega_1, \dots, \omega_{|\mathcal{D}_m|}\}$, can be either constructed online by selecting previous datapoints that satisfy some sparsification criterion (in which case $\omega_j \in \{1, \dots, i-1\}$) [18], or can be set a priori (in this case, we denote the sample indexes of the dictionary elements as $\omega_j < 0$, with a slight abuse of notation) [20].

Considering dictionaries \mathcal{D}_m , an efficient strategy for the online estimation of f_{nm} is the KLMS algorithm. The KLMS aims to minimize the following cost function at every time instant i using the stochastic gradient descent method:

$$J_n(\boldsymbol{\alpha}_n) = \frac{1}{2} \mathbb{E} \left\{ \left[y_n(i) - \boldsymbol{\alpha}_n^\top \mathbf{k}(i) \right]^2 \middle| \{\mathcal{D}_m\}_{m \in \mathcal{N} \setminus \{n\}} \right\}. \quad (4)$$

where $\boldsymbol{\alpha}_n = [\tilde{\boldsymbol{\alpha}}_{n1}^\top, \dots, \tilde{\boldsymbol{\alpha}}_{nN}^\top]^\top$, $\tilde{\boldsymbol{\alpha}}_{nm} = \text{col}\{\alpha_{nmp}\}_{p=1}^{|\mathcal{D}_m|}$, $\mathbf{k}(\ell) = [\tilde{\mathbf{k}}_1^\top(\ell), \dots, \tilde{\mathbf{k}}_N^\top(\ell)]^\top$, $\tilde{\mathbf{k}}_m(\ell) = \text{col}\{\kappa_m(\mathbf{y}_{L_m}(\ell), \mathbf{y}_{L_m}(\omega_j))\}_{j=1}^{|\mathcal{D}_m|}$.

This leads to the following coefficient update rule:

$$\hat{\boldsymbol{\alpha}}_{n(i+1)} = \hat{\boldsymbol{\alpha}}_{n(i)} + \mu \mathbf{k}(i) \varepsilon(i), \quad (5)$$

where $\varepsilon(i) \triangleq y_n(i) - \mathbf{k}^\top(i) \hat{\boldsymbol{\alpha}}_{n(i)}$ represents the instantaneous error conditioned on the dictionaries $\{\mathcal{D}_m\}_{m=1}^N$, and $\mu > 0$ denotes a small step size. The estimated coefficients $\hat{\boldsymbol{\alpha}}_{n(i)}$ can be related to the adjacency matrix \mathbf{A} at instant i by using a threshold τ_n by setting $a_{nm}(i) = 1$ if $\|\hat{\boldsymbol{\alpha}}_{nm(i)}\|_2 > \tau_n$ and $a_{nm}(i) = 0$ otherwise [27].

3. ALGORITHM ANALYSIS

Let us denote the optimal coefficients which minimize cost function (4) by $\boldsymbol{\alpha}_n^* = \mathbf{R}_{kk}^{-1} \mathbf{r}_{ky}$, with $\mathbf{R}_{kk} \triangleq \mathbb{E}\{\mathbf{k}(i)\mathbf{k}^\top(i)\}$ and $\mathbf{r}_{ky} \triangleq \mathbb{E}\{y_n(i)\mathbf{k}(i)\}$, and the difference between the current available estimate and the optimal solution by:

$$\mathbf{d}_{(i)} \triangleq \hat{\boldsymbol{\alpha}}_{n(i)} - \boldsymbol{\alpha}_n^*. \quad (6)$$

The remainder of this analysis concerns the use of the Gaussian kernel, chosen due to its universal approximating capabilities [28]. This kernel is defined as $k_n^G(\mathbf{y}_a, \mathbf{y}_b) = \exp(-\|\mathbf{y}_a - \mathbf{y}_b\|^2 / 2\sigma^2)$, where σ^2 is the kernel bandwidth.

Assumption 1: Dictionaries $\{\mathcal{D}_m\}$ are set beforehand. Inputs $\mathbf{y}(i)$ are assumed independent, zero-mean Gaussian random vectors with auto-correlation matrix $\mathbf{R}_{yy} \triangleq \mathbb{E}\{\mathbf{y}(i)\mathbf{y}^\top(i)\}$.

Assumption 2: Quantity $\mathbf{k}(i)\mathbf{k}^\top(i)$ is statistically independent of the error vector $\mathbf{d}_{(i)}$. A justification for the feasibility of the latter assumption is presented in [29].

Assumption 3: The optimal estimation error $\varepsilon_0(i) \triangleq y_n(i) - \mathbf{k}^\top(i)\boldsymbol{\alpha}_n^*$ given by the finite order model is close to the one by the infinite length model, such that $\mathbb{E}\{\varepsilon_0(i)\} \approx 0$.

Error $\varepsilon(i)$ can be expressed in terms of the error vector $\mathbf{d}_{(i)}$:

$$\varepsilon(i) = y_n(i) - \mathbf{k}^\top(i)\mathbf{d}_{(i)} - \mathbf{k}^\top(i)\boldsymbol{\alpha}_n^*. \quad (7)$$

Replacing (7) into (5) leads to the following recursion:

$$\mathbf{d}_{(i+1)} = \mathbf{d}_{(i)} - \mu\mathbf{k}(i)\mathbf{k}^\top(i)\mathbf{d}_{(i)} + \mu\mathbf{k}(i)\varepsilon_0(i). \quad (8)$$

3.1. Mean error behavior

Before proceeding to the mean behavior analysis, consider the quadratic form ξ of a Gaussian vector ζ , given by $\xi = \zeta\mathbf{B}\zeta^\top + \mathbf{b}^\top\zeta$, with $\mathbb{E}\{\zeta\} = \mathbf{0}$, $\mathbf{R}_{\zeta\zeta} = \mathbb{E}\{\zeta\zeta^\top\}$. For $t \in \mathbb{R}$, the moment-generating function of the random variable ξ is [30, p. 101]:

$$\begin{aligned} \Psi_\xi(t) &\triangleq \mathbb{E}\{\exp(t\xi)\} = \det\{\mathbf{I} - 2t\mathbf{B}\mathbf{R}_{\zeta\zeta}\}^{-\frac{1}{2}} \\ &\times \exp\left(\frac{t^2}{2}\mathbf{b}^\top\mathbf{R}_{\zeta\zeta}(\mathbf{I} - 2t\mathbf{B}\mathbf{R}_{\zeta\zeta})^{-1}\mathbf{b}\right). \end{aligned} \quad (9)$$

Taking the expectation of relation (8), and employing Assumptions 2 and 3, leads to the mean error (ME) behavior:

$$\mathbb{E}\{\mathbf{d}_{(i+1)}\} = (\mathbf{I} - \mu\mathbf{R}_{kk})\mathbb{E}\{\mathbf{d}_{(i)}\}. \quad (10)$$

Block matrix \mathbf{R}_{kk} contains blocks $\mathbf{R}_{kk}^{(m_1, m_2)} \triangleq \mathbb{E}\{\mathbf{k}_{m_1}(i)\mathbf{k}_{m_2}^\top(i)\}$, $\forall m_1, m_2 \in \mathcal{N}_{\setminus n}$. Each entry (u, v) , $u = 1, \dots, |\mathcal{D}_{m_1}|$, $v = 1, \dots, |\mathcal{D}_{m_2}|$ of every block $\mathbf{R}_{kk}^{(m_1, m_2)}$ is:

$$\begin{aligned} [\mathbf{R}_{kk}^{(m_1, m_2)}]_{uv} &= \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{y}^{(uv)}\|^2\right) \\ &\times \mathbb{E}\left\{\exp\left(-\frac{1}{\sigma^2}\left(\frac{1}{2}\|\mathbf{G}\mathbf{y}^{(ii)}\|^2 - (\mathbf{y}^{(uv)})^\top\mathbf{G}\mathbf{y}^{(ii)}\right)\right)\right\}, \end{aligned} \quad (11)$$

where $\mathbf{y}^{(uv)} = [\mathbf{y}_{L_{m_1}}^\top(\omega_u), \mathbf{y}_{L_{m_2}}^\top(\omega_v)]^\top$, and $\mathbf{y}^{(ii)}$ and \mathbf{G} are given by:

$$\mathbf{y}^{(ii)} = \begin{cases} [\mathbf{y}_{L_{m_1}}^\top(i), \mathbf{y}_{L_{m_2}}^\top(i)]^\top, & m_1 \neq m_2 \\ \mathbf{y}_{L_{m_1}}^\top(i), & m_1 = m_2 \end{cases}, \quad (12)$$

$$\mathbf{G} = \begin{cases} \mathbf{I}, & m_1 \neq m_2 \\ [\mathbf{I}, \mathbf{I}]^\top, & m_1 = m_2 \end{cases}. \quad (13)$$

Making use of (9) with $\mathbf{B} = \frac{1}{2}\mathbf{G}^\top\mathbf{G}$, $\mathbf{b} = -\mathbf{G}^\top\mathbf{y}^{(uv)}$ and $t = -\frac{1}{\sigma^2}$, we obtain:

$$\begin{aligned} [\mathbf{R}_{kk}^{(m_1, m_2)}]_{uv} &= \exp\left(\frac{1}{2\sigma^4}(\mathbf{y}^{(uv)})^\top\mathbf{G}\mathbf{H}^{(m_1 m_2)}\mathbf{G}^\top\mathbf{y}^{(uv)}\right) \\ &\times \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{y}^{(uv)}\|^2\right) \det\left\{\mathbf{I} + \frac{1}{\sigma^2}\mathbf{G}^\top\mathbf{G}\mathbf{R}_{yy}^{(m_1 m_2)}\right\}^{-\frac{1}{2}}, \end{aligned} \quad (14)$$

where \mathbf{I} is of size $(L_{m_1} + L_{m_2}) \times (L_{m_1} + L_{m_2})$, $\mathbf{H}^{(m_1 m_2)} = \mathbf{R}_{yy}^{(m_1 m_2)}\left(\mathbf{I} + \frac{1}{\sigma^2}\mathbf{G}^\top\mathbf{G}\mathbf{R}_{yy}^{(m_1 m_2)}\right)^{-1}$, with $\mathbf{R}_{yy}^{(m_1 m_2)}$ a block matrix. Each of its blocks $(k, \ell) \in \{1, 2\}^2$ contains $[\mathbf{R}_{yy}]_{m_k m_\ell}$ on all entries of its main diagonal, and zeros elsewhere.

From (10), convergence of the coefficients in the mean is assured if the step size satisfies $0 < \mu < 2/\lambda_{\max}(\mathbf{R}_{kk})$, where $\lambda_{\max}(\cdot)$ represents the maximum eigenvalue of its matrix argument.

3.2. Mean square error behavior

Let us denote $\mathbf{D}_{(i)} \triangleq \mathbb{E}\{\mathbf{d}_{(i)}\mathbf{d}_{(i)}^\top\}$. Using (8) and the previous assumptions, we obtain the mean square error (MSE) behavior:

$$\mathbf{D}_{(i+1)} = \mathbf{D}_{(i)} - \mu(\mathbf{D}_{(i)}\mathbf{R}_{kk} + \mathbf{R}_{kk}\mathbf{D}_{(i)}) + \mu^2\mathbf{Q} + \mu^2\mathbf{R}_{kk}J_{n, \min}, \quad (15)$$

with $J_{n, \min}$ being the minimum value of the cost function (4), and:

$$J_{n, \min} = J_n(\boldsymbol{\alpha}_n^*) = \mathbb{E}\{y_n^2(i)\} - \mathbf{r}_{ky}^\top\mathbf{R}_{kk}^{-1}\mathbf{r}_{ky}, \quad (16)$$

$$\mathbf{Q} = \mathbb{E}\left\{\mathbf{k}(i)\mathbf{k}^\top(i)\mathbf{d}_{(i)}\mathbf{v}_{(i)}^\top\mathbf{k}(i)\mathbf{k}^\top(i)\right\}. \quad (17)$$

These second order moments relate to the MSE via [31]:

$$J_{n, \text{MSE}}(i) \triangleq \mathbb{E}\{\varepsilon^2(i)\} = J_{n, \min} + \text{Tr}\{\mathbf{R}_{kk}\mathbf{D}_{(i)}\}, \quad (18)$$

and to the MSD through:

$$\text{MSD}(i) \triangleq \mathbb{E}\{\|\mathbf{d}_{(i)}\|^2\} = \text{Tr}\{\mathbf{D}_{(i)}\}. \quad (19)$$

Let us note $k_D = \sum_{m \in \mathcal{N}_{\setminus n}} |\mathcal{D}_m|$, the total number of dictionary entries. We make use of Assumption 2, leading to the writing of the (u, v) -th entry of \mathbf{Q} as:

$$[\mathbf{Q}]_{uv} = \sum_{a=1}^{k_D} \sum_{b=1}^{k_D} \mathbb{E}\{[\mathbf{k}(i)]_u [\mathbf{k}(i)]_v [\mathbf{k}(i)]_a [\mathbf{k}(i)]_b\} [\mathbf{D}_{(i)}]_{ab}. \quad (20)$$

We introduce matrix $\mathbf{K}^{(u, v)}$, whose (a, b) -th entry is:

$$[\mathbf{K}^{(u, v)}]_{ab} = \mathbb{E}\{[\mathbf{k}(i)]_u [\mathbf{k}(i)]_v [\mathbf{k}(i)]_a [\mathbf{k}(i)]_b\}. \quad (21)$$

Now we can write relation (20) as:

$$[\mathbf{Q}(i)]_{uv} = \text{Tr}\{\mathbf{K}^{(u, v)}\mathbf{D}_{(i)}\}. \quad (22)$$

To compute the expectation in (21), we need to find which block of $\mathbf{k}(i)$ (and which element therein) is being indexed by values u, v, a and b . Let us define functions $\varsigma, \varrho: \mathbb{N}_* \rightarrow \mathbb{N}_*$, which relate an index u in block vector $\mathbf{k}(i)$ to its corresponding constituent block $\varsigma(u)$ and entry $\varrho(u): [\mathbf{k}(i)]_u = [\tilde{\mathbf{k}}_{L_{\varsigma(u)}}(i)]_{\varrho(u)}$ (i.e., $[\mathbf{k}(i)]_u$ is the $\varrho(u)$ -th entry of the $\varsigma(u)$ -th block). Then, by denoting $m_1 = \varsigma(u)$, $m_2 = \varsigma(v)$, $m_3 = \varsigma(a)$, $m_4 = \varsigma(b)$ and $\omega_p = \omega_{\varrho(u)}$, $\omega_q = \omega_{\varrho(v)}$, $\omega_r = \omega_{\varrho(a)}$, $\omega_s = \omega_{\varrho(b)}$, we can write the following:

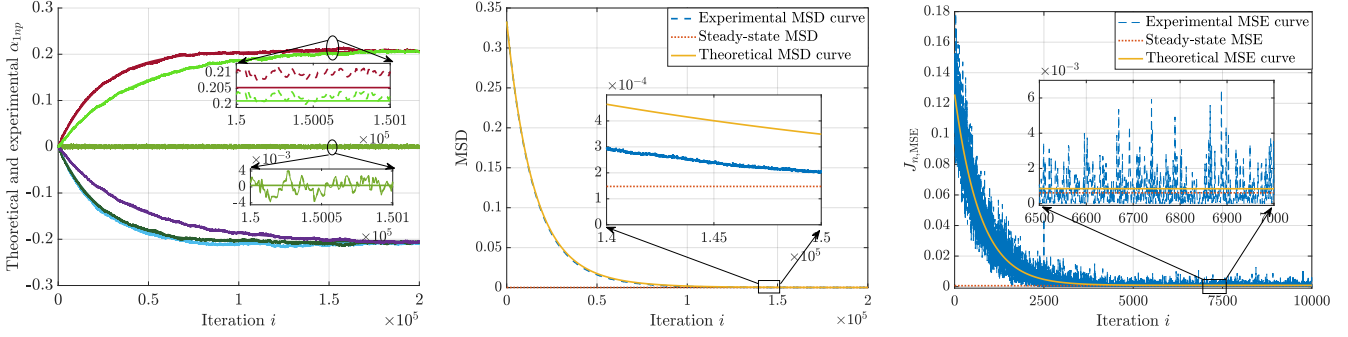
$$\begin{aligned} [\mathbf{K}^{(u, v)}]_{ab} &= \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{y}^d\|^2\right) \\ &\times \mathbb{E}\left\{\exp\left(-\frac{1}{\sigma^2}\left(\frac{1}{4}\|\tilde{\mathbf{G}}\mathbf{y}^i\|^2 - (\mathbf{y}^d)^\top\tilde{\mathbf{G}}\mathbf{y}^i\right)\right)\right\}, \end{aligned}$$

with $\mathbf{y}^d = [\mathbf{y}_{L_{m_1}}^\top(\omega_p), \mathbf{y}_{L_{m_2}}^\top(\omega_q), \mathbf{y}_{L_{m_3}}^\top(\omega_r), \mathbf{y}_{L_{m_4}}^\top(\omega_s)]^\top$ collecting the dictionary entries and $\mathbf{y}^i = \text{col}\{\mathbf{y}_{L_{m_\ell}}(i)\}_{m_\ell \in \cup_{\ell=1}^4 \{m_\ell\}}$ collecting the unique instantaneous measurements (i.e., without the repeated elements when there is $m_r = m_s$ for $r \neq s$), respectively. Matrix $\tilde{\mathbf{G}}$, which depends on m_1, \dots, m_4 , duplicates any repeated inputs on its image, such that

$$\tilde{\mathbf{G}}\mathbf{y}^i = [\mathbf{y}_{L_{m_1}}^\top(i), \mathbf{y}_{L_{m_2}}^\top(i), \mathbf{y}_{L_{m_3}}^\top(i), \mathbf{y}_{L_{m_4}}^\top(i)]^\top.$$

We now use relation (9), with $\mathbf{B} = \frac{1}{4}\tilde{\mathbf{G}}^\top\tilde{\mathbf{G}}$, $\mathbf{b} = -\tilde{\mathbf{G}}^\top\mathbf{y}^d$ and $t = -\frac{1}{\sigma^2}$. Thus, we obtain:

$$\begin{aligned} [\mathbf{K}^{(u, v)}]_{ab} &= \det\left\{\mathbf{I} + \frac{1}{2\sigma^2}\tilde{\mathbf{G}}^\top\tilde{\mathbf{G}}\mathbf{R}_{yy}^{(m_1 \rightarrow 4)}\right\}^{-\frac{1}{2}} \\ &\times \exp\left(\frac{1}{2\sigma^4}(\mathbf{y}^d)^\top\tilde{\mathbf{G}}\mathbf{H}^{(m_1 \rightarrow 4)}\tilde{\mathbf{G}}^\top\mathbf{y}^d - \frac{1}{2\sigma^2}\|\mathbf{y}^d\|^2\right), \end{aligned} \quad (23)$$



(a) Theoretical and experimental entries of α . (b) Experimental, steady-state, and theoretical MSD curves. Continuous lines are the theoretical curves, while dashed lines are experimental ones. (c) Experimental, steady-state, and theoretical MSE curves. For reasons of better visibility, a different initialization $\hat{\alpha}_n(0)$ was used.

Fig. 1: Analysis validation in the mean and mean square sense.

where $\mathbf{H}^{(m_1 \rightarrow 4)} = \mathbf{R}_{yy}^{(m_1 \rightarrow 4)} \left(\mathbf{I} + \frac{1}{2\sigma^2} \tilde{\mathbf{G}}^\top \tilde{\mathbf{G}} \mathbf{R}_{yy}^{(m_1 \rightarrow 4)} \right)^{-1}$, identity \mathbf{I} is of size $\sum_{\ell=1}^4 L_{m_\ell} \times \sum_{\ell=1}^4 L_{m_\ell}$, and $\mathbf{R}_{yy}^{(m_1 \rightarrow 4)}$ is a block matrix formed similarly to $\mathbf{R}_{yy}^{(m_1 m_2)}$: each of its blocks $(k, \ell) \in \{1, 2, 3, 4\}^2$ contains $[\mathbf{R}_{yy}]_{m_k m_\ell}$ on all entries of its main diagonal, and zeros elsewhere.

By stacking the columns of $\mathbf{D}_{(i)}$ on top of each other, i.e., $\bar{\mathbf{d}}_{(i)} = \text{vec} \{ \mathbf{D}_{(i)} \}$ and making use of the properties of the vectorization operator, recursion (15) can now be computed as:

$$\bar{\mathbf{d}}_{(i+1)} = \mathbf{F}_0 \bar{\mathbf{d}}_{(i)} + \mu^2 J_{n, \min} \bar{\mathbf{r}}_{kk}, \quad (24)$$

where $\bar{\mathbf{r}}_{kk} = \text{vec} \{ \mathbf{R}_{kk} \}$, and:

$$\mathbf{F}_0 = \mathbf{I}_2 - \mu (\mathbf{I} \otimes \mathbf{R}_{kk} + \mathbf{R}_{kk} \otimes \mathbf{I}) + \mu^2 \mathbf{F}_1. \quad (25)$$

We remark upon the fact that the identity matrix \mathbf{I}_2 is of size $k_D^2 \times k_D^2$, while \mathbf{I} is of size $k_D \times k_D$. Also, entries of the matrix \mathbf{F}_1 are $[\mathbf{F}_1]_{u+(v-1)k_D, a+(b-1)k_D} = [\mathbf{K}^{(u,v)}]_{ab}$.

Assuming a small enough step size μ , the algorithm is mean-square stable as $i \rightarrow \infty$, and converges towards:

$$\lim_{i \rightarrow \infty} \bar{\mathbf{d}}_{(i)} = \mu^2 J_{n, \min} (\mathbf{I} - \mathbf{F}_0)^{-1} \bar{\mathbf{r}}_{kk} = \bar{\mathbf{d}}_{(\infty)}. \quad (26)$$

4. EXPERIMENTAL VALIDATION

We consider a simulation scenario with i.i.d. Gaussian data $\mathbf{y}(i)$ generated using a correlation matrix \mathbf{R}_{yy} , depicted in (27). We note that this particular correlation matrix also corresponds to the data correlation of the linear model $\mathbf{y}(i) = \mathbf{A}\mathbf{y}(i) + \mathbf{v}(i)$, with $\mathbf{v}(i)$ zero-mean Gaussian noise with covariance $\sigma_v^2 \mathbf{I}$, where $\sigma_v = 0.05$. This model, although not inherently nonlinear, allows direct knowledge of the 5-node ground truth matrix \mathbf{A} given in (28). Moreover, it offers exact knowledge of the statistical properties of the input \mathbf{y} , necessary in the analysis. See [17] for the methods' behavior in nonlinear settings. For the algorithm we selected $L_m = 1$, for all m , a step size $\mu = 5 \cdot 10^{-2}$, kernel bandwidths of $\sigma = 1$, and each node stored a dictionary \mathcal{D}_m with 3 entries, chosen in a uniform grid on $[-1, 1]$. We compare the theoretical model for the ME given by (10) and for the MSE given by (15) with the empirical performance of the algorithm, averaged over 100 Monte Carlo runs. Fig. 1a shows both the theoretical and experimental values for a subset of non-zero coefficients α_{nmp} , for $n = 1$. Fig. 1b and 1c show both the theoretical

and experimental MSD and MSE curves, as well as their steady-state values, computed using relations (26), (19), (18) and (15). It can be seen that the theoretical model was able to predict the behavior of the algorithm very accurately, both in the mean and mean-square sense. Moreover, the resulting theoretical curves serve to show that an adequate step-size, as well as other parameters, can be selected in order to attain a desirable application-dependent performance. In particular, with a small dictionary of only 3 entries for each node, it was possible to obtain a very small MSE at steady-state, which means that the estimated \hat{f}_{nm} were able to predict the observations y_n accurately. A similar behavior can also be observed for the MSD. In the recovery of the network topology \mathbf{A} , we obtained an average error of around 10%, corresponding to only two links estimated incorrectly. This indicates that the considered algorithm can achieve a good trade-off between performance and computational complexity.

$$\mathbf{R}_{yy} = 10^{-4} \cdot \begin{pmatrix} 8.52 & 1.70 & -2.84 & -2.84 & 1.70 \\ 1.70 & 8.52 & 1.70 & -2.84 & -2.84 \\ -2.84 & 1.70 & 8.52 & 1.70 & -2.84 \\ -2.84 & -2.84 & 1.70 & 8.52 & 1.70 \\ 1.70 & -2.84 & -2.84 & 1.70 & 8.52 \end{pmatrix} \quad (27)$$

$$\mathbf{A} = \begin{pmatrix} 0 & 1 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 & 0 \end{pmatrix} \quad (28)$$

5. CONCLUSION

Designing online, adaptive network topology identification strategies which account for nonlinear interactions between network agents is an important topic in network analysis. In this paper, we presented a statistical convergence analysis of the kernel LMS algorithm applied to graph topology identification. The derived model characterizes the performance of the algorithm as a function of parameters such as the kernel bandwidth and the dictionary. In turn, this allows for precise tuning of such parameters in order to obtain a desired performance in transient or in steady state. Moreover, the theoretical model quantifies the typical statistical behavior of the algorithm both in the presence and in the absence of a link between a given node pair. Future work can harness this fact in order to design an automatic thresholding rule by comparing the behavior of the empirical algorithm coefficients to the one dictated by the model.

6. REFERENCES

- [1] R. H. Heiberger, “Predicting economic growth with stock networks,” *Physica A: Statistical Mechanics and its Applications*, vol. 489, pp. 102–111, 2018.
- [2] M. Kramer, E. D. Kolaczyk, and H. Kirsch, “Emergent network topology at seizure onset in humans,” *Epilepsy research*, vol. 79, pp. 173–86, 2008.
- [3] P. M. Djurić and C. Richard, *Cooperative and Graph Signal Processing: Principles and Applications*, Academic Press, 2018.
- [4] G. B. Giannakis, Y. Shen, and G. V. Karanikolas, “Topology Identification and Learning over Graphs: Accounting for Non-linearities and Dynamics,” *Proceedings of the IEEE*, vol. 106, no. 5, pp. 787–807, 2018.
- [5] J. Friedman, T. Hastie, and R. Tibshirani, “Sparse inverse covariance estimation with the graphical Lasso,” *Biostatistics*, vol. 9, no. 3, pp. 432–41, 2008.
- [6] S. Segarra, A. G. Marques, G. Mateos, and A. Ribeiro, “Network topology inference from spectral templates,” *IEEE Transactions on Signal and Information Processing over Networks*, vol. 3, no. 3, pp. 467–483, 2017.
- [7] S. Vlaski, H. P. Maretić, R. Nassif, P. Frossard, and A. H. Sayed, “Online graph learning from sequential data,” in *Proc. IEEE Data Science Workshop*, Lausanne, Switzerland, 2018, pp. 190–194.
- [8] M. Moscu, R. Nassif, F. Hua, and C. Richard, “Learning causal networks topology from streaming graph signals,” in *Proc. 27th European Signal Processing Conference (EUSIPCO)*, 2019, pp. 1–5.
- [9] R. Shafipour, A. Hashemi, G. Mateos, and H. Vikalo, “Online topology inference from streaming stationary graph signals,” *IEEE Data Science Workshop (DSW)*, pp. 140–144, 2019.
- [10] J. Harring, B. Weiss, and J.-C. Hsu, “A comparison of methods for estimating quadratic effects in nonlinear structural equation models,” *Psychological methods*, vol. 17, pp. 193–214, 2012.
- [11] W. Holmes Finch, “Modeling nonlinear structural equation models: A comparison of the two-stage generalized additive models and the finite mixture structural equation model,” *Structural Equation Modeling: A Multidisciplinary Journal*, vol. 22, no. 1, pp. 60–75, 2015.
- [12] N. Lim, F. D’Alché-Buc, C. Auliac, and G. Michailidis, “Operator-valued kernel-based vector autoregressive models for network inference,” *Machine learning*, vol. 99, no. 3, pp. 489–513, 2015.
- [13] V. I. Paulsen and M. Raghupathi, *An introduction to the theory of reproducing kernel Hilbert spaces*, vol. 152, Cambridge University Press, 2016.
- [14] Y. Shen, B. Baingana, and G. B. Giannakis, “Topology inference of directed graphs using nonlinear structural vector autoregressive models,” in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 6513–6517.
- [15] L. Zhang, G. Wang, and G. B. Giannakis, “Going beyond linear dependencies to unveil connectivity of meshed grids,” in *Proc. IEEE International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*, Curaçao, Dutch Antilles, 2017, pp. 1–5.
- [16] Y. Shen and G. B. Giannakis, “Online identification of directional graph topologies capturing dynamic and nonlinear dependencies,” in *Proc. IEEE Data Science Workshop (DSW)*, Lausanne, Switzerland, 2018, pp. 195–199.
- [17] M. Moscu, R. Borsoi, and C. Richard, “Online graph topology inference with kernels for brain connectivity estimation,” in *Proc. IEEE International Conference on Acoustics Speech and Signal Processing (ICASSP)*, Barcelona, Spain, 2020, pp. 1200–1204.
- [18] C. Richard, J.-C. M. Bermudez, and P. Honeine, “Online prediction of time series data with kernels,” *IEEE Transactions on Signal Processing*, vol. 57, no. 3, pp. 1058–1067, 2009.
- [19] M. Moscu, R. Borsoi, and C. Richard, “Online kernel-based graph topology identification with partial-derivative-imposed sparsity,” in *Proc. 28th European Conference on Signal Processing (EUSIPCO)*, Amsterdam, The Netherlands, 2020, pp. 2190–2194.
- [20] J. Chen, W. Gao, C. Richard, and J.-C. M. Bermudez, “Convergence analysis of kernel LMS algorithm with pre-tuned dictionary,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Florence, Italy, 2014, pp. 7243–7247.
- [21] W. Gao, J. Chen, C. Richard, J. Huang, and R. Flamary, “Kernel LMS algorithm with forward-backward splitting for dictionary learning,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2013.
- [22] A. Koppel, S. Paternain, C. Richard, and A. Ribeiro, “Decentralized online learning with kernels,” *Signal Processing, IEEE Transactions on*, vol. 66, no. 12, pp. 3240–3255, 2018.
- [23] S. Petkoski and V. K. Jirsa, “Transmission time delays organize the brain network synchronization,” *A Philosophical Transactions of the Royal Society*, vol. 377, no. 2153, 2019.
- [24] N. Biggs, *Algebraic Graph Theory*, Cambridge University Press, 1993.
- [25] K. G. Jöreskog, F. Yang, G. Marcoulides, and R. Schumacker, “Nonlinear structural equation models: The Kenny-Judd model with interaction effects,” *Advanced structural equation modeling: Issues and techniques*, pp. 57–88, 1996.
- [26] B. Schölkopf, R. Herbrich, and A. J. Smola, “A generalized representer theorem,” in *International conference on computational learning theory*. Springer, 2001, pp. 416–426.
- [27] Y. Shen, B. Baingana, and G. B. Giannakis, “Kernel-based structural equation models for topology identification of directed networks,” *IEEE Transactions on Signal Processing*, vol. 65, no. 10, pp. 2503–2516, 2017.
- [28] W. Liu, J. C. Principe, and S. Haykin, *Kernel adaptive filtering: a comprehensive introduction*, chapter 1, pp. 1–26, John Wiley & Sons, Ltd, 2010.
- [29] J. Minkoff, “Comment on the ”Unnecessary assumption of statistical independence between reference signal and filter weights in feedforward adaptive systems”,” *IEEE Transactions on Signal Processing*, vol. 49, no. 5, pp. 1109–, 2001.
- [30] J. Omura and Thomas Kailath, *Useful Probability Distributions*, vol. 2, Stanford Electronics, September 1965.
- [31] W. D. Parreira, J. C. M. Bermudez, C. Richard, and J.-Y. Tourneret, “Stochastic behavior analysis of the Gaussian kernel least-mean-square algorithm,” *IEEE Transactions on Signal Processing*, vol. 60, no. 5, pp. 2208–2222, 2012.