# Autoregressive Kernels for Time Series

Marco Cuturi
Graduate School of Informatics
Kyoto University
mcuturi@i.kyoto-u.ac.jp

Arnaud Doucet
Department of Computer Science & Department of Statistics
University of British Columbia
arnaud@cs.ubc.ca

January 5, 2011

### Abstract

We propose in this work a new family of kernels for variable-length time series. Our work builds upon the vector autoregressive (VAR) model for multivariate stochastic processes: given a multivariate time series $\mathbf{x}$, we consider the likelihood function $p_\theta(\mathbf{x})$ of different parameters $\theta$ in the VAR model as features to describe $\mathbf{x}$. To compare two time series $\mathbf{x}$ and $\mathbf{x}'$, we form the product of their features $p_\theta(\mathbf{x}) \cdot p_\theta(\mathbf{x}')$ which is integrated out w.r.t $\theta$ using a matrix normal-inverse Wishart prior. Among other properties, this kernel can be easily computed when the dimension $d$ of the time series is much larger than the lengths of the considered time series $\mathbf{x}$ and $\mathbf{x}'$. It can also be generalized to time series taking values in arbitrary state spaces, as long as the state space itself is endowed with a kernel $\kappa$. In that case, the kernel between $\mathbf{x}$ and $\mathbf{x}'$ is a a function of the Gram matrices produced by $\kappa$ on observations and subsequences of observations enumerated in $\mathbf{x}$ and $\mathbf{x}'$. We describe a computationally efficient implementation of this generalization that uses low-rank matrix factorization techniques. These kernels are compared to other known kernels using a set of benchmark classification tasks carried out with support vector machines.

## 1 Introduction

Kernel methods [Hofmann et al., 2008] have proved useful to handle and analyze structured data. A non-exhaustive list of such data types includes images [Chapelle et al., 1999, Grauman and Darrell, 2005, Cuturi and Fukumizu, 2007, Harchaoui and Bach, 2007], graphs [Kashima et al., 2003, Mahe et al., 2005, Vishwanathan et al., 2008, Shervashidze and Borgwardt, 2009], texts [Joachims, 2002, Moschitti and Zanzotto, 2007] and strings on finite alphabets [Leslie et al., 2002, Cortes et al., 2004, Vert et al., 2004, Cuturi and Vert, 2005, Sonnenburg et al., 2007], which have all drawn much attention in recent years. Time series, although ubiquitous in science and engineering, have been comparatively the subject of less research in the kernel literature.

Numerous similarity measures and distances for time series have been proposed in the past decades [Schreiber and Schmitz, 1997]. These similarities are not, however, always well suited to the kernel methods framework. First, most available similarity measures are not positive definite [Haasdonk and Bahlmann, 2004]. Likewise, most distances are not negative definite.

The positive definiteness of similarity measures (alternatively the negative definiteness of distances) is needed to use the convex optimization algorithms that underly most kernel machines. Positive definiteness is also the cornerstone of the reproducing kernel Hilbert space (RKHS) framework which supports these techniques [Berlinet and Thomas-Agnan, 2003]. Second, most similarities measures are only defined for 'standard' multivariate time series, that is time series of finite dimensional vectors. Yet, some of the main application fields of kernel methods include bioinformatics, natural language processing and computer vision, where the analysis of time series of structured objects (images, texts, graphs) remains a very promising field of study. Ideally, a useful kernel for time series should be both positive definite and able to handle time series of structured data. An oft-quoted example [Bahlmann et al., 2002, Shimodaira et al., 2002] of a non-positive definite similarity for time series is the Dynamic Time Warping (DTW) score [Sakoei and Chiba, 1978], arguably the most popular similarity score for variable-length multivariate time series [Rabiner and Juang, 1993, §4.7]. Hence the DTW can only be used in a kernel machine if it is altered through ad hoc modifications such as diagonal regularizations [Zhou et al., 2010]. Some extensions of the DTW score have addressed this issue: Hayashi et al. [2005] propose to embed time series in Euclidean spaces such that the distance of such representations approximates the distance induced by the DTW score. Cuturi et al. [2007] consider the soft-max of the alignment scores of all possible alignments to compute a positive definite kernel for two time series. This kernel can also be used on two time series $\mathbf{x} = (x_1, \cdots, x_n)$ and $\mathbf{x}' = (x'_1, \cdots, x'_{n'})$ of structured objects since the kernel between $\mathbf{x}$ and $\mathbf{x}'$ can be expressed as a function of the Gram matrix $\mathcal{K} = \left[ \kappa(x_i, x'_j) \right]_{i \leq n, j \leq n'}$ where $\kappa$ is a given kernel on the structured objects of interest.

A few alternative kernels have been proposed for multivariate time series. Kumara et al. [2008] consider a non-parametric approach to interpolate time series using splines, and define directly kernels on these interpolated representations. In a paragraph of their broad work on probability product kernels, Jebara et al. [2004, §4.5] briefly mention the idea of using the Bhattacharyya distance on suitable representations as normal densities of two time series using state-space models. Vishwanathan et al. [2007] as well as Borgwardt et al. [2006] use the family of Binet-Cauchy kernels [Vishwanathan and Smola, 2004], originally defined by the coefficients of the characteristic polynomial of kernel matrices such as the matrix $\mathcal{K}$ described above (when $n = n'$). Unlike other techniques listed above, these two proposals rely on a probabilistic modeling of the time series to define a kernel. Namely, in both the probability product and Binet-Cauchy approaches the kernel value is the result of a two step computation: each time series is first mapped onto a set of parameters that summarizes their dynamic behavior; the kernel is then defined as a kernel between these two sets of parameters.

The kernels we propose in this paper also rely on a probabilistic modeling of time series to define kernels but do away with the two step approach detailed above. This distinction is discussed later in the paper in Remark 2. Our contribution builds upon the the *covariance kernel* framework proposed by Seeger [2002], and whose approach can be traced back to the work of Haussler [1999] and Jaakkola et al. [1999], who advocate the use of probabilistic models to extract features from structured objects. Given a measurable space $\mathcal{X}$ and a model, that is a parameterized family of probability distributions on $\mathcal{X}$ of the form $\{p_\theta, \theta \in \Theta\}$, a kernel for two objects $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ can be defined as

$$k(\mathbf{x}, \mathbf{x}') = \int_{\theta \in \Theta} p_\theta(\mathbf{x}) \, p_\theta(\mathbf{x}') \, \omega(d\theta),$$

where $\omega$ is a prior on the parameter space. Our work can be broadly characterized as the implementation of this idea for time series data, using the VAR model to define the space

2

of densities $p_\theta$ and a specific prior for $\omega$, the matrix-normal inverse-Wishart prior. This conjugate prior allows us to obtain a closed-form expression for the kernel which admits useful properties.

The rest of this paper is organized as follows. Section 2 starts with a brief review of the Bayesian approach to dynamic linear modeling for multivariate stochastic processes, which provides the main tool to define autoregressive kernels on multivariate time series. We follow by detailing a few of the appealing properties of autoregressive kernels for multivariate time series, namely their infinite divisibility and their ability to handle high-dimensional time series of short length. We show in Section 3 that autoregressive kernels can not only be used on multivariate time series but also on time series taking values in any set endowed with a kernel. The kernel is then computed as a function of the Gram matrices of subsets of shorter time series found in $\mathbf{x}$ and $\mathbf{x}'$. This computation requires itself the computation of large Gram matrices in addition to a large number of operations that grows cubicly with the lengths of $\mathbf{x}$ and $\mathbf{x}'$. We propose in Section 3.2 to circumvent this computational burden by using low-rank matrix factorization of these Gram matrices. We present in Section 4 different experimental results using toy and real-life datasets.

## 2 Autoregressive Kernels

We introduce in this section the crux of our contribution, that is a family of kernels that can handle variable-length multivariate time series. Sections 2.1, 2.2 and 2.3 detail the construction of such kernels, while Sections 2.4 and 2.5 highlight some of their properties.

### 2.1 Autoregressive Kernels as an Instance of Covariance Kernels

A vector autoregressive model of order $p$ henceforth abbreviated as VAR($p$) is a family of densities for $\mathbb{R}^d$-valued stochastic processes parameterized by $p$ matrices $A_i$ of size $d \times d$ and a positive definite matrix $V$. Given a parameter $\theta = (A_1, \cdots, A_p, V)$ the conditional probability density that an observed time series $\mathbf{x} = (x_1, x_2, \ldots, x_n)$ has been drawn from model $\theta$ given the $p$ first observations $(x_1, \ldots, x_p)$, assuming $p < n$, is equal to

$$p_\theta(\mathbf{x}|x_1, \cdots, x_p) = \frac{1}{(2\pi|V|)^{\frac{d(n-p)}{2}}} \prod_{i=p+1}^{n} \exp\left(-\frac{1}{2}\left\| x_i - \sum_{i=1}^{p} A_i x_{t-i} \right\|_V^2\right),$$

where for a vector $x$ and a positive definite matrix $V$ the Mahalanobis norm $\|x\|_V^2$ is equal to $x^T V^{-1} x$. We write $|\mathbf{x}|$ for the length of the time series $\mathbf{x}$, $n$ in the example above. We abbreviate the conditional probability $p_\theta(\mathbf{x}|x_1, \cdots, x_p)$ as $p_\theta(\mathbf{x})$ and take for granted that the $p$ first observations of the time series are not taken into account to compute the probability of $\mathbf{x}$. We consider the set of parameters

$$\Theta = \underbrace{\mathbb{R}^{d \times d} \times \cdots \times \mathbb{R}^{d \times d}}_{p} \times \mathbf{S}_d^{++},$$

where $\mathbf{S}_d^{++}$ is the cone of positive definite matrices of size $d \times d$, to define a kernel $k$ that computes a weighted sum of the product of features of two time series $\mathbf{x}$ and $\mathbf{x}'$ as $\theta$ varies in $\Theta$,

$$k(\mathbf{x}, \mathbf{x}') = \int_{\theta \in \Theta} p_\theta(\mathbf{x})^{c(|\mathbf{x}|)} p_\theta(\mathbf{x}')^{c(|\mathbf{x}'|)} \omega(d\theta), \tag{1}$$

where the exponential weight factor $c(|\mathbf{x}|)$ is used to normalize the probabilities $p_\theta(\mathbf{x})$ by the length of the considered time series $\mathbf{x}$.

**Remark 1.** *The feature map $\mathbf{x} \to \{p_\theta(\mathbf{x})\}_{\theta \in \Theta}$ produces strictly positive features. Features will be naturally closer to 0 for longer time series. We follow in this work the common practice in the kernel methods literature, as well as the signal processing literature, to normalize to some extent these features so that their magnitude is independent of the size of the input object, namely the length of the input time series. We propose to do so by normalizing the probabilities $p_\theta(\mathbf{x})$ by the lengths of the considered series. The weight $c(|\mathbf{x}|)$ introduced in Equation 1 is defined to this effect in section 2.3.*

**Remark 2.** *The formulation of Equation (1) is somehow orthogonal to kernels that map structured objects, in this case time series, to a single density in a given model and compare directly these densities using a kernel between densities, such as probability product kernels [Jebara et al., 2004], Binet-Cauchy kernels [Vishwanathan et al., 2007], structural kernels [Hein and Bousquet, 2005] or information diffusion kernels [Lebanon, 2006]. Indeed, such kernels rely first on the estimation of a parameter $\hat\theta_\mathbf{x}$ in a parameterized model class to summarize the properties of an object $\mathbf{x}$, and then compare two objects $\mathbf{x}$ and $\mathbf{x}'$ by using a proxy kernel on $\hat\theta_\mathbf{x}$ and $\hat\theta_{\mathbf{x}'}$. These approaches do require that the model class, the VAR model in the context of this paper, properly models the considered objects, namely that $\mathbf{x}$ is well summarized by $\hat\theta_\mathbf{x}$. The kernels we propose do not suffer from this key restriction: the VAR model considered in this work is never used to infer likely parameters for a time series $\mathbf{x}$ but is used instead to generate an infinite family of features.*

The kernel $k$ is mainly defined by the prior $\omega$ on the parameter space. We present in the next section a possible choice for this prior, the matrix-normal inverse-Wishart prior, which has been extensively studied in the framework of Bayesian linear regression applied to dynamic linear models.

## 2.2 The Matrix-Normal Inverse-Wishart Prior in the Bayesian Linear Regression Framework

Consider the regression model between a vector of explanatory variables $x$ of $\mathbb{R}^m$, and an output variable $y \in \mathbb{R}^d$, $y = Ax + \varepsilon$, where $A$ is a $d \times m$ coefficient matrix and $\varepsilon$ is a centered Gaussian noise with $d \times d$ covariance matrix $V$. Accordingly, $y$ follows conditionally to $x$ a normal density

$$y|(x, A, V) \sim \mathcal{N}(Ax, V).$$

Given $n$ pairs of explanatory and response vectors $(x_i, y_i)_{1 \le i \le n}$ weighted by $n$ nonnegative coefficients $t = (t_1, \ldots, t_n) \ge 0$ such that $\sum_{i=1}^n t_i = 1$, the weighted likelihood of this sample of $n$ observations is defined by the expression

$$
\begin{aligned}
\rho(Y|X, A, V, \Delta) &= \prod_{i=1}^n p(y_j|x_j, A, V)^{t_j}, \\
&= \frac{1}{|2\pi V|^{1/2}} \exp\left(-\frac{1}{2}\operatorname{tr}\Delta(Y-AX)^T V^{-1}(Y-AX)\right), \\
&= \frac{1}{|2\pi V|^{1/2}} \exp\left(-\frac{1}{2}\operatorname{tr} V^{-1}(Y-AX)\Delta(Y-AX)^T\right).
\end{aligned}
$$

4

where the matrices $\Delta, Y$ and $X$ stand respectively for $\mathbf{diag}(t_1, \ldots, t_n) \in \mathbb{R}^{n \times n}$, $[y_1, \cdots, y_n] \in \mathbb{R}^{d \times n}$ and $[x_1, \ldots, x_n] \in \mathbb{R}^{m \times n}$.

The matrix-normal inverse Wishart joint distribution West and Harrison [1997, §16.4.3] is a natural choice to model the randomness for $(A, V)$. The prior assumes that the $d \times m$ matrix $A$ is distributed following a centered matrix-normal density with left variance matrix parameter $V$ and right variance matrix parameter $\Omega$,

$$p(A) = \mathcal{MN}(A; 0, V, K) = \frac{1}{|2\pi V|^{m/2}|\Omega|^{d/2}} \exp\left(-\frac{1}{2} \operatorname{tr} A^T V^{-1} A \Omega^{-1}\right)$$

where $\Omega$ is a $m \times m$ positive definite matrix. Using the following notations,

$$S_{xx} = X\Delta X^T + \Omega^{-1}, \quad S_{yx} = Y\Delta X^T, S_{yy} = Y\Delta Y^T, \quad S_{y|x} = S_{yy} - S_{yx}S_{xx}^{-1}S_{yx}^T.$$

we can integrate out $A$ in $\rho(Y|X, A, V, \Delta)$ to obtain

$$\rho(Y|X, V) = \frac{1}{|\Omega|^{d/2}|S_{xx}|^{d/2}|2\pi V|^{1/2}} \exp\left(-\frac{1}{2} \operatorname{tr}(V^{-1}S_{y|x})\right). \tag{2}$$

The matrix-normal inverse Wishart prior for $(A, V)$ also assumes that $V$ is distributed with inverse-Wishart density $\mathcal{W}_\lambda^{-1}(\Sigma)$ of inverse scale matrix $\Sigma$ and degrees of freedom $\lambda > 0$. The posterior obtained by multiplying this prior by Equation (2) is itself proportional to an inverse-Wishart density with parameters $\mathcal{W}^{-1}(\Sigma + S_{y|x}, 1 + \lambda)$ which can be integrated to obtain the marginal weighted likelihood,

$$\rho(Y|X) = \prod_{i=1}^{n} \frac{\Gamma(\frac{\lambda+2-i}{2})}{\Gamma(\frac{\lambda+1-i}{2})} \frac{1}{|\Omega|^{d/2}|S_{xx}|^{d/2}} \frac{|\Sigma|^{\lambda/2}}{|S_{y|x} + \Sigma|^{\frac{1+\lambda}{2}}}$$

Using for $\Sigma$ the prior $I_d$, for $\Omega$ the matrix $I_m$, and discarding all constants independent of $X$ and $Y$ yields the expression

$$\rho(Y|X) \propto \frac{1}{|X\Delta X^T + I_m|^{d/2}|Y(\Delta - \Delta X^T (X\Delta X^T + I_m)^{-1} X\Delta)Y^T + I_d|^{\frac{1+\lambda}{2}}},$$

Note that the matrix $H_\Delta \overset{\text{def}}{=} \Delta X^T (X\Delta X^T + I_m)^{-1} X\Delta$ in the denominator is known as the hat-matrix of the orthogonal least-squares regression of $Y$ versus $X$. The right term in the denominator can be interpreted as the determinant of the weighted cross-covariance matrix of $Y$ with the residues $(\Delta - H_\Delta)Y$ regularized by the identity matrix.

## 2.3   Bayesian Averaging over VAR Models

Given a VAR($p$) model, we represent a time series $\mathbf{x} = (x_1, \cdots, x_n)$ as a sample $X$ of $n-p$ pairs of explanatory variables in $\mathbb{R}^{pd}$ and response variables in $\mathbb{R}^d$, namely $\{([x_i, \cdots, x_{p+i-1}], x_{p+i}), i = 1, \cdots, n - p\}$. Following a standard practice in the study of VAR models [Lütkepohl, 2005,

§3], this set is better summarized by matrices

$$
X = \begin{bmatrix} \begin{bmatrix} \vdots \\ x_1 \\ \vdots \end{bmatrix} & \cdots & \begin{bmatrix} \vdots \\ x_{n-p+1} \\ \vdots \end{bmatrix} \\ \vdots & \cdots & \vdots \\ \begin{bmatrix} \vdots \\ x_p \\ \vdots \end{bmatrix} & \cdots & \begin{bmatrix} \vdots \\ x_{n-1} \\ \vdots \end{bmatrix} \end{bmatrix} \in \mathbb{R}^{pd \times n-p}, \quad \text{and } Y = \begin{bmatrix} \vdots & \cdots & \vdots \\ x_{p+1} & \cdots & x_n \\ \vdots & \cdots & \vdots \end{bmatrix} \in \mathbb{R}^{d \times n-p}.
$$

Analogously, we use the corresponding notations $X', Y'$ for a second time series $\mathbf{x}'$ of length $n'$. Using the notation

$$
N \overset{\text{def}}{=} n + n' - 2p,
$$

these samples are aggregated in the $\mathbb{R}^{N \times N}$, $\mathbb{R}^{pd \times N}$ and $\mathbb{R}^{d \times N}$ matrices

$$
\Delta = \mathbf{diag}\left( \frac{1}{2} \left[ \underbrace{\tfrac{1}{n-p}, \cdots, \tfrac{1}{n-p}}_{n-p \text{ times}}, \underbrace{\tfrac{1}{n'-p}, \cdots, \tfrac{1}{n'-p}}_{n'-p \text{ times}} \right] \right), \quad \mathbf{X} = \begin{bmatrix} X\ X' \end{bmatrix}, \quad \mathbf{Y} = \begin{bmatrix} Y\ Y' \end{bmatrix}. \tag{3}
$$

Note that by setting $A = [A_1, \cdots, A_p]$ and $\theta = (A, V)$, the integrand that appears in Equation (1) can be cast as the following probability,

$$
p_\theta(\mathbf{x})^{\frac{1}{2(n-p)}} \, p_\theta(\mathbf{x}')^{\frac{1}{2(n'-p)}} = \rho(\mathbf{Y}|\mathbf{X}, A, V, \Delta).
$$

Integrating out $\theta$ using the matrix-normal inverse Wishart prior $\mathcal{MW}_\lambda^{-1}(I_d, I_{pd})$ for $(A, V)$ yields the following definition:

**Definition 1.** *Given two time series $\boldsymbol{x}, \boldsymbol{x}'$ and using the notations introduced in Equation (3), the autoregressive kernel $k$ of order $p$ and degrees of freedom $\lambda$ is defined as*

$$
k(\mathbf{x},\mathbf{x}') = \frac{1}{|\mathbf{X}\Delta\mathbf{X}^T + I_{pd}|^{\frac{d}{2}} \, |\mathbf{Y}(\Delta - \Delta\mathbf{X}^T \left( \mathbf{X}\Delta\mathbf{X}^T + I_{pd} \right)^{-1} \mathbf{X}\Delta)\mathbf{Y}^T + I_d|^{\frac{1+\lambda}{2}}}. \tag{4}
$$

## 2.4 Variance and Gram Based Formulations

We show in this section that $k$ can be reformulated in terms of Gram matrices of subsequences of $\mathbf{x}$ and $\mathbf{x}'$ rather than variance-covariance matrices. For two square matrices $C \in \mathbb{R}^{q \times q}$ and $D \in \mathbb{R}^{r \times r}$ we write $C \smallfrown D$ when the spectrums of $C$ and $D$ coincide, taking into account multiplicity, except for the value 0. Recall first the following trivial lemma.

**Lemma 1.** *For two matrices $A$ and $B$ in $\mathbb{R}^{q \times r}$ and $\mathbb{R}^{r \times q}$ respectively, $AB \smallfrown BA$, and as a consequence $|AB + I_q| = |BA + I_r|$*

Based on this lemma, it is possible to establish the following result.

**Proposition 1.** *Let $\alpha \overset{\text{def}}{=} \frac{1+\lambda}{d}$ then the autogressive kernel $k$ of order $p$ and degrees of freedom $\lambda$ given in Equation (4) is equal to*

$$
k(\boldsymbol{x}, \boldsymbol{x}') = \left( |\mathbf{X}^T\mathbf{X}\Delta + I_N|^{1-\alpha} \, |\mathbf{X}^T\mathbf{X}\Delta + \mathbf{Y}^T\mathbf{Y}\Delta + I_N|^\alpha \right)^{-\frac{d}{2}}, \tag{5}
$$

*Proof.* We use Lemma 1 to rewrite the first term of the denominator of Equation (4) using the Gram matrix $\mathbf{X}^T\mathbf{X}$,

$$|\mathbf{X}\Delta\mathbf{X}^T + I_{pd}| = |\mathbf{X}^T\mathbf{X}\Delta + I_N|.$$

Taking a closer look at the denominator, the matrix inversion lemma[1] yields the equality

$$(\Delta - \Delta\mathbf{X}^T\left(\mathbf{X}\Delta\mathbf{X}^T + I_{pd}\right)^{-1}\mathbf{X}\Delta) = (\mathbf{X}^T\mathbf{X} + \Delta^{-1})^{-1}.$$

Using again Lemma 1 the denominator of Equation (4) can be reformulated as

$$\left|\mathbf{Y}(\Delta - \Delta\mathbf{X}^T\left(\mathbf{X}\Delta\mathbf{X}^T + I_{pd}\right)^{-1}\mathbf{X}\Delta)\mathbf{Y}^T + I_d\right| = \frac{|\mathbf{X}^T\mathbf{X} + \Delta^{-1} + \mathbf{Y}^T\mathbf{Y}|}{|\mathbf{X}^T\mathbf{X} + \Delta^{-1}|}.$$

Factoring in these two results, we obtain Equation (5). ∎

We call Equation (4) the **Variance** formulation and Equation (5) the **Gram** formulation of the autoregressive kernel $k$ as it only depends on the Gram matrices of $\mathbf{X}$ and $\mathbf{Y}$. Although both the Variance and Gram formulations of $k$ are equal, their computational cost is different as detailed in the remark below.

**Remark 3.** *In a **high $N$-low** $d$ **setting**, the computation of $k$ requires $O(N(pd)^2)$ operations to compute the denominator's matrices and $O(p^3d^3 + d^3)$ to compute their inverse, which yields an overall cost of the order of $O(Np^2d^2 + (p^3 + 1)d^3)$. This may seem reasonable for applications where the cumulated time series lengths' $N$ is much larger than the dimension $d$ of these time series, such as speech signals or EEG data. In a **low $N$-high** $d$ **setting**, which frequently appears in bioinformatics or video processing applications, autoregressive kernels can be computed using Equation (5) in $O((p + 1)dN^2 + N^3)$ operations.*

## 2.5   Infinite Divisibility of Autoregressive Kernels

We recall that a positive definite kernel function $k$ is infinitely divisible if for all $n \in \mathbb{N}$, $k^{1/n}$ is also positive definite [Berg et al., 1984, §3.2.6]. We prove in this section that under certain conditions on $\lambda$, the degrees-of-freedom parameter of the inverse Wishart law, the autoregressive kernel is infinitely divisible. This result builds upon [Cuturi et al., 2005, Proposition 3].

Proving the infinite divisibility of $k$ is useful for the following two reasons: First, following a well-known result Berg et al. [1984, §3.2.7], the infinite divisibility of $k$ implies the negative definiteness of $-\log k$. Using Berg et al. [1984, §3.3.2] for instance, there exists a mapping $\Phi$ of $\mathcal{X}^{\mathbb{N}}$ onto a Hilbert space such that

$$\|\Phi(\mathbf{x}) - \Phi(\mathbf{x}')\|^2 = \frac{\log k(\mathbf{x},\mathbf{x}) + \log k(\mathbf{x}',\mathbf{x}')}{2} - \log k(\mathbf{x},\mathbf{x}').$$

and hence $k$ defines a Hilbertian metric for time series which can be used with distance-based tools such as nearest-neighbors.

Second, on a more practical note, the exponent $d/2$ in Equation (5) is numerically problematic when $d$ is large. In such a situation, the kernel matrices produced by $k$ would be diagonally dominant. This is analogous to selecting a bandwidth parameter $\sigma$ which is too small when using the Gaussian kernel on high-dimensional data. By proving the infinite divisibility of $k$, the exponent $d$ can be removed and substituted by any arbitrary exponent.

---

[1] $(A + UCV)^{-1} = A^{-1} - A^{-1}U\left(C^{-1} + VA^{-1}U\right)^{-1}VA^{-1}$

To establish the infinite divisibility result, we need a few additional notation. Let $\mathcal{M}(\mathbb{R}^d)$ be the set of positive measures on $\mathbb{R}^d$ with finite second-moment $\mu[xx^T] \stackrel{\text{def}}{=} \mathbb{E}_\mu[xx^T] \in \mathbb{R}^{d \times d}$. This set is a semigroup [Berg et al., 1984] when endowed with the usual addition of measures.

**Lemma 2.** *For two measures $\mu$ and $\mu'$ of $\mathcal{M}(\mathbb{R}^d)$, the kernel*

$$\tau : (\mu, \mu') \mapsto \frac{1}{\sqrt{|(\mu + \mu')[xx^T] + I_d|}},$$

*is an infinitely divisible positive definite kernel.*

*Proof.* The following identity is valid for any $d \times d$ positive-definite matrix $\Sigma$

$$\frac{1}{\sqrt{|\Sigma + I_d|}} = \frac{1}{\sqrt{(2\pi)^d}} \int_{\mathbb{R}^d} e^{-\frac{1}{2}y^T(\Sigma + I_d)\,y} dy$$

Given a measure $\mu$ with finite second-moment on $\mathbb{R}^d$, we thus have

$$\frac{1}{\sqrt{|\mu[xx^T] + I_d|}} = \int_{\mathbb{R}^d} e^{-\frac{1}{2}\langle \mu[xx^T],\, yy^T \rangle} \, p_{\mathcal{N}(0,I_d)}(dy) \tag{6}$$

where $\langle\,,\,\rangle$ stands for the Frobenius dot-product between matrices and $p_{\mathcal{N}(0,I_d)}$ is the standard multivariate normal density. In the formalism of [Berg et al., 1984] the integral of Equation (6) is an integral of bounded semicharacters on the semigroup $(\mathcal{M}(\mathbb{R}^d), +)$ equipped with autoinvolution. Each semicharacter $\rho_y$ is indexed by a vector $y \in \mathbb{R}^d$ as

$$\rho_y : \mu \mapsto e^{-\frac{1}{2}\langle \mu[xx^T],\, yy^T \rangle}.$$

To verify that $\rho_y$ is a semicharacter notice that $\rho_y(0) = 1$, where $0$ is the zero measure, and $\rho_y(\mu + \mu') = \rho_y(\mu)\,\rho_y(\mu')$ for two measures of $\mathcal{M}(\mathbb{R}^d)$. Now, using the fact that the multivariate normal density is a stable distribution, one has that for any $t \in \mathbb{N}$,

$$\frac{1}{\sqrt{|\mu[xx^T] + I_d|}} = \int_{\mathbb{R}^d} \left( e^{-\frac{1}{2t}\langle \mu[xx^T],\, yy^T \rangle} \right)^t p_{\mathcal{N}(0,I_d/t)}^{\otimes t}(dy),$$

$$= \left( \int_{\mathbb{R}^d} e^{-\frac{1}{2t}\langle \mu[xx^T],\, yy^T \rangle} \, p_{\mathcal{N}(0,I_d/t)}(dy) \right)^t,$$

where $p^{\otimes t}$ is the $t$-th convolution of density $p$, which proves the result. ∎

**Theorem 2.** *For $0 \leq \alpha \leq 1$, equivalently for $0 < \lambda \leq d - 1$, $\varphi \stackrel{\text{def}}{=} -\frac{2}{d} \log k$ is a negative definite kernel*

*Proof.* $k$ is infinitely divisible as the product of two infinitely divisible kernels, $\tau^{d(1-\alpha)}$ and $\tau^{d\alpha}$ computed on two different representations of $\mathbf{x}$ and $\mathbf{x}'$: first as empirical measures on $\mathbb{R}^{pd}$ with locations enumerated in the columns of $X$ and $X'$ respectively, and as empirical measures on $\mathbb{R}^{pd+d}$ with locations enumerated in the columns of the stacked matrices $[X; Y]$ and $[X'; Y']$. The set of weights for both representations are the uniform weights $\frac{1}{2(n-p)}$ and $\frac{1}{2(n'-p)}$. The negative definiteness of $\varphi$ follows from, and is equivalent to, the infinite divisibility of $k$ Berg et al. [1984, §3.2.7]. ∎

**Remark 4.** *The infinite divisibility of the joint distribution matrix normal-inverse Wishart distribution would be a sufficient condition to obtain directly the infinite divisibility of $k$ using for instance Berg et al. [1984, §3.3.7]. Unfortunately we have neither been able to prove this property nor found it in the literature. The inverse Wishart distribution alone is known to be not infinitely divisible in the general case [Lévy, 1948]. We do not know either whether $k$ can be proved to be infinitely divisible when $\lambda > d - 1$. The condition $0 < \lambda \le d - 1$, and hence $0 \le \alpha \le 1$ also plays an important role in Proposition 4.*

In the next sections, we will usually refer to the (negative definite) autoregressive kernel

$$\varphi(\mathbf{x}, \mathbf{x}') = C_{n,n'} + (1 - \alpha) \log |\mathbf{X}^T \mathbf{X} + \Delta^{-1}| + \alpha \log |\mathbf{X}^T \mathbf{X} + \mathbf{Y}^T \mathbf{Y} + \Delta^{-1}| \qquad (7)$$

where the constant $C_{n,n'} = (n-p) \log(2(n-p)) + (n'-p) \log(2(n'-p))$, rather than considering $k$ itself.

# 3 Extension to Time Series Valued in a Set Endowed with a Kernel

We show in Section that autoregressive kernels can be extended quite simply to time series valued in arbitrary spaces by considering Gram matrices. This extension is interesting but can be very computationally expensive. We propose a way to mitigate this computational cost by using low-rank matrix factorization techniques in Section 3.2

## 3.1 Autoregressive Kernels Defined Through Arbitrary Gram Matrices

Using again notation introduced in Equation (3), we write $K_{\mathbf{X}} = \mathbf{X}^T \mathbf{X}$ for the $N \times N$ Gram matrix of all explanatory variables contained in the joint sample $\mathbf{X}$ and $K_{\mathbf{Y}} = \mathbf{Y}^T \mathbf{Y}$ for the Gram matrix of all outputs of the local regression formulations. As stated in Equation (7) above, $\varphi$ and $k$ by extension can be defined as a function of the Gram matrices $K_{\mathbf{X}}$ and $K_{\mathbf{Y}}$. To alleviate notations, we introduce two functions $g$ and $f$ defined respectively on the cone of positive semidefinite matrices $\mathbf{S}_N^+$ and on $(\mathbf{S}_N^+)^2$:

$$g : Q \mapsto \log |Q + \Delta^{-1}|, \quad f : (Q, R) \mapsto (1 - \alpha) g(Q) + \alpha g(R). \qquad (8)$$

Using these notations, we have

$$\varphi(\mathbf{x}, \mathbf{x}') = C_{n,n'} + f(K_{\mathbf{X}}, K_{\mathbf{X}} + K_{\mathbf{Y}}),$$

which highlights the connection between $\varphi(\mathbf{x}, \mathbf{x}')$ and the Gram matrices $K_{\mathbf{X}}$ and $K_{\mathbf{X}} + K_{\mathbf{Y}}$. In the context of kernel methods, the natural question brought forward by this reformulation is whether the linear dot-product matrices $K_{\mathbf{X}}$ and $K_{\mathbf{Y}}$ in $\varphi$ or $k$ can be replaced by arbitrary kernel matrices $\mathbf{K_X}$ and $\mathbf{K_Y}$ between the vectors in $\mathbb{R}^{pd}$ and $\mathbb{R}^d$ enumerated in $\mathbf{X}$ and $\mathbf{Y}$, and the resulting quantity still be a valid positive definite kernel between $\mathbf{x}$ and $\mathbf{x}'$. More generally, suppose that $\mathbf{x}$ and $\mathbf{x}'$ are time series of structured objects, graphs for instance. In such a case, can Equation (7) be used to define a kernel between time series of graphs $\mathbf{x}$ and $\mathbf{x}'$ by using directly Gram matrices that measure the similarities between graphs observed in $\mathbf{x}$ and $\mathbf{x}'$? We prove here this is possible.

Let us redefine and introduce some notations to establish this result. Given a $k$-uple of points $\mathbf{u} = (u_1, \cdots, u_k)$ taken in an arbitrary set $\mathcal{U}$, and a positive kernel $\kappa$ on $\mathcal{U} \times \mathcal{U}$ we write $\mathbf{K}^\kappa(\mathbf{u})$ for the $k \times k$ Gram matrix

$$\mathbf{K}^\kappa(\mathbf{u}) \overset{\text{def}}{=} \left[ \kappa(u_i, u_j) \right]_{1 \leq i,j \leq k}.$$

For two lists $\mathbf{u}$ and $\mathbf{u}'$, we write $\mathbf{u} \cdot \mathbf{u}'$ for the concatenation of $\mathbf{u}$ and $\mathbf{u}'$. Recall that an empirical measure $\mu$ on a measurable set $\mathcal{X}$ is a finite sum of weighted Dirac masses, $\mu = \sum_{i=1}^n t_i \delta_{u_i}$, where the $u_i \in \mathcal{X}$ are the locations and the $t_i \in \mathbb{R}^+$ the weights of such masses.

**Lemma 3.** *For two empirical measures $\mu$ and $\mu'$ defined on a set $\mathcal{X}$ by locations $\boldsymbol{u} = (u_1, \cdots, u_k)$ and $\boldsymbol{u}' = (u'_1, \cdots, u'_l)$ and weights $\mathbf{t} = (t_1, \cdots, t_k) \in \mathbb{R}^k_+$ and $\mathbf{t}' = (t'_1, \ldots, t'_l) \in \mathbb{R}^l_+$ respectively, the function*

$$\xi : (\mu, \mu') \mapsto \log \left| \mathbf{K}^\kappa(\boldsymbol{u} \cdot \boldsymbol{u}') \, \mathbf{diag}(\mathbf{t} \cdot \mathbf{t}') + I_{k+l} \right|$$

*is a negative definite kernel.*

*Proof.* We follow the approach of the proof of Cuturi et al. [2005, Theorem 7]. Consider $m$ measures $\mu_1, \cdots, \mu_m$ and $m$ real weights $c_1, \cdots, c_m$ such that $\sum_{i=1}^m c_i = 0$. We prove that the quantity

$$\sum_{i,j=1}^m c_i c_j \xi(\mu_i, \mu_j), \tag{9}$$

is necessarily non-positive. Consider the finite set $\mathcal{S}$ of all locations in $\mathcal{X}$ enumerated in all measures $\mu_i$. For each point $u$ in $\mathcal{S}$, we consider the function $\kappa(u, \cdot)$ in the reproducing kernel Hilbert space $\mathcal{H}$ of functions defined by $\kappa$. Let $\mathcal{H}_\mathcal{S} \overset{\text{def}}{=} \text{span}\{\kappa(u), u \in \mathcal{S}\}$ be the finite dimensional subspace of $\mathcal{H}$ spanned by all images in $\mathcal{H}$ of elements of $\mathcal{S}$ by this mapping. For each empirical measures $\mu_i$ we consider its counterpart $\nu_i$, the empirical measure in $\mathcal{H}_\mathcal{S}$ with the same weights and locations defined as the mapped locations of $\mu_i$ in $\mathcal{H}_\mathcal{S}$. Since for two points $u_1, u_2$ in $\mathcal{S}$ we have by the reproducing property that $\langle \kappa(u_1, \cdot), \kappa(u_2, \cdot) \rangle = \kappa(u_1, u_2)$, we obtain that $\xi(\mu_i, \mu_j) = -\frac{1}{2} \log \tau(\nu_i, \nu_j)$ where $\tau$ is in this case cast as a positive definite kernel on the Euclidean space $\mathcal{H}_\mathcal{S}$. Hence the left hand side of Equation (9) is nonnegative by negative definiteness of $-\frac{1}{2} \log \tau$. ∎

We now consider two time series $\mathbf{x}$ and $\mathbf{x}'$ taking values in an arbitrary space $\mathcal{X}$. For any sequence $\mathbf{x} = (x_1, \cdots, x_n)$ we write $\mathbf{x}_i^j$ where $1 \leq i < j \leq n$ for the sequence $(x_i, x_{i+1}, \cdots, x_j)$. To summarize the transitions enumerated in $\mathbf{x}$ and $\mathbf{x}'$ we consider the sequences of subsequences

$$X = (\mathbf{x}_1^p, \mathbf{x}_2^{p+1}, \cdots, \mathbf{x}_{n-p+1}^{n-1}), \quad X' = (\mathbf{x}'^p_1, \mathbf{x}'^{p+1}_2, \cdots, \mathbf{x}'^{n'-1}_{n'-p+1}),$$

and

$$Y = (x_{p+1}, \cdots, x_n), \quad Y' = (x'_{p+1}, \cdots, x'_{n'}).$$

Considering now a p.d. kernel $\kappa_1$ on $\mathcal{X}^p$ and $\kappa_2$ on $\mathcal{X}$ we can build Gram matrices,

$$\mathbf{K}_1 = \mathbf{K}^{\kappa_1}(X \cdot X'), \quad \mathbf{K}_2 = \mathbf{K}^{\kappa_2}(Y \cdot Y').$$

**Theorem 3.** *Given two time series $\boldsymbol{x}, \boldsymbol{x}'$ in $\mathcal{X}^{\mathbb{N}}$, the autoregressive negative definite kernel $\varphi_\kappa$ of order $p$, parameter $0 < \alpha \leq 1$ and base kernels $\kappa_1$ and $\kappa_2$ defined as*

$$\varphi_\kappa(\mathbf{x}, \mathbf{x}') = C_{n,n'} + f(\mathbf{K}_1, \mathbf{K}_1 + \mathbf{K}_2),$$

*is negative definite.*

*Proof.* $\varphi_\kappa$ is negative definite as the sum of three negative definite kernels: $C_{n,n'}$ is an additive function in $n$ and $n'$ and is thus trivially negative definite. The term $(1 - \alpha) \log |K_\mathbf{X} + \Delta^{-1}|$ can be cast as $(1 - \alpha)$ times the negative definite kernel $\xi$ defined on measures of $\mathcal{X}^p$ with kernel $\kappa_1$ while the term $\alpha \log |K_\mathbf{X} + K_\mathbf{Y} + \Delta^{-1}|$ is $\alpha$ times the negative definite kernel $\xi$ defined on measures of $\mathcal{X}^p \times \mathcal{X}$ with kernel $\kappa_1 + \kappa_2$. ∎

## 3.2   Approximations Using Low-Rank Factorizations

We consider in this section matrix factorization techniques to approximate the kernel matrices $\mathbf{K}_1$ and $\mathbf{K}_1 + \mathbf{K}_2$ used to compute $\varphi_\kappa(\mathbf{x}, \mathbf{x}')$ by low rank matrices. Theorem 5 provides a useful tool to control the tradeoff between the accuracy and the computational speed of this approximation.

### 3.2.1   Computing $f$ using low-rank matrices

Consider an $N \times m_1$ matrix $\mathbf{g}_1$ and an $N \times m_2$ matrix $\mathbf{g}_2$ such that $\mathbf{G}_1 \stackrel{\text{def}}{=} \mathbf{g}_1 \mathbf{g}_1^T$ approximates $\mathbf{K}_1$ and $\mathbf{G}_2 \stackrel{\text{def}}{=} \mathbf{g}_2 \mathbf{g}_2^T$ approximates $\mathbf{K}_1 + \mathbf{K}_2$. Namely, such that the Frobenius norms of the differences

$$\varepsilon_1 \stackrel{\text{def}}{=} \mathbf{K}_1 - \mathbf{G}_1, \quad \varepsilon_2 \stackrel{\text{def}}{=} \mathbf{K}_1 + \mathbf{K}_2 - \mathbf{G}_2,$$

are small, where the Frobenius norm of a matrix $M$ is $\|M\| \stackrel{\text{def}}{=} \sqrt{\operatorname{tr} M^T M}$.

Computing $f(\mathbf{G}_1, \mathbf{G}_2)$ requires an order of $O(N(m_1 + m_2)^2 + m_1^3 + m_2^3)$ operations. Techniques to obtain such matrices $\mathbf{g}_1$ and $\mathbf{g}_2$ range from standard truncated eigenvalue decompositions, such as the power method, to incomplete Cholesky decompositions [Fine and Scheinberg, 2002, Bach and Jordan, 2005] and Nystrm methods [Williams and Seeger, 2001, Drineas and Mahoney, 2005] which are arguably the most popular in the kernel methods literature. The analysis we propose below is valid for any factorization method.

**Proposition 4.** *Let $0 \leq \alpha \leq 1$, then $f$ defined in Equation (8) is a strictly concave function of $(\mathbf{S}_N^+)^2$ which is strictly increasing in the sense that $f(Q_1, R_1) < f(Q_2, R_2)$ if $Q_2 \succ Q_1$ and $R_2 \succ R_1$.*

*Proof.* The gradient of $g : Q \mapsto \log |Q + \Delta^{-1}|$ is $\nabla g(Q) = (Q + \Delta^{-1})^{-1}$ which is thus a positive definite matrix. As a consequence, $f$ is a strictly increasing function. The Jacobian of this gradient evaluated at $Q$ is the linear map $\varepsilon \in \mathbf{S}_n^+ \mapsto -\operatorname{tr}(Q + \Delta - 1)^{-1} \varepsilon (Q + \Delta^{-1})^{-1}$. For any matrix $C \succ 0$ the Hessian of $g$ computed at $Q$ is thus the quadratic form

$$\nabla^2 g(Q) : (\varepsilon, \nu) \to -\operatorname{tr}(Q + \Delta^{-1})^{-1} \varepsilon (Q + \Delta^{-1})^{-1} \nu$$

Since $\operatorname{tr} UVUV = \operatorname{tr}((\sqrt{U} V \sqrt{U})^2) > 0$ for any two matrices $U, V \succ 0$, $\nabla^2 g(Q)(\varepsilon, \varepsilon)$ is negative for any positive definite matrix $\varepsilon$. Hence the Hessian of $f$ is minus a positive definite quadratic form on $(\mathbf{S}_N^+)^2$ and thus $f$ is strictly concave. ∎

We use a first order argument to bound the difference between the approximation and the true value of $f(\mathbf{K}_1, \mathbf{K}_1 + \mathbf{K}_2)$ using terms in $\|\varepsilon_1\|$ and $\|\varepsilon_2\|$:

$$f(\mathbf{K}_1, \mathbf{K}_1 + \mathbf{K}_2) - (1 - \alpha)\langle \nabla g(\mathbf{G}_1), \varepsilon_1 \rangle - \alpha \langle \nabla g(\mathbf{G}_2), \varepsilon_2 \rangle$$
$$\leq f(\mathbf{G}_1, \mathbf{G}_2) \leq f(\mathbf{K}_1, \mathbf{K}_1 + \mathbf{K}_2).$$

**Theorem 5.** *Given two time series* $\mathbf{x}$,$\mathbf{x}$'*, for any low rank approximations* $\mathbf{G}_1$ *and* $\mathbf{G}_2$ *in* $\mathbf{S}_N^+$ *such that* $\mathbf{G}_1 \preceq \mathbf{K}_1$ *and* $\mathbf{G}_2 \preceq \mathbf{K}_1 + \mathbf{K}_2$ *we have that*

$$e^{-\varphi_\kappa(\mathbf{x},\mathbf{x}')} \leq e^{-C_{n,n'} - f(\mathbf{G}_1, \mathbf{G}_2)} \leq (1 + \rho)e^{-\varphi_\kappa(\mathbf{x},\mathbf{x}')},$$

*where* $\rho \overset{\text{def}}{=} \exp\left((1 - \alpha)\|\nabla g(\mathbf{G}_1)\|\|\varepsilon_1\| + \alpha\|\nabla g(\mathbf{G}_2)\|\|\varepsilon_2\|\right) - 1.$

*Proof.* Immediate given that $f$ is concave and increasing ∎.

### 3.2.2 Early stopping criterion

Incomplete Cholesky decomposition and Nystrm methods can build iteratively a series of matrices $\mathbf{g}_{1,t}$ and $\mathbf{g}_{2,t} \in \mathbb{R}^{N \times t}, 1 \leq t \leq N$ such that $\mathbf{G}_{1,t} \overset{\text{def}}{=} \mathbf{g}_{1,t}\mathbf{g}_{1,t}^T$ and $\mathbf{G}_{2,t} \overset{\text{def}}{=} \mathbf{g}_{2,t}\mathbf{g}_{2,t}^T$ increase respectively towards $\mathbf{K}_1$ and $\mathbf{K}_1 + \mathbf{K}_2$ as $t$ goes to $N$. The series $\mathbf{g}_{1,t}$ and $\mathbf{g}_{2,t}$ can be obtained without having to compute explicitly the whole of $\mathbf{K}_1$ nor $\mathbf{K}_1 + \mathbf{K}_2$ except for their diagonal.

The iterative computations of $\mathbf{G}_{1,t}$ and $\mathbf{G}_{2,t}$ can be halted whenever an upper bound for each of the norms $\|\varepsilon_{1,t}\|$ and $\|\varepsilon_{2,t}\|$ of the residues $\varepsilon_{1,t} \overset{\text{def}}{=} \mathbf{K}_1 - \mathbf{G}_{1,t}$ and $\varepsilon_{2,t} \overset{\text{def}}{=} \mathbf{K}_1 + \mathbf{K}_2 - \mathbf{G}_{2,t}$ goes below an approximation threshold.

Theorem 5 can be used to produce such a stopping criterion by a rule which combines an upper bound on $\|\varepsilon_{1,t}\|$ and $\|\varepsilon_{2,t}\|$ and the exact norm of the gradients of $g$ at $\mathbf{G}_{1,t}$ and $\mathbf{G}_{2,t}$. This would require computing Frobenius norms of the matrices $(\mathbf{G}_{i,t} + \Delta^{-1})^{-1}, i = 1, 2$. These matrices can be updated iteratively using rank-one updates. A simpler alternative which we consider is to bound $\|\nabla g(\mathbf{G}_{i,t})\|$ uniformly between 0 and $\mathbf{K}$ using the inequality

$$(\mathbf{G}_{i,t} + \Delta^{-1})^{-1} \preceq \Delta, i = 1, 2.$$

which yields the following bound:

$$e^{-\varphi_\kappa(\mathbf{x},\mathbf{x}')} \leq e^{-C_{n,n'} - f(\mathbf{G}_1, \mathbf{G}_2)} \leq e^{-\varphi_\kappa(\mathbf{x},\mathbf{x}')}e^{\frac{1}{2}\sqrt{\frac{N}{(n-p)(n'-p)}}((1-\alpha)\|\varepsilon_1\| + \alpha\|\varepsilon_2\|)}$$

We consider in the experimental section the positive definite kernel $e^{-t\varphi_\kappa}$, that is the scaled exponentiation of $\varphi_\kappa$ multiplied by a bandwidth parameter $t > 0$. Setting a target tolerance $\sigma > 0$ on the ratio between the approximation of $e^{-t\varphi_\kappa}$ and its true value, namely requiring that

$$e^{-t\varphi_\kappa(\mathbf{x},\mathbf{x}'))} \leq e^{-t\left(C_{n,n'} + f(\mathbf{G}_1, \mathbf{G}_2)\right)} \leq (1 + \sigma)e^{-t\varphi_\kappa(\mathbf{x},\mathbf{x}'))},$$

can be ensured by stopping the factorizations at an iteration $t$ such that

$$(1 - \alpha)\|\varepsilon_{1,t}\| + \alpha\|\varepsilon_{2,t}\| \leq \frac{2\log(1 + \tau)}{t}\sqrt{\frac{(n-p)(n'-p)}{N}}.$$

12

which we simplify to performing the factorizations separately, and stopping at the lowest iterations $t_1$ and $t_2$ such that

$$
\begin{aligned}
\|\varepsilon_{1,t_1}\| &\leq \frac{\log(1+\tau)}{(1-\alpha)t}\sqrt{\frac{(n-p)(n'-p)}{N}}, \\
\|\varepsilon_{2,t_2}\| &\leq \frac{\log(1+\tau)}{\alpha t}\sqrt{\frac{(n-p)(n'-p)}{N}}.
\end{aligned}
\tag{10}
$$

We provide in Figure 2 of Section 4.4 an experimental assessment of this speed/accuracy tradeoff when computing the value of $\varphi_\kappa$.

# 4 Experiments

We provide in this section a fair assessment of the performance and efficiency of autoregressive kernels on different tasks. We detail in Section 4.1 the different kernels we consider in this benchmark. Section 4.2 and 4.3 introduce the toy and real-life datasets of this benchmark, results are presented in Section 4.4 before reaching the conclusion of this paper.

## 4.1 Kernels and parameter tuning

The kernels we consider in this experimental section are all of the form $K = e^{-\frac{1}{t}\Phi}$, where $\Phi$ is a negative definite kernel. We select for each kernel $K$ the value of the bandwidth $t$ as the median value $\hat{m}_\Phi$ of $\Phi$ on all pairs of time series observed in the training fold times 0.5, 1 or 2, namely $t \in \{.5\,\hat{m}_\Phi, \hat{m}_\Phi, 2\,\hat{m}_\Phi\}$. The selection is based on the cross validation error on the training fold for each (kernel,dataset) pair. Some kernels described below bear the superscript $\cdot^\kappa$, which means that they are parameterized by a base kernel $\kappa$. Given two times series $\mathbf{x} = (x_1, \cdots, x_n)$ and $\mathbf{x}' = (x'_1, \cdots, x'_{n'})$, this base kernel $\kappa$ is used to computed similarities between single components $\kappa(x_i, x'_j)$ or p-uples of components $\kappa((x_{i+1}, \cdots, x_{i+p}), (x'_{j+1}, \cdots, x'_{j+p}))$. For all superscripted kernels below, $\kappa$ is set to be the Gaussian kernel between two vectors $\kappa(x, y) = e^{-\|x-y\|^2/(2\sigma^2)}$, where the dimension is obvious from the context and is either $d$ or $pd$. The variance parameter $\sigma^2$ is arbitrarily set to be the median value of all Euclidean distances $\|x_i^{(r)} - x_j^{(s)}\|$ where $i \leq |\mathbf{x}^{(r)}|$, $j \leq |\mathbf{x}^{(s)}|$, $(r, s) \in \mathcal{R}^2$, where $\mathcal{R}$ is a random subset of $\{1, 2, \cdots, \#\text{training points}\}$

**Autoregressive kernels** $k_{\text{ar}}, k_{\text{ar}}^\kappa$: We consider the kernels

$$
k_{\text{ar}} = e^{-t_{\text{ar}}\varphi}, \quad k_{\text{ar}}^\kappa = e^{-t_{\text{ar}}^\kappa \varphi_\kappa},
$$

parameterized by the bandwidths $t_{\text{ar}}$ and $t_{\text{ar}}^\kappa$. We set parameters $\alpha$ and $p$ as $\alpha = 1/2$ and $p = 5$ in all experiments. The matrix factorizations used to compute approximations to $k_{\text{ar}}^\kappa$ (with $\tau$ set to $10^{-4}$) can be performed using the `chol_gauss` routine proposed by Shen et al. [2009]. Since running separately this routine for $\mathbf{K}_1$ and $\mathbf{K}_1 + \mathbf{K}_2$ results in duplicate computations of portions of $\mathbf{K}_1$, we have added our modifications to this routine in order to cache values of $\mathbf{K}_1$ that can be reused when evaluating $\mathbf{K}_1 + \mathbf{K}_2$. The routine `TwoCholGauss` is available on our website, as well as other pieces of code. We insist on the fact that, other than $\alpha, p$ and the temperature $t_{\text{ar}}$, the autoregressive kernel $k_{\text{ar}}$ does not require any parameter tuning.

**Bag of vectors kernel** $k_{\text{BoV}}^\kappa$: A time series $(x_1, \cdots, x_n)$ can be considered as a bag of vectors $\{x_1, \cdots, x_n\}$ where the time-dependent information from each state's timestamp is deliberately ignored. Two time series $\mathbf{x}$ and $\mathbf{x}'$ can be compared through their respective bags $\{x_1, \cdots, x_n\}$ and $\{x'_1, \cdots, x'_{n'}\}$ by setting

$$\psi_\kappa(\mathbf{x}, \mathbf{x}') \stackrel{\text{def}}{=} \frac{1}{nn'} \sum_{i \leq n, j \leq n'} \kappa(x_i, x'_j),$$

and defining the kernel $k_{\text{BoV}}^\kappa(\mathbf{x}, \mathbf{x}') = \exp\left(-t_{\text{BoV}}\left(\psi_\kappa(\mathbf{x}, \mathbf{x}) + \psi_\kappa(\mathbf{x}', \mathbf{x}') - 2\psi_\kappa(\mathbf{x}, \mathbf{x}')\right)\right)$ where $t_{\text{BoV}} > 0$, see for instance [Hein and Bousquet, 2005]. This relatively simple kernel will act as the baseline of our experiments, both for performance and computational time.

**Global alignment kernel** $k_{\text{GA}}^\kappa$: The global alignment kernel [Cuturi et al., 2007] is a positive definite kernel that builds upon the dynamic time warping framework, by considering the soft-maximum of the alignment score of all possible alignments between two time series. We use an implementation of this kernel distributed on the web, and consider the kernel $k_{\text{GA}}^\kappa = \exp(-t_{\text{GA}}\text{GlobalAlignment}_\kappa(\mathbf{x}, \mathbf{x}'))$, parameterized by the bandwidth $t_{\text{GA}}$. Note that the global alignment kernel has not been proved to be infinitely divisible. Namely, $k_{\text{GA}}^\kappa$ is known to be positive definite for $t_{\text{GA}} = 1$, and as a consequence for $t_{\text{GA}} \in \mathbb{N}$, but not for all positive values. However, the Gram matrices that were generated in these experiments have been found to be positive definite for all values of $t_{\text{GA}}$ as discussed earlier in this section. This suggests that $k_{\text{GA}}^\kappa$ might indeed be infinitely divisible.

**Splines Smoothing kernel** $k_{\text{S}}$: Kumara et al. [2008] use spline smoothing techniques to map each time series $(x_1, \cdots, x_n)$ onto a multivariate polynomial function $p_\mathbf{x}$ defined on $[0, 1]$. As a pre-processing step, each time series is mapped onto a multivariate time series of arbitrary length $\tilde{\mathbf{x}}$ (set to 200 in our experiments) such that $\tilde{\mathbf{x}}^T \tilde{\mathbf{x}}'$ corresponds to a relevant dot-product for these polynomials. We have modified an implementation that we received from the authors in email correspondence. Kumara et al. [2008] consider a linear kernel in their original paper on such representations. We have found that a Gaussian kernel between these two vector representations performs better and use $k_{\text{S}} = \exp(-t_{\text{S}}\|\tilde{\mathbf{x}} - \tilde{\mathbf{x}}'\|^2)$.

**Remark 5.** *Although promising, the kernel proposed by Jebara et al. [2004, §4.5] is embryonic and leaves many open questions on its practical implementation. A simple implementation using VAR models would not work with these experiments, since for many datasets the dimension d of the considered time series is comparable to or larger than their lengths' and would prevent any estimation of the $(pd^2 + d(d+1)/2)$ parameters of a VAR(p) model. A more advanced implementation not detailed in the original paper would be beyond the scope of this work. We have also tried to implement the fairly complex families of kernels described by Vishwanathan et al. [2007], namely Equations (10) and (16) in that reference, but our implementations performed very poorly on the datasets we considered, and we hence decided not to report these results out of concerns for the validity of our codes. Despite repeated attempts, we could not obtain computer codes for these kernels from the authors, a problem also reported in [Lin et al., 2008].*

## 4.2 Toy dataset

We study the performance of these kernels in a simple binary classification toy experiment that illustrates some of the merits of autoregressive kernels. We consider high dimensional

time series ($d = 1000$) dimensional of short length ($n = 10$) generated randomly using one of two VAR(1) models,

$$x_{t+1} = A_i x_t + \varepsilon_t, \quad i = 1, 2,$$

where the process $\varepsilon_t$ is a white Gaussian noise with covariance matrix $.1I_{1000}$. Each time series' initial point is a random vector whose components are each distributed randomly following the uniform distribution in $[-5, 5]$. The two matrices $A_i$, $i = 1, 2$, are sparse (10% of non-zero values, that is 100.000 non zero entries out of potentially one million) and have entries that are randomly distributed following a standard Gaussian law[2]. These matrices are divided by their spectral radius to ensure that their largest eigenvalue has norm smaller than one to ensure their stationarity.

We draw randomly 10 time series with transition matrix $A_1$ and 10 times with transition matrix $A_2$ and use these 20 time series as a training set to learn an SVM that can discriminate between time series of type 1 or 2. We draw 100 test time series for each class $i = 1, 2$, that is a total of 200 time series, and test the performance of all kernels following the protocol outlined above. The test error is represented in the leftmost bar plot of Figure 1. The autoregressive kernel $k_{\mathrm{ar}}$ achieves a remarkable test error of 0, whereas other kernels, including $k_{\mathrm{ar}}^{\kappa}$, make a not-so-surprisingly larger number of mistakes, given the difficulty of this task. One of the strongest appeals of the autoregressive kernel $k_{\mathrm{ar}}$ is that it manages to quantify a dynamic similarity between two time series (something that neither the Kumara kernel or any kernel based on alignments may achieve with so few samples) without resorting to the actual estimation of a density, which would of course be impossible given the samples' length.

## 4.3   Real-life datasets

We assess the performance of the kernels proposed in this paper using different benchmark datasets and other known kernels for time series. The datasets are all taken from the UCI Machine Learning repository [Frank and Asuncion, 2010], except for the PEMS dataset which we have compiled. The datasets characteristics' are summarized in Table 1.

**Japanese Vowels**: The database records utterances by nine male speakers of two Japanese vowels 'a' and 'e' successively. Each utterance is described as a time series of LPC cepstrum coefficients. The length of each time series lies within a range of 7 to 29 observations, each observation being a vector of $\mathbb{R}^{12}$. The task is to guess which of the nine speakers pronounces a new utterance of 'a' or 'e'. We use the original split proposed by the authors, namely 270 utterances for training and 370 for testing.

**Libras Movement Data Set**: LIBRAS is the acronym for the brazilian sign language. The observations are 2-dimensional time series of length 45. Each time series describes the location of the gravity center of a hand's coordinates in the visual plane. 15 different signs are considered, the training set has 24 instances of each class, for a total $360 = 24 \times 15$ time series. We consider another dataset of 585 time series for the test set.

**Handwritten characters**: 2858 recordings of a pen tip trajectory were taken from the same writer. Each trajectory, a $3 \times n$ matrix where $n$ varies between 60 and 182 records the location of the pen and its tip force. Each trajectory belongs to one out of 20 different classes. The data is split into 2 balanced folds of 600 examples for training and 2258 examples for testing.

**Australian Language of Signs**: Sensors are set on the two hands of a native signer communicating with the AUSLAN sign language. There are 11 sensors on each hand and

---

[2]In Matlab notation, $A = $ `sprandn(1000,.1)`

| Database | $d$ | $n$ | classes | # train | #test |
|---|---|---|---|---|---|
| Toy dataset | 1000 | 10 | 2 | 20 | 200 |
| Japanese Vowels | 12 | 7-29 | 9 | 270 | 370 |
| Libras | 2 | 45 | 15 | 360 | 585 |
| Handwritten Characters | 3 | 60-182 | 20 | 600 | 2258 |
| AUSLAN | 22 | 45-136 | 95 | 600 | 1865 |
| PEMS | 963 | 144 | 7 | 267 | 173 |

**Table 1:** Characteristics of the different databases considered in the benchmark test

hence 22 coordinates for each observation of the time series. The length of each time series ranges from 45 to 136 measurements. A sample of 27 distinct recordings is performed for each of the 95 considered signs, which totals 2565 time series. These are split between balanced train and tests sets of size 600 and 1865 respectively. Each time series in both test and training sets is centered individually, that is $\mathbf{x}^{(i)}$ is replaced by $\mathbf{x}^{(i)} - \bar{\mathbf{x}}^{(i)}$. Without such a centering the performance of all kernels is seriously degraded, except for the autoregressive kernel which remains very competitive with an error below 10%.

**PEMS Database**: We have downloaded 15 months worth of daily data from the California Department of Transportation PEMS website[3]. The data describes measurements at 10 minute intervals of occupancy rate, between 0 and 1, of different car lanes of the San Francisco bay area (D04) freeway system. The measurements cover the period from January 1st 2008 to March 30th 2009. We consider each day of measurements as a single time series of dimension 963 (the number of sensors which functioned consistently throughout the studied period) and length $6 \times 24 = 144$. The task is to classify each day as the correct day of the week, from Monday to Sunday, e.g. label it with an integer between 1 and 7. We remove public holidays from the dataset, as well as two days with anomalies (March 8 2009 and March 9 2008) where all sensors have been seemingly turned off between 2:00 and 3:00 AM. This leaves 440 time series in total, which are shuffled and split between 267 training observations and 173 test observations. We plan to donate this dataset to the UCI repository, and it should be available shortly. In the meantime, the dataset can be accessed in Matlab format on our website.

## 4.4   Results and computational speed

The kernels introduced in the Section above are paired with a standard multiclass SVM implementation using a one-versus-rest approach. For each kernel and training set pair, the SVM constant $C$ to be used on the test set was chosen as either 1, 10 or 100, whichever gave the lowest cross-validation mean-error on the training fold. We report the test errors in Figure 1. The errors on the test sets can be also compared with the average computation time per kernel evaluation graph displayed for 4 datasets in Figure2. In terms of performance, the autoregressive kernels perform favorably with respect to other kernels, notably the Global Alignment kernel, which is usually very difficult to beat. Their computational time offer a diametrically opposed perspective since for these benchmarks datasets the flexibility of using a kernel $\kappa$ to encode the inputs in a RKHS does not yield practical gains in performance but has a tremendously high computational price. On the contrary, $k_{\mathrm{ar}}$ is both efficient and usually very fast compared to the other kernels.
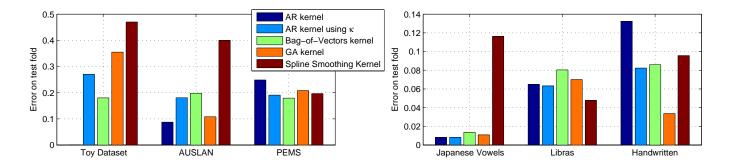
---

[3] http://pems.dot.ca.gov

**Figure 1:** Test error of the 5 considered kernels on 5 different tasks split into two panels for better legibility of the error rates (notice the difference in scale). The AR kernel has a test-error of 0 on the toy dataset's test fold.

## 4.5 Conclusion and Discussion

We have proposed in this work two infinitely divisible kernels $k_{ar}$ and $k_{ar}^{\kappa}$ for time series. These kernels can be used within the framework of kernel machines, *e.g.* the SVM or kernel-PCA, or more generally as Hilbertian distances by using directly their logarithms $\varphi_{var}$ and $\varphi_{var}^{\kappa}$ once properly normalized. The first kernel $k_{ar}$ computes a similarity between two multivariate time series with a low computational cost. This similarity is easy to implement, easy to tune given its infinite divisibility, and often performs similarly or better than more costly alternatives. The second kernel, $k_{ar}^{\kappa}$, is a generalization of $k_{ar}$ that can handle structured data by considering a local kernel $\kappa$ on the structures. Its computation requires the computation of all or a part of large Gram matrices as well as the determinant of these. Given its computational drawbacks, the experimental evidence gathered in this paper is not sufficient to advocate its use on vectorial data. Moreover, El Karoui [2010] has recently shown that the spectrum of a Gram matrix of high-dimensional points using the Gaussian kernel may, under certain assumptions, be very similar to the spectrum of the standard Gram matrix of these same points using the linear dot-product. In such a case, the sophistication brought forward by $k_{ar}^{\kappa}$ might be gratuitous and yield similar results to the direct use of $k_{ar}$. However, we believe that $k_{ar}^{\kappa}$ may prove particularly useful when considering time series of structured data. For instance, we plan to apply the kernel $k_{ar}^{\kappa}$ to the classification of video segments, where each segment would be represented as a time varying histogram of features and $\kappa$ a suitable kernel on histograms that can take into account the similarity of features themselves.

Our contribution follows the blueprint laid down by Seeger [2002] which can be effectively applied to other exponential models. However, we believe that the infinite divisibility of such kernels, which is crucial in practical applications, had not been considered before this work. Our result in this respect is not as general as we would wish for, in the sense that we do not know whether $k_{ar}$ remains infinitely divisible when the degrees of freedom $\lambda$ of the inverse Wishart prior exceed $d - 1$. In such a case, the concavity of $f$ in Equation (8) would not be given either. Finally, although the prior that we use to define $k_{ar}$ is non informative, it might be of interest to learn the hyperparameters for these priors based on a data corpus of interest.
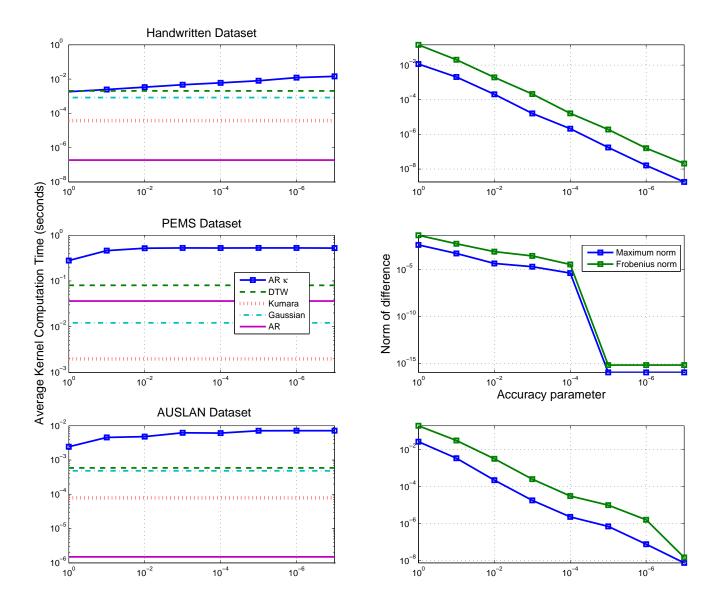
**Figure 2:** These graphs provide on the left side the average time needed to compute one evaluation of each of the 5 kernels on the largest datasets. The average speeds (computed over a sample of $50 \times 50$ calculations) for each kernel are quantified in the y-axis. The x-axis is only effective for the kernel $k_{ar}^{\kappa}$ and shows the influence of the accuracy para meter $\tau$ on that speed using the low-rank matrix factorization expression used for $\varphi$ as described in Equation (10). This parameter is set between $10^0$ (poor approximation) to $10^{-7}$ (high accuracy). The accuracy is measured by the maximum norm and the Frobenius norm of the difference between the two $50 \times 50$ $\varphi$-Gram matrices. Note that $k_{GA}^{\kappa}$ is fully implemented in mex-C code, $k_{ar}^{\kappa}$ uses mex-C subroutines for Cholesky decomposition while all other kernels are implemented using standard algebra in Matlab. Preprocessing times are not counted in these averages. The simulations were run using an iMac 2.66 GhZ Intel Core with 4Gb of memory.

# References

F. R. Bach and M. I. Jordan. Predictive low-rank decomposition for kernel methods. In *Proceedings of ICML '05: Twenty-second international conference on Machine learning*. ACM Press, 2005.

C. Bahlmann, B. Haasdonk, and H. Burkhardt. Online handwriting recognition with support vector machines-a kernel approach. In *Frontiers in Handwriting Recognition, 2002. Proceedings. Eighth International Workshop on*, pages 49–54, 2002.

C. Berg, J. P. R. Christensen, and P. Ressel. *Harmonic Analysis on Semigroups*. Number 100 in Graduate Texts in Mathematics. Springer Verlag, 1984.

A. Berlinet and C. Thomas-Agnan. *Reproducing Kernel Hilbert Spaces in Probability and Statistics*. Kluwer Academic Publishers, 2003.

K. Borgwardt, S. Vishwanathan, and H. Kriegel. Class prediction from time series gene expression profiles using dynamical systems kernels. In *Proceedings of the 11th Pacific Symposium on Biocomputing*, pages 547–558, 2006.

O. Chapelle, P. Haffner, and V. Vapnik. SVMs for histogram based image classification. *IEEE Transactions on Neural Networks*, 10(5):1055, Sept. 1999.

C. Cortes, P. Haffner, and M. Mohri. Rational kernels: Theory and algorithms. *The Journal of Machine Learning Research*, 5:1035–1062, 2004.

M. Cuturi and K. Fukumizu. Kernels on structured objects through nested histograms. In B. Schölkopf, J. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems 19*. MIT Press, Cambridge, MA, 2007.

M. Cuturi and J.-P. Vert. The context-tree kernel for strings. *Neural Networks*, 18(8), 2005.

M. Cuturi, K. Fukumizu, and J.-P. Vert. Semigroup kernels on measures. *Journal of Machine Learning Research*, 6:1169–1198, 2005.

M. Cuturi, J.-P. Vert, Øystein. Birkenes, and T. Matsui. A kernel for time series based on global alignments. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, volume II, pages 413 – 416, 2007.

P. Drineas and M. W. Mahoney. On the Nyström method for approximating a gram matrix for improved kernel-based learning. *Journal of Machine Learning Research*, 6:2153–2175, 2005. ISSN 1532-4435.

N. El Karoui. The spectrum of kernel random matrices. *The Annals of Statistics*, 38(1):1–50, 2010. ISSN 0090-5364.

S. Fine and K. Scheinberg. Efficient SVM training using low-rank kernel representations. *Journal of Machine Learning Research*, 2:243–264, 2002.

A. Frank and A. Asuncion. UCI machine learning repository, http://archive.ics.uci.edu/ml, 2010.

K. Grauman and T. Darrell. The pyramid match kernel: Discriminative classification with sets of image features. In *ICCV*, pages 1458–1465. IEEE Computer Society, 2005.

B. Haasdonk and C. Bahlmann. Learning with distance substitution kernels. *Pattern Recognition, Proc. of the 26th DAGM Symposium*, pages 220–227, 2004.

Z. Harchaoui and F. Bach. Image classification with segmentation graph kernels. In *CVPR*, 2007.

D. Haussler. Convolution kernels on discrete structures. Technical report, UC Santa Cruz, 1999. USCS-CRL-99-10.

A. Hayashi, Y. Mizuhara, and N. Suematsu. Embedding time series data for classification. *Machine Learning and Data Mining in Pattern Recognition*, pages 356–365, 2005.

M. Hein and O. Bousquet. Hilbertian metrics and positive definite kernels on probability measures. In Z. Ghahramani and R. Cowell, editors, *Proceedings of AISTATS 2005*, January 2005.

T. Hofmann, B. Scholkopf, and A. Smola. Kernel methods in machine learning. *Annals of Statistics*, 36(3):1171, 2008.

T. Jaakkola, M. Diekhaus, and D. Haussler. Using the fisher kernel method to detect remote protein homologies. *7th Intell. Sys. Mol. Biol.*, pages 149–158, 1999.

T. Jebara, R. Kondor, and A. Howard. Probability product kernels. *Journal of Machine Learning Research*, 5:819–844, 2004.

T. Joachims. *Learning to Classify Text Using Support Vector Machines: Methods, Theory, and Algorithms*. Kluwer Academic Publishers, 2002.

H. Kashima, K. Tsuda, and A. Inokuchi. Marginalized kernels between labeled graphs. In T. Faucett and N. Mishra, editors, *Proceedings of the Twentieth International Conference on Machine Learning*, pages 321–328. AAAI Press, 2003.

K. Kumara, R. Agrawal, and C. Bhattacharyya. A large margin approach for writer independent online handwriting classification. *Pattern Recognition Letters*, 29(7):933–937, 2008.

G. Lebanon. Metric learning for text documents. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(4):497–508, 2006. ISSN 0162-8828.

C. Leslie, E. Eskin, and W. S. Noble. The spectrum kernel: a string kernel for svm protein classific ation. In *Proc. of PSB 2002*, pages 564–575, 2002.

P. Lévy. The arithmetical character of the Wishart distribution. In *Mathematical Proceedings of the Cambridge Philosophical Society*, volume 44, pages 295–297. Cambridge Univ Press, 1948.

T. Lin, N. Kaminski, and Z. Bar-Joseph. Alignment and classification of time series gene expression in clinical studies. *Bioinformatics*, 24(13):i147, 2008.

H. Lütkepohl. *New Introduction to Multiple Time Series Analysis*. Springer, 2005.

P. Mahe, N. Ueda, Akutsu, J.-L. T., Perret, and J.-P. Vert. Graph kernels for molecular structure-activity relationship analysis with support vector machines. *Journal of Chemical Information and Modeling*, 45(4):939–951, 2005.

A. Moschitti and F. Zanzotto. Fast and effective kernels for relational learning from texts. In *Proceedings of the 24th international conference on Machine learning*, pages 649–656. ACM, 2007.

L. Rabiner and B. Juang. *Fundamentals of speech recognition*, volume 103. Prentice hall Englewood Cliffs, New Jersey, 1993.

H. Sakoei and S. Chiba. Dynamic programming algorithm optimization for spoken word recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 26:43–49, 1978.

T. Schreiber and A. Schmitz. Classification of time series data with nonlinear similarity measures. *Physical Review Letters*, 79(8):1475–1478, 1997.

M. Seeger. Covariance kernels from bayesian generative models. In *Advances in Neural Information Processing Systems 14*, pages 905–912. MIT Press, 2002.

H. Shen, S. Jegelka, and A. Gretton. Fast kernel-based independent component analysis. *IEEE Transactions on Signal Processing*, 57(9):3498–3511, 2009.

N. Shervashidze and K. Borgwardt. Fast subtree kernels on graphs. *Advances in Neural Information Processing Systems 22*, 2009.

H. Shimodaira, K.-I. Noma, M. Nakai, and S. Sagayama. Dynamic time-alignment kernel in support vector machine. In T. G. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems 14*, Cambridge, MA, 2002. MIT Press.

S. Sonnenburg, K. Rieck, F. F. Ida, and G. Rtsch. Large scale learning with string kernels. In *Large Scale Kernel Machines*, pages 73–103. MIT Press, 2007.

J.-P. Vert, H. Saigo, and T. Akutsu. Local alignment kernels for protein sequences. In B. Schölkopf, K. Tsuda, and J.-P. Vert, editors, *Kernel Methods in Computational Biology*. MIT Press, 2004.

S. Vishwanathan and A. Smola. Binet-cauchy kernels. *Advances in Neural Information Processing Systems*, 17, 2004.

S. Vishwanathan, A. Smola, and R. Vidal. Binet-cauchy kernels on dynamical systems and its application to the analysis of dynamic scenes. *International Journal of Computer Vision*, 73(1):95–119, 2007.

S. Vishwanathan, K. Borgwardt, I. Kondor, and N. Schraudolph. Graph kernels. *Journal of Machine Learning Research*, 9:1–37, 2008.

M. West and J. Harrison. *Bayesian forecasting and dynamic models*. Springer Verlag, 1997.

C. Williams and M. Seeger. Using the Nyström method to speed up kernel machines. In *Advances in Neural Information Processing Systems 13*, pages 682–688. MIT Press, 2001.

F. Zhou, F. De la Torre, and J. Cohn. Unsupervised discovery of facial events. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 2574–2581. IEEE, 2010.