

Reducing statistical time-series problems to binary classification

Daniil Ryabko, Jérémie Mary

► **To cite this version:**

Daniil Ryabko, Jérémie Mary. Reducing statistical time-series problems to binary classification. NIPS, Dec 2012, Lake Tahoe, United States. pp.2069–2077, 2012. <hal-00675637v5>

HAL Id: hal-00675637

<https://hal.inria.fr/hal-00675637v5>

Submitted on 7 Jun 2013

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Reducing statistical time-series problems to binary classification

Daniil Ryabko
SequeL-INRIA/LIFL-CNRS,
Université de Lille, France
daniil@ryabko.net

Jérémie Mary
SequeL-INRIA/LIFL-CNRS,
Université de Lille, France
Jeremie.Mary@inria.fr

Abstract

We show how binary classification methods developed to work on i.i.d. data can be used for solving statistical problems that are seemingly unrelated to classification and concern highly-dependent time series. Specifically, the problems of time-series clustering, homogeneity testing and the three-sample problem are addressed. The algorithms that we construct for solving these problems are based on a new metric between time-series distributions, which can be evaluated using binary classification methods. Universal consistency of the proposed algorithms is proven under most general assumptions. The theoretical results are illustrated with experiments on synthetic and real-world data.

1 Introduction

Binary classification is one of the most well-understood problems of machine learning and statistics: a wealth of efficient classification algorithms has been developed and applied to a wide range of applications. Perhaps one of the reasons for this is that binary classification is conceptually one of the simplest statistical learning problems. It is thus natural to try and use it as a building block for solving other, more complex, newer or just different problems; in other words, one can try to obtain efficient algorithms for different learning problems by reducing them to binary classification. This approach has been applied to many different problems, starting with multi-class classification, and including regression and ranking [3, 16], to give just a few examples. However, all of these problems are formulated in terms of independent and identically distributed (i.i.d.) samples. This is also the assumption underlying the theoretical analysis of most of the classification algorithms.

In this work we consider learning problems that concern time-series data for which independence assumptions do not hold. The series can exhibit arbitrary long-range dependence, and different time-series samples may be interdependent as well. Moreover, the learning problems that we consider — the three-sample problem, time-series clustering, and homogeneity testing — at first glance seem completely unrelated to classification.

We show how the considered problems can be reduced to binary classification methods. The results include asymptotically consistent algorithms, as well as finite-sample analysis. To establish the consistency of the suggested methods, for clustering and the three-sample problem the only assumption that we make on the data is that the distributions generating the samples are stationary ergodic; this is one of the weakest assumptions used in statistics. For homogeneity testing we have to make some mixing assumptions in order to obtain consistency results (this is indeed unavoidable [22]). Mixing conditions are also used to obtain finite-sample performance guarantees for the first two problems.

The proposed approach is based on a new distance between time-series distributions (that is, between probability distributions on the space of infinite sequences), which we call *telescope distance*. This distance can be evaluated using binary classification methods, and its finite-sample estimates are shown to be asymptotically consistent. Three main building blocks are used to construct the tele-

scope distance. The first one is a distance on finite-dimensional marginal distributions. The distance we use for this is the following: $d_{\mathcal{H}}(P, Q) := \sup_{h \in \mathcal{H}} |\mathbf{E}_P h - \mathbf{E}_Q h|$ where P, Q are distributions and \mathcal{H} is a set of functions. This distance can be estimated using binary classification methods, and thus can be used to reduce various statistical problems to the classification problem. This distance was previously applied to such statistical problems as homogeneity testing and change-point estimation [14]. However, these applications so far have only concerned i.i.d. data, whereas we want to work with highly-dependent time series. Thus, the second building block are the recent results of [1, 2], that show that empirical estimates of $d_{\mathcal{H}}$ are consistent (under certain conditions on \mathcal{H}) for arbitrary stationary ergodic distributions. This, however, is not enough: evaluating $d_{\mathcal{H}}$ for (stationary ergodic) time-series distributions means measuring the distance between their finite-dimensional marginals, and not the distributions themselves. Finally, the third step to construct the distance is what we call *telescoping*. It consists in summing the distances for all the (infinitely many) finite-dimensional marginals with decreasing weights.

We show that the resulting distance (telescope distance) indeed can be consistently estimated based on sampling, for arbitrary stationary ergodic distributions. Further, we show how this fact can be used to construct consistent algorithms for the considered problems on time series. Thus we can harness binary classification methods to solve statistical learning problems concerning time series.

To illustrate the theoretical results in an experimental setting, we chose the problem of time-series clustering, since it is a difficult unsupervised problem which seems most different from the problem of binary classification. Experiments on both synthetic and real-world data are provided. The real-world setting concerns brain-computer interface (BCI) data, which is a notoriously challenging application, and on which the presented algorithm demonstrates competitive performance.

A related approach to address the problems considered here, as well some related problems about stationary ergodic time series, is based on (consistent) empirical estimates of the distributional distance, see [23, 21, 13] and [8] about the distributional distance. The empirical distance is based on counting frequencies of bins of decreasing sizes and “telescoping.” A similar telescoping trick is used in different problems, e.g. sequence prediction [19]. Another related approach to time-series analysis involves a different reduction, namely, that to data compression [20].

Organisation. Section 2 is preliminary. In Section 3 we introduce and discuss the telescope distance. Section 4 explains how this distance can be calculated using binary classification methods. Sections 5 and 6 are devoted to the three-sample problem and clustering, respectively. In Section 7, under some mixing conditions, we address the problems of homogeneity testing, clustering with unknown k , and finite-sample performance guarantees. Section 8 presents experimental evaluation.

2 Notation and definitions

Let $(\mathcal{X}, \mathcal{F}_1)$ be a measurable space (the domain), and denote $(\mathcal{X}^k, \mathcal{F}_k)$ and $(\mathcal{X}^{\mathbb{N}}, \mathcal{F})$ the product probability space over \mathcal{X}^k and the induced probability space over the one-way infinite sequences taking values in \mathcal{X} . Time-series (or process) distributions are probability measures on the space $(\mathcal{X}^{\mathbb{N}}, \mathcal{F})$. We use the abbreviation $X_{1..k}$ for X_1, \dots, X_k . A set \mathcal{H} of functions is called *separable* if there is a countable set \mathcal{H}' of functions such that any function in \mathcal{H} is a pointwise limit of a sequence of elements of \mathcal{H}' .

A distribution ρ is stationary if $\rho(X_{1..k} \in A) = \rho(X_{n+1..n+k} \in A)$ for all $A \in \mathcal{F}_k$, $k, n \in \mathbb{N}$. A stationary distribution is called (stationary) ergodic if $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1..n-k+1} \mathbb{I}_{X_{i..i+k} \in A} = \rho(A)$ ρ -a.s. for every $A \in \mathcal{F}_k$, $k \in \mathbb{N}$. (This definition, which is more suited for the purposes of this work, is equivalent to the usual one expressed in terms of invariant sets, see, e.g., [8].)

3 A distance between time-series distributions

We start with a distance between distributions on \mathcal{X} , and then we will extend it to distributions on $\mathcal{X}^{\mathbb{N}}$. For two probability distributions P and Q on $(\mathcal{X}, \mathcal{F}_1)$ and a set \mathcal{H} of measurable functions on \mathcal{X} , one can define the distance

$$d_{\mathcal{H}}(P, Q) := \sup_{h \in \mathcal{H}} |\mathbf{E}_P h - \mathbf{E}_Q h|.$$

This metric has been studied since at least [26]; its special cases include Kolmogorov-Smirnov [15], Kantorovich-Rubinstein [11] and Fortet-Mourier [7] metrics. Note that the distance function so defined may not be measurable; however, it is measurable under mild conditions which we assume when necessary. In particular, separability of \mathcal{H} is a sufficient condition (separability is required in most of the results below).

We will be interested in the cases where $d_{\mathcal{H}}(P, Q) = 0$ implies $P = Q$. Note that in this case $d_{\mathcal{H}}$ is a metric (the rest of the properties are easy to see). For reasons that will become apparent shortly (see Remark below), we will be mainly interested in the sets \mathcal{H} that consist of indicator functions. In this case we can identify each $f \in \mathcal{H}$ with the indicator set $\{x : f(x) = 1\} \subset \mathcal{X}$ and (by a slight abuse of notation) write $d_{\mathcal{H}}(P, Q) := \sup_{h \in \mathcal{H}} |P(h) - Q(h)|$. In this case it is easy to check that the following statement holds true.

Lemma 1. $d_{\mathcal{H}}$ is a metric on the space of probability distributions over \mathcal{X} if and only if \mathcal{H} generates \mathcal{F}_1 .

The property that \mathcal{H} generates \mathcal{F}_1 is often easy to verify directly. First of all, it trivially holds for the case where \mathcal{H} is the set of halfspaces in a Euclidean \mathcal{X} . It is also easy to check that it holds if \mathcal{H} is the set of halfspaces in the feature space of most commonly used kernels (provided the feature space is of the same or higher dimension than the input space), such as polynomial and Gaussian kernels.

Based on $d_{\mathcal{H}}$ we can construct a distance between time-series probability distributions. For two time-series distributions ρ_1, ρ_2 we take the $d_{\mathcal{H}}$ between k -dimensional marginal distributions of ρ_1 and ρ_2 for each $k \in \mathbb{N}$, and sum them all up with decreasing weights.

Definition 1 (telescope distance $D_{\mathbf{H}}$). For two time series distributions ρ_1 and ρ_2 on the space $(\mathcal{X}^{\mathbb{N}}, \mathcal{F})$ and a sequence of sets of functions $\mathbf{H} = (\mathcal{H}_1, \mathcal{H}_2, \dots)$ define the telescope distance

$$D_{\mathbf{H}}(\rho_1, \rho_2) := \sum_{k=1}^{\infty} w_k \sup_{h \in \mathcal{H}_k} |\mathbf{E}_{\rho_1} h(X_1, \dots, X_k) - \mathbf{E}_{\rho_2} h(Y_1, \dots, Y_k)|, \quad (1)$$

where $w_k, k \in \mathbb{N}$ is a sequence of positive summable real weights (e.g., $w_k = 1/k^2$ or $w_k = 2^{-k}$).

Lemma 2. $D_{\mathbf{H}}$ is a metric if and only if $d_{\mathcal{H}_k}$ is a metric for every $k \in \mathbb{N}$.

Proof. The statement follows from the fact that two process distributions are the same if and only if all their finite-dimensional marginals coincide. \square

Definition 2 (empirical telescope distance \hat{D}). For a pair of samples $X_{1..n}$ and $Y_{1..m}$ define empirical telescope distance as

$$\hat{D}_{\mathbf{H}}(X_{1..n}, Y_{1..m}) := \sum_{k=1}^{\min\{m, n\}} w_k \sup_{h \in \mathcal{H}_k} \left| \frac{1}{n-k+1} \sum_{i=1}^{n-k+1} h(X_{i..i+k-1}) - \frac{1}{m-k+1} \sum_{i=1}^{m-k+1} h(Y_{i..i+k-1}) \right|. \quad (2)$$

All the methods presented in this work are based on the empirical telescope distance. The key fact is that it is an asymptotically consistent estimate of the telescope distance, that is, the latter can be consistently estimated based on sampling.

Theorem 1. Let $\mathbf{H} = (\mathcal{H}_k)_{k \in \mathbb{N}}$ be a sequence of separable sets \mathcal{H}_k of indicator functions (over \mathcal{X}^k) of finite VC dimension such that \mathcal{H}_k generates \mathcal{F}_k . Then, for every stationary ergodic time series distributions ρ_X and ρ_Y generating samples $X_{1..n}$ and $Y_{1..m}$ we have

$$\lim_{n, m \rightarrow \infty} \hat{D}_{\mathbf{H}}(X_{1..n}, Y_{1..m}) = D_{\mathbf{H}}(\rho_X, \rho_Y) \quad (3)$$

Note that $\hat{D}_{\mathbf{H}}$ is a biased estimate of $D_{\mathbf{H}}$, and, unlike in the i.i.d. case, the bias may depend on the distributions; however, the bias is $o(n)$.

Remark. The condition that the sets \mathcal{H}_k are sets of indicator function of finite VC dimension comes from [2], where it is shown that for any stationary ergodic distribution ρ , under these conditions, $\sup_{h \in \mathcal{H}_k} \frac{1}{n-k+1} \sum_{i=1}^{n-k+1} h(X_{i..i+k-1})$ is an asymptotically consistent estimate of $\sup_{h \in \mathcal{H}_k} \mathbf{E}_{\rho} h(X_1, \dots, X_k)$. This fact implies that $d_{\mathcal{H}_k}$ can be consistently estimated, from which the theorem is derived.

Proof of Theorem 1. As established in [2], under the conditions of the theorem we have

$$\lim_{n \rightarrow \infty} \sup_{h \in \mathcal{H}_k} \frac{1}{n-k+1} \sum_{i=1}^{n-k+1} h(X_{i..i+k-1}) = \sup_{h \in \mathcal{H}_k} \mathbf{E}_{\rho_X} h(X_1, \dots, X_k) \quad \rho_X\text{-a.s.} \quad (4)$$

for all $k \in \mathbb{N}$, and likewise for ρ_Y . Fix an $\varepsilon > 0$. We can find a $T \in \mathbb{N}$ such that

$$\sum_{k>T} w_k \leq \varepsilon. \quad (5)$$

Note that T depends only on ε . Moreover, as follows from (4), for each $k = 1..T$ we can find an N_k such that

$$\left| \sup_{h \in \mathcal{H}_k} \frac{1}{n-k+1} \sum_{i=1}^{n-k+1} h(X_{i..i+k-1}) - \sup_{h \in \mathcal{H}_k} \mathbf{E}_{\rho_X} h(X_{1..k}) \right| \leq \varepsilon/T \quad (6)$$

Let $N_k := \max_{i=1..T} N_i$ and define analogously M for ρ_Y . Thus, for $n \geq N, m \geq M$ we have

$$\begin{aligned} & \hat{D}_{\mathbf{H}}(X_{1..n}, Y_{1..m}) \\ & \leq \sum_{k=1}^T w_k \sup_{h \in \mathcal{H}_k} \left| \frac{1}{n-k+1} \sum_{i=1}^{n-k+1} h(X_{i..i+k-1}) - \frac{1}{m-k+1} \sum_{i=1}^{m-k+1} h(Y_{i..i+k-1}) \right| + \varepsilon \\ & \leq \sum_{k=1}^T w_k \sup_{h \in \mathcal{H}_k} \left\{ \left| \frac{1}{n-k+1} \sum_{i=1}^{n-k+1} h(X_{i..i+k-1}) - \mathbf{E}_{\rho_1} h(X_{1..k}) \right| \right. \\ & \quad \left. + |\mathbf{E}_{\rho_1} h(X_{1..k}) - \mathbf{E}_{\rho_2} h(Y_{1..k})| \right. \\ & \quad \left. + \left| \mathbf{E}_{\rho_2} h(Y_{1..k}) - \frac{1}{m-k+1} \sum_{i=1}^{m-k+1} h(Y_{i..i+k-1}) \right| \right\} + \varepsilon \\ & \leq 3\varepsilon + D_{\mathbf{H}}(\rho_X, \rho_Y), \end{aligned}$$

where the first inequality follows from the definition (2) of $\hat{D}_{\mathbf{H}}$ and from (5), and the last inequality follows from (6). Since ε was chosen arbitrary the statement follows. \square

4 Calculating $\hat{D}_{\mathbf{H}}$ using binary classification methods

The methods for solving various statistical problems that we suggest are all based on $\hat{D}_{\mathbf{H}}$. The main appeal of this approach is that $\hat{D}_{\mathbf{H}}$ can be calculated using binary classification methods. Here we explain how to do it.

The definition (2) of $D_{\mathbf{H}}$ involves calculating l summands (where $l := \min\{n, m\}$), that is

$$\sup_{h \in \mathcal{H}_k} \left| \frac{1}{n-k+1} \sum_{i=1}^{n-k+1} h(X_{i..i+k-1}) - \frac{1}{m-k+1} \sum_{i=1}^{m-k+1} h(Y_{i..i+k-1}) \right| \quad (7)$$

for each $k = 1..l$. Assuming that $h \in \mathcal{H}_k$ are indicator functions, calculating each of the summands amounts to solving the following k -dimensional binary classification problem. Consider $X_{i..i+k-1}$, $i = 1..n-k+1$ as class-1 examples and $Y_{i..i+k-1}$, $i = 1..m-k+1$ as class-0 examples. The supremum (7) is attained on $h \in \mathcal{H}_k$ that minimizes the empirical risk, with examples weighted with respect to the sample size. Indeed, we can define the weighted empirical risk of any $h \in \mathcal{H}_k$ as

$$\left| \frac{1}{n-k+1} \sum_{i=1}^{n-k+1} (1 - h(X_{i..i+k-1})) + \frac{1}{m-k+1} \sum_{i=1}^{m-k+1} h(Y_{i..i+k-1}) \right|,$$

which is obviously minimized by any $h \in \mathcal{H}_k$ that attains (7).

Thus, as long as we have a way to find $h \in \mathcal{H}_k$ that minimizes empirical risk, we have a consistent estimate of $D_{\mathcal{H}}(\rho_X, \rho_Y)$, under the mild conditions on \mathbf{H} required by Theorem 1. Since the dimension of the resulting classification problems grows with the length of the sequences, one should prefer methods that work in high dimensions, such as soft-margin SVMs [6].

A particularly remarkable feature is that *the choice of \mathcal{H}_k is much easier* for the problems that we consider than in the binary classification problem. Specifically, if (for some fixed k) the classifier that achieves the minimal (Bayes) error for the classification problem is not in \mathcal{H}_k , then obviously the error of an empirical risk minimizer will not tend to zero, no matter how much data we have. In contrast, all we need to achieve asymptotically 0 error in estimating \hat{D} (and therefore, in the learning problems considered below) is that the sets \mathcal{H}_k generate \mathcal{F}_k and have a finite VC dimension (for each k). This is the case already for the set of half-spaces in \mathbb{R}_k . In other words, the *approximation* error of the binary classification method (the classification error of the best f in \mathcal{H}_k) is not important. What is important is the estimation error; for asymptotic consistency results it has to go to 0 (hence the requirement on the VC dimension); for non-asymptotic results, it will appear in the error bounds, see Section 7. Thus, we have the following statement.

Claim 1. *The approximation error $|D_{\mathbf{H}}(P, Q) - \hat{D}_{\mathbf{H}}(X, Y)|$, and thus the error of the algorithms below, can be much smaller than the error of classification algorithms used to calculate $D_{\mathbf{H}}(X, Y)$.*

We can conclude that, beyond the requirement that \mathcal{H}_k generate \mathcal{F}_k for each $k \in \mathbb{N}$, the choice of H_k (or, say, of the kernel to use in SVM) is entirely up to the needs and constraints of specific applications.

Finally, we remark that while in the definition of the empirical distributional distance (2) the number of summands is l (the length of the shorter of the two samples), it can be replaced with any γ_l such that $\gamma_l \rightarrow \infty$, without affecting any asymptotic consistency results. In other words, Theorem 1, as well as all the consistency statements below, hold true for l replaced with any function γ_l that increases to infinity. A practically viable choice is $\gamma_l = \log l$; in fact, there is no reason to choose faster growing γ_n since the estimates for higher-order summands will not have enough data to converge. This is also the value we use in the experiments.

5 The three-sample problem

We start with a conceptually simple problem known in statistics as the three-sample problem (some times also called time-series classification). We are given three samples $X = (X_1, \dots, X_n)$, $Y = (Y_1, \dots, Y_m)$ and $Z = (Z_1, \dots, Z_l)$. It is known that X and Y were generated by different time-series distributions, whereas Z was generated by the same distribution as either X or Y . It is required to find out which one is the case. Both distributions are assumed to be stationary ergodic, but no further assumptions are made about them (no independence, mixing or memory assumptions). The three sample-problem for dependent time series has been addressed in [9] for Markov processes and in [23] for stationary ergodic time series. The latter work uses an approach based on the distributional distance.

Indeed, to solve this problem it suffices to have consistent estimates of some distance between time series distributions. Thus, we can use the telescope distance. The following statement is a simple corollary of Theorem 1.

Theorem 2. *Let the samples $X = (X_1, \dots, X_n)$, $Y = (Y_1, \dots, Y_m)$ and $Z = (Z_1, \dots, Z_l)$ be generated by stationary ergodic distributions ρ_X, ρ_Y and ρ_Z , with $\rho_X \neq \rho_Y$ and either (i) $\rho_Z = \rho_X$ or (ii) $\rho_Z = \rho_Y$. Let the sets \mathcal{H}_k , $k \in \mathbb{N}$ be separable sets of indicator functions over \mathcal{X}^k . Assume that each set \mathcal{H}_k , $k \in \mathbb{N}$ has a finite VC dimension and generates \mathcal{F}_k . A test that declares that (i) is true if $\hat{D}_{\mathbf{H}}(Z, X) \leq \hat{D}_{\mathbf{H}}(Z, Y)$ and that (ii) is true otherwise, makes only finitely many errors with probability 1 as $n, m, l \rightarrow \infty$.*

It is straightforward to extend this theorem to more than two classes; in other words, instead of X and Y one can have an arbitrary number of samples from different stationary ergodic distributions. A further generalization of this problem is the problem of time-series clustering, considered in the next section.

6 Clustering time series

We are given N time-series samples $X^1 = (X_1^1, \dots, X_{n_1}^1), \dots, X^N = (X_1^N, \dots, X_{n_N}^N)$, and it is required to cluster them into K groups, where, in different settings, K may be either known or unknown. While there may be many different approaches to define what should be considered a

good clustering, and, thus, what it means to have a consistent clustering algorithm, for the problem of clustering time-series samples there is a natural choice, proposed in [21]: Assume that each of the time-series samples $X^1 = (X_1^1, \dots, X_{n_1}^1), \dots, X^N = (X_1^N, \dots, X_{n_N}^N)$ was generated by one out of K different time-series distributions ρ_1, \dots, ρ_K . These distributions are unknown. The *target clustering* is defined according to whether the samples were generated by the same or different distributions: the samples belong to the same cluster if and only if they were generated by the same distribution. A clustering algorithm is called *asymptotically consistent* if with probability 1 from some n on it outputs the target clustering, where n is the length of the shortest sample $n := \min_{i=1..N} n_i \geq n'$.

Again, to solve this problem it is enough to have a metric between time-series distributions that can be consistently estimated. Our approach here is based on the telescope distance, and thus we use \hat{D} .

The clustering problem is relatively simple if the target clustering has what is called the *strict separation property* [4]: every two points in the same target cluster are closer to each other than to any point from a different target cluster. The following statement is an easy corollary of Theorem 1.

Theorem 3. *Let the sets $\mathcal{H}_k, k \in \mathbb{N}$ be separable sets of indicator functions over \mathcal{X}^k . Assume that each set $\mathcal{H}_k, k \in \mathbb{N}$ has a finite VC dimension and generates \mathcal{F}_k . If the distributions ρ_1, \dots, ρ_K generating the samples $X^1 = (X_1^1, \dots, X_{n_1}^1), \dots, X^N = (X_1^N, \dots, X_{n_N}^N)$ are stationary ergodic, then with probability 1 from some $n := \min_{i=1..N} n_i$ on the target clustering has the strict separation property with respect to $\hat{D}_{\mathbf{H}}$.*

With the strict separation property at hand, if the number of clusters K is known, it is easy to find asymptotically consistent algorithms. Here we give some simple examples, but the theorem below can be extended to many other distance-based clustering algorithms.

The *average linkage* algorithm works as follows. The distance between clusters is defined as the average distance between points in these clusters. First, put each point into a separate cluster. Then, merge the two closest clusters; repeat the last step until the total number of clusters is K . The *farthest point* clustering works as follows. Assign $c_1 := X^1$ to the first cluster. For $i = 2..K$, find the point $X^j, j \in \{1..N\}$ that maximizes the distance $\min_{t=1..i} \hat{D}_{\mathbf{H}}(X^j, c_t)$ (to the points already assigned to clusters) and assign $c_i := X^j$ to the cluster i . Then assign each of the remaining points to the nearest cluster. The following statement is a corollary of Theorem 3.

Theorem 4. *Under the conditions of Theorem 3, average linkage and farthest point clusterings are asymptotically consistent, provided the correct number of clusters K is given to the algorithm.*

Note that we do not require the samples to be independent; the joint distributions of the samples may be completely arbitrary, as long as the marginal distribution of each sample is stationary ergodic. These results can be extended to the online setting in the spirit of [13].

For the case of unknown number of clusters, the situation is different: one has to make stronger assumptions on the distributions generating the samples, since there is no algorithm that is consistent for all stationary ergodic distributions [22]; such stronger assumptions are considered in the next section.

7 Speed of convergence

The results established so far are asymptotic out of necessity: they are established under the assumption that the distributions involved are stationary ergodic, which is too general to allow for any meaningful finite-time performance guarantees. While it is interesting to be able to establish consistency results under such general assumptions, it is also interesting to see what results can be obtained under stronger assumptions. Moreover, since it is usually not known in advance whether the data at hand satisfies given assumptions or not, it appears important to have methods that have *both* asymptotic consistency in the general setting and finite-time performance guarantees under stronger assumptions. It turns out that this is possible: for the methods based on \hat{D} one can establish both the asymptotic performance guarantees for all stationary ergodic distributions and finite-sample performance guarantees under stronger assumptions, namely the uniform mixing conditions introduced below.

Another reason to consider stronger assumptions on the distributions generating the data is that some statistical problems, such as homogeneity testing or clustering when the number of clusters is unknown, are provably impossible to solve under the only assumption of stationary ergodic distributions, as shown in [22].

Thus, in this section we analyse the speed of convergence of \hat{D} under certain mixing conditions, and use it to construct solutions for the problems of homogeneity and clustering with an unknown number of clusters, as well as to establish finite-time performance guarantees for the methods presented in the previous sections.

A stationary distribution on the space of one-way infinite sequences $(\mathcal{X}^{\mathbb{N}}, \mathcal{F})$ can be uniquely extended to a stationary distribution on the space of two-way infinite sequences $(\mathcal{X}^{\mathbb{Z}}, \mathcal{F}_{\mathbb{Z}})$ of the form $\dots, X_{-1}, X_0, X_1, \dots$.

Definition 3 (β -mixing coefficients). *For a process distribution ρ define the mixing coefficients*

$$\beta(\rho, k) := \sup_{\substack{A \in \sigma(X_{-\infty..0}), \\ B \in \sigma(X_{k..\infty})}} |\rho(A \cap B) - \rho(A)\rho(B)|$$

where $\sigma(\dots)$ denotes the sigma-algebra of the random variables in brackets.

When $\beta(\rho, k) \rightarrow 0$ the process ρ is called uniformly β -mixing (with coefficients $\beta(\rho, k)$); this condition is much stronger than ergodicity, but is much weaker than the i.i.d. assumption.

7.1 Speed of convergence of \hat{D}

Assume that a sample $X_{1..n}$ is generated by a distribution ρ that is uniformly β -mixing with coefficients $\beta(\rho, k)$. Assume further that \mathcal{H}_k is a set of indicator functions with a finite VC dimension d_k , for each $k \in \mathbb{N}$.

Since in this section we are after finite-time bounds, we fix a concrete choice of the weights w_k in the definition 1 of \hat{D} ,

$$w_k := 2^{-k}. \quad (8)$$

The general tool that we use to obtain performance guarantees in this section is the following bound that can be obtained from the results of [12].

$$q_n(\rho, \mathcal{H}_k, \varepsilon) := \rho \left(\sup_{h \in \mathcal{H}_k} \left| \frac{1}{n-k+1} \sum_{i=1}^{n-k+1} h(X_{i..i+k-1}) - \mathbf{E}_{\rho} h(X_{1..k}) \right| > \varepsilon \right) \leq n\beta(\rho, t_n - k) + 8t_n^{d_k+1} e^{-l_n \varepsilon^2/8}, \quad (9)$$

where t_n are any integers in $1..n$ and $l_n = n/t_n$. The parameters t_n should be set according to the values of β in order to optimize the bound.

One can use similar bounds for classes of finite Pollard dimension [18] or more general bounds expressed in terms of covering numbers, such as those given in [12]. Here we consider classes of finite VC dimension only for the ease of the exposition and for the sake of continuity with the previous section (where it was necessary).

Furthermore, for the rest of this section we assume geometric β -mixing distributions, that is, $\beta(\rho, t) \leq \gamma^t$ for some $\gamma < 1$. Letting $l_n = t_n = \sqrt{n}$ the bound (9) becomes

$$q_n(\rho, \mathcal{H}_k, \varepsilon) \leq n\gamma^{\sqrt{n}-k} + 8n^{(d_k+1)/2} e^{-\sqrt{n}\varepsilon^2/8}. \quad (10)$$

Lemma 3. *Let two samples $X_{1..n}$ and $Y_{1..m}$ be generated by stationary distributions ρ_X and ρ_Y whose β -mixing coefficients satisfy $\beta(\rho, t) \leq \gamma^t$ for some $\gamma < 1$. Let \mathcal{H}_k , $k \in \mathbb{N}$ be some sets of indicator functions on \mathcal{X}^k whose VC dimension d_k is finite and non-decreasing with k . Then*

$$P(|\hat{D}_{\mathbf{H}}(X_{1..n}, Y_{1..m}) - D_{\mathbf{H}}(\rho_X, \rho_Y)| > \varepsilon) \leq 2\Delta(\varepsilon/4, n') \quad (11)$$

where $n' := \min\{n, m\}$, the probability is with respect to $\rho_X \times \rho_Y$ and

$$\Delta(\varepsilon, n) := -\log \varepsilon (n\gamma^{\sqrt{n}+\log(\varepsilon)} + 8n^{(d_{-\log \varepsilon}+1)/2} e^{-\sqrt{n}\varepsilon^2/8}). \quad (12)$$

Proof. From (8) we have $\sum_{k=-\log \varepsilon/2}^{\infty} w_k < \varepsilon/2$. Using this and the definitions (1) and (2) of $D_{\mathbf{H}}$ and $\hat{D}_{\mathbf{H}}$ we obtain

$$\begin{aligned} P(|\hat{D}_{\mathbf{H}}(X_{1..n_1}, Y_{1..n_2}) - D_{\mathbf{H}}(\rho_X, \rho_Y)| > \varepsilon) \\ \leq \sum_{k=1}^{-\log(\varepsilon/2)} (q_n(\rho_X, \mathcal{H}_k, \varepsilon/4) + q_n(\rho_Y, \mathcal{H}_k, \varepsilon/4)), \end{aligned}$$

which, together with (6) implies the statement. \square

7.2 Homogeneity testing

Given two samples $X_{1..n}$ and $Y_{1..m}$ generated by distributions ρ_X and ρ_Y respectively, the problem of homogeneity testing (or the two-sample problem) consists in deciding whether $\rho_X = \rho_Y$. A test is called (asymptotically) consistent if its probability of error goes to zero as $n' := \min\{m, n\}$ goes to infinity. As mentioned above, in general, for stationary ergodic time series distributions there is no asymptotically consistent test for homogeneity [22] (even for binary-valued time series); thus, stronger assumptions are in order.

Homogeneity testing is one of the classical problems of mathematical statistics, and one of the most studied ones. Vast literature exists on homogeneity testing for i.i.d. data, and for dependent processes as well. We do not attempt to survey this literature here. Our contribution to this line of research is to show that this problem can be reduced (via the telescope distance) to binary classification, in the case of strongly dependent processes satisfying some mixing conditions.

It is easy to see that under the mixing conditions of Lemma 1 a consistent test for homogeneity exists, and finite-sample performance guarantees can be obtained. It is enough to find a sequence $\varepsilon_n \rightarrow 0$ such that $\Delta(\varepsilon_n, n) \rightarrow 0$ (see (12)). Then the test can be constructed as follows: say that the two sequences $X_{1..n}$ and $Y_{1..m}$ were generated by the same distribution if $\hat{D}_{\mathbf{H}}(X_{1..n}, Y_{1..m}) < \varepsilon_{\min\{n, m\}}$; otherwise say that they were generated by different distributions.

Theorem 5. *Under the conditions of Lemma 3 the probability of Type I error (the distributions are the same but the test says they are different) of the described test is upper-bounded by $2\Delta(\varepsilon/4, n')$. The probability of Type II error (the distributions are different but the test says they are the same) is upper-bounded by $2\Delta((\delta - \varepsilon)/4, n')$ where $\delta := D_{\mathbf{H}}(\rho_X, \rho_Y)$.*

Proof. The statement is an immediate consequence of Lemma 3. Indeed, for the Type I error, the two sequences are generated by the same distribution, so the probability of error of the test is given by (11) with $D_{\mathbf{H}}(\rho_X, \rho_Y) = 0$. The probability of Type II error is given by $P(D_{\mathbf{H}}(\rho_X, \rho_Y) - \hat{D}_{\mathbf{H}}(X_{1..n_1}, Y_{1..n_2}) > \delta - \varepsilon)$, which is upper-bounded by $2\Delta((\delta - \varepsilon)/4, n')$ as follows from (11). \square

The optimal choice of ε_n may depend on the speed at which d_k (the VC dimension of \mathcal{H}_k) increases; however, for most natural cases (recall that \mathcal{H}_k are also parameters of the algorithm) this growth is polynomial, so the main term to control is $e^{-\sqrt{n}\varepsilon^2/8}$.

For example, if \mathcal{H}_k is the set of halfspaces in $\mathcal{X}^k = \mathbb{R}^k$ then $d_k = k + 1$ and one can chose $\varepsilon_n := n^{-1/8}$. The resulting probability of Type I error decreases as $\exp(-n^{1/4})$.

7.3 Clustering with a known or unknown number of clusters

If the distributions generating the samples satisfy certain mixing conditions, then we can augment Theorems 3 and 4 with finite-sample performance guarantees.

Theorem 6. *Let the distributions ρ_1, \dots, ρ_k generating the samples $X^1 = (X_1^1, \dots, X_{n_1}^1), \dots, X^N = (X_1^N, \dots, X_{n_N}^N)$ satisfy the conditions of Lemma 3. Define $\delta := \min_{i,j=1..N, i \neq j} D_{\mathbf{H}}(\rho_i, \rho_j)$ and $n := \min_{i=1..N} n_i$. Then with probability at least*

$$1 - N(N-1)\Delta(\delta/12, n')$$

the target clustering of the samples has the strict separation property. In this case single linkage and farthest point algorithms output the target clustering.

Proof. Note that a sufficient condition for the strict separation property to hold is that for every pair i, j of samples generated by the same distribution we have $\hat{D}_{\mathbf{H}}(X^i, X^j) \leq \delta/3$, and for every pair i, j of samples generated by different distributions we have $\hat{D}_{\mathbf{H}}(X^i, X^j) \geq 2\delta/3$. Using Lemma 3, the probability of such an even (for each pair) is upper-bounded by $2\Delta(\delta/12, n')$, which, multiplied by the total number $N(N-1)/2$ of pairs gives the statement. The second statement is obvious. \square

As with homogeneity testing, while in the general case of stationary ergodic distributions it is impossible to have a consistent clustering algorithm when the number of clusters k is unknown, the situation changes if the distributions satisfy certain mixing conditions. In this case a consistent clustering algorithm can be obtained as follows. Assign to the same cluster all samples that are at most ε_n -far from each other, where the threshold ε_n is selected the same way as for homogeneity testing: $\varepsilon_n \rightarrow 0$ and $\Delta(\varepsilon_n, n) \rightarrow 0$. The optimal choice of this parameter depends on the choice of \mathcal{H}_k through the speed of growth of the VC dimension d_k of these sets.

Theorem 7. *Given N samples generated by k different stationary distributions $\rho_i, i = 1..k$ (unknown k) all satisfying the conditions of Lemma 3, the probability of error (misclustering at least one sample) of the described algorithm is upper-bounded by*

$$N(N-1) \max\{\Delta(\varepsilon/4, n'), \Delta((\delta - \varepsilon)/4, n')\}$$

where $\delta := \min_{i,j=1..k, i \neq j} D_{\mathbf{H}}(\rho_i, \rho_j)$ and $n = \min_{i=1..N} n_i$, with $n_i, i = 1..N$ being lengths of the samples.

Proof. The statement follows from Theorem 5. \square

8 Experiments

For experimental evaluation we chose the problem of time-series clustering. Average-linkage clustering is used, with the telescope distance between samples calculated using an SVM, as described in Section 4. In all experiments, SVM is used with radial basis kernel, with default parameters of libsvm [5]. The parameters w_k in the definition of the telescope distance (Definition 1) are set to $w_k := k^{-2}$.

8.1 Synthetic data

For the artificial setting we have chosen highly-dependent time series distributions which have the same single-dimensional marginals and which cannot be well approximated by finite- or countable-state models. The distributions $\rho(\alpha), \alpha \in (0, 1)$, are constructed as follows. Select $r_0 \in [0, 1]$ uniformly at random; then, for each $i = 1..n$ obtain r_i by shifting r_{i-1} by α to the right, and removing the integer part. The time series (X_1, X_2, \dots) is then obtained from r_i by drawing a point from a distribution law \mathcal{N}_1 if $r_i < 0.5$ and from \mathcal{N}_2 otherwise. \mathcal{N}_1 is a 3-dimensional Gaussian with mean of 0 and covariance matrix $\text{Id} \times 1/4$. \mathcal{N}_2 is the same but with mean 1. If α is irrational¹ then the distribution $\rho(\alpha)$ is stationary ergodic, but does not belong to any simpler natural distribution family [25]. The single-dimensional marginal is the same for all values of α . The latter two properties make all parametric and most non-parametric methods inapplicable to this problem.

In our experiments, we use two process distributions $\rho(\alpha_i), i \in \{1, 2\}$, with $\alpha_1 = 0.31\dots, \alpha_2 = 0.35\dots$. The dependence of error rate on the length of time series is shown on Figure 1. One clustering experiment on sequences of length 1000 takes about 5 min. on a standard laptop.

8.2 Real data

To demonstrate the applicability of the proposed methods to realistic scenarios, we chose the brain-computer interface data from BCI competition III [17]. The dataset consists of (pre-processed) BCI recordings of mental imagery: a person is thinking about one of three subjects (left foot, right foot, a random letter). Originally, each time series consisted of several consecutive sequences of different classes, and the problem was supervised: three time series for training and one for testing.

¹in experiments simulated by a `longdouble` with a long mantissa

We split each of the original time series into classes, and then used our clustering algorithm in a completely unsupervised setting. The original problem is 96-dimensional, but we used only the first 3 dimensions (using all 96 gives worse performance). The typical sequence length is 300. The performance is reported in Table 1, labeled TS_{SVM} . All the computation for this experiment takes approximately 6 minutes on a standard laptop.

The following methods were used for comparison. First, we used dynamic time wrapping (DTW) [24] which is a popular base-line approach for time-series clustering. The other two methods in Table 1 are from [10]. The comparison is not fully relevant, since the results in [10] are for different settings; the method KCpA was used in change-point estimation method (a different but also unsupervised setting), and SVM was used in a supervised setting. The latter is of particular interest since the classification method we used in the telescope distance is also SVM, but our setting is unsupervised (clustering).

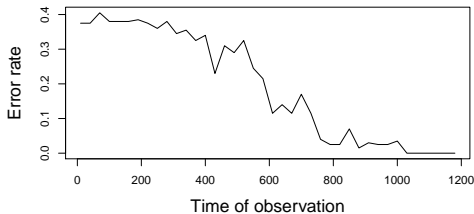


Figure 1: Error of two-class clustering using TS_{SVM} ; 10 time series in each target cluster, averaged over 20 runs.

	s_1	s_2	s_3
TS_{SVM}	84%	81%	61%
DTW	46%	41%	36%
KCpA	79%	74%	61%
SVM	76%	69%	60%

Table 1: Clustering accuracy in the BCI dataset. 3 subjects (columns), 4 methods (rows). Our method is TS_{SVM} .

Acknowledgments. This research was funded by the Ministry of Higher Education and Research, Nord-Pas-de-Calais Regional Council and FEDER (Contrat de Projets Etat Region CPER 2007-2013), ANR projects EXPLO-RA (ANR-08-COSI-004), Lampada (ANR-09-EMER-007) and CoAdapt, and by the European Community’s FP7 Program under grant agreements n° 216886 (PASCAL2) and n° 270327 (CompLACS).

References

- [1] Terrence M. Adams and Andrew B. Nobel. Uniform convergence of Vapnik-Chervonenkis classes under ergodic sampling. *The Annals of Probability*, 38:1345–1367, 2010.
- [2] Terrence M. Adams and Andrew B. Nobel. Uniform approximation of Vapnik-Chervonenkis classes. *Bernoulli*, 18(4):1310–1319, 2012.
- [3] Maria-Florina Balcan, Nikhil Bansal, Alina Beygelzimer, Don Coppersmith, John Langford, and Gregory Sorkin. Robust reductions from ranking to classification. In Nader Bshouty and Claudio Gentile, editors, *Learning Theory*, volume 4539 of *Lecture Notes in Computer Science*, pages 604–619. 2007.
- [4] M.F. Balcan, A. Blum, and S. Vempala. A discriminative framework for clustering via similarity functions. In *Proceedings of the 40th annual ACM symposium on Theory of computing*, pages 671–680. ACM, 2008.
- [5] Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [6] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Mach. Learn.*, 20(3):273–297, 1995.
- [7] R. Fortet and E. Mourier. Convergence de la répartition empirique vers la répartition théorique. *Ann. Sci. Ec. Norm. Super., III. Ser.*, 70(3):267–285, 1953.
- [8] R. Gray. *Probability, Random Processes, and Ergodic Properties*. Springer Verlag, 1988.
- [9] M. Gutman. Asymptotically optimal classification for multiple tests with empirically observed statistics. *IEEE Transactions on Information Theory*, 35(2):402–408, 1989.

- [10] Zaïd Harchaoui, Francis Bach, and Eric Moulines. Kernel change-point analysis. In *NIPS*, pages 609–616, 2008.
- [11] L. V. Kantorovich and G. S. Rubinstein. On a function space in certain extremal problems. *Dokl. Akad. Nauk USSR*, 115(6):1058–1061, 1957.
- [12] R.L. Karandikar and M. Vidyasagar. Rates of uniform convergence of empirical means with mixing processes. *Statistics and Probability Letters*, 58:297–307, 2002.
- [13] A. Khaleghi, D. Ryabko, J. Mary, and P. Preux. Online clustering of processes. In *AISTATS*, JMLR W&CP 22, pages 601–609, 2012.
- [14] Daniel Kifer, Shai Ben-David, and Johannes Gehrke. Detecting change in data streams. In *Proceedings of the Thirtieth international conference on Very large data bases - Volume 30*, VLDB’04, pages 180–191, 2004.
- [15] A.N. Kolmogorov. Sulla determinazione empirica di una legge di distribuzione. *G. Inst. Ital. Attuari*, pages 83–91, 1933.
- [16] John Langford, Roberto Oliveira, and Bianca Zadrozny. Predicting conditional quantiles via reduction to classification. In *UAI*, 2006.
- [17] José del R. Millán. On the need for on-line learning in brain-computer interfaces. In *Proc. of the Int. Joint Conf. on Neural Networks*, 2004.
- [18] D. Pollard. *Convergence of Stochastic Processes*. Springer, 1984.
- [19] B. Ryabko. Prediction of random sequences and universal coding. *Problems of Information Transmission*, 24:87–96, 1988.
- [20] B. Ryabko. Compression-based methods for nonparametric prediction and estimation of some characteristics of time series. *IEEE Transactions on Information Theory*, 55:4309–4315, 2009.
- [21] D. Ryabko. Clustering processes. In *Proc. the 27th International Conference on Machine Learning (ICML 2010)*, pages 919–926, Haifa, Israel, 2010.
- [22] D. Ryabko. Discrimination between B-processes is impossible. *Journal of Theoretical Probability*, 23(2):565–575, 2010.
- [23] D. Ryabko and B. Ryabko. Nonparametric statistical inference for ergodic processes. *IEEE Transactions on Information Theory*, 56(3):1430–1435, 2010.
- [24] H. Sakoe and S. Chiba. Dynamic programming algorithm optimization for spoken word recognition. *IEEE Transactions on Acoustics, Speech and Signal Processing*, 26(1):43–49, 1978.
- [25] P. Shields. *The Ergodic Theory of Discrete Sample Paths*. AMS Bookstore, 1996.
- [26] V. M. Zolotarev. Metric distances in spaces of random variables and their distributions. *Math. USSR-Sb*, 30(3):373–401, 1976.