

# Linear Distances between Markov Chains\*

Przemysław Daca<sup>1</sup>, Thomas A. Henzinger<sup>2</sup>, Jan Křetínský<sup>3</sup>, and Tatjana Petrov<sup>4</sup>

1 IST Austria, Klosterneuburg, Austria

2 IST Austria, Klosterneuburg, Austria

3 Institut für Informatik, Technische Universität München, Germany

4 IST Austria, Klosterneuburg, Austria

---

## Abstract

We introduce a general class of distances (metrics) between Markov chains, which are based on linear behaviour. This class encompasses distances given topologically (such as the total variation distance or trace distance) as well as by temporal logics or automata. We investigate which of the distances can be approximated by observing the systems, i.e. by black-box testing or simulation, and we provide both negative and positive results.

**1998 ACM Subject Classification** G.3 Probability and Statistics, F.4.3 Formal Languages

**Keywords and phrases** probabilistic systems, verification, statistical model checking, temporal logic, behavioural equivalence

**Digital Object Identifier** 10.4230/LIPIcs.CONCUR.2016.20

## 1 Introduction

Behaviour of processes is traditionally compared using various notions of equivalence, such as trace equivalence, bisimulation, etc. However, the concept of equivalence is often too coarse for quantitative systems, such as Markov chains. For instance, probabilities of failures of particular hardware components are typically only empirically estimated and the slightest imprecision in the estimate may result in breaking the equivalence between processes. Moreover, if the (possibly black-box) processes are indeed different we would like to measure *how much they differ*. This has led to lifting the Boolean idea of behavioural equivalence to a finer, quantitative notion of behavioural *distance* between processes. The distance between processes  $s$  and  $t$  is typically formalized as  $\sup_{p \in \mathcal{C}} |p(s) - p(t)|$  where  $\mathcal{C}$  is a class of properties of interest and  $p(s)$  is a quantitative value of the property  $p$  in process  $s$  [13]. This notion has been introduced in [13] for Markov chains and further developed in various settings, such as Markov decision processes [16], quantitative transition systems [12], or concurrent games [11].

Several kinds of distances have been investigated for Markov chains. On the one hand, branching distances, e.g. [1, 13, 26, 25, 4, 3, 2, 18], lift the equivalence given by the probabilistic bisimulation of Larsen and Skou [22]. On the other hand, there are *linear* distances, in particular the total variation distance [8, 6] and trace distances [20, 5]. Linear distances are particularly appropriate when (i) we are interested in linear-time properties, and (ii) we want to estimate the distance based only on simulation runs from the initial distribution of

---

\* This research was funded in part by the European Research Council (ERC) under grant agreement 267989 (QUAREM), the Austrian Science Fund (FWF) under grants project S11402-N23 (RiSE and SHiNE) and Z211-N23 (Wittgenstein Award), by the Czech Science Foundation Grant No. P202/12/G061, and by the SNSF Advanced Postdoc. Mobility Fellowship – grant number P300P2\_161067.



the system, i.e. in a black-box setting. (Recall that for branching distances, the underlying probabilistic bisimulation corresponds to testing equivalence where not only runs from the initial distribution can be observed, but it is also possible to dump the current state of the system, and later restart the simulation from this state [22].)

In this paper, we introduce a simple framework for linear distances between Markov chains, using the formula above, where  $p(s)$  is the probability of satisfying  $p$  when starting a simulation run in state  $s$  (when  $p$  is seen as a language of  $\omega$ -words it is the probability to generate a trace belonging to  $p$ ). We consider several classes  $\mathcal{C}$  of languages of interest, characterized from several points of view, e.g. topologically, by linear-time logics, or by automata, thus rendering our framework versatile.

We investigate when a given distance can be estimated in a black-box setting, i.e. only from simulations. One of the main difficulties is that the class  $\mathcal{C}$  typically includes properties with arbitrarily long horizon or even infinite-horizon properties, whereas every simulation run is necessarily finite. Note that we do not employ any simplifications such as imposed fixed horizon or discounting, typically used for obtaining efficient algorithms, e.g., [13, 26, 3], and the undiscounted setting is fundamentally more complex [25]. Since even simpler tasks are impossible for unbounded horizon in the black-box setting without any further knowledge, we assume we only know a lower bound on the minimum transition probability  $p_{\min}$ . Note that knowledge of  $p_{\min}$  has been justified in [10].

Our contribution is the following:

- We introduce a systematic linear-distance framework and illustrate it with several examples, including distances previously investigated in the literature.
- The main technical contributions are (i) a negative result stating that the total variation distance cannot be estimated by simulating the systems, and (ii) a positive result that the trace distance can be estimated.
- These results are further exploited to provide both negative and positive results for each of the settings where the language class is given topologically, by LTL (linear temporal logic) fragments, and by automata. We also show that the negative result on the total variation distance can be turned into a positive result if the transition probabilities have finite precision.

## 1.1 Related work

There are two main linear distances considered for Markov chains: the total variation distance and trace distance. Several algorithms have been proposed for both of them in the case when the Markov chains are known (white-box setting). We are not aware of any work where the distances are estimated only from simulating the systems (black-box setting).

Firstly, for the *total variation distance* in the white-box setting, [8] shows that deciding whether it equals one can be done in polynomial time, but computing it is NP-hard and not known to be decidable, however, it can be approximated; [6] considers this distance more generally for semi-Markov processes, provides a different approximation algorithm, and shows it coincides with distances based on (i) metric temporal logic, and (ii) timed automata languages.

Secondly, the *trace distance* is based on the notion of trace equivalence, which can be decided in polynomial time [15] (however, trace refinement of Markov decision processes is already undecidable [17]). Several variants of trace distance are considered in [20] where it is taken as a limit of finite-trace distances, possibly using discounting or averaging. In [5] the finite-trace distance is shown to coincide with distances based on (i) LTL, and (ii)

LTL without the U operator, i.e., only using the X operator and Boolean connectives. This distances is also shown to be NP-hard and not known to be decidable, similarly to the total variation distance. Finally, an approximation algorithm is shown (again in the white-box setting), where the over-approximates are branching-time distances, showing an interesting connection between the branching and linear distances.

In [21] the distinguishability problem is considered, i.e. given two Markov chains whether there is a monitor that reads a single sample and with high probability decides which chain produced the sequence. This is indeed possible when the total variation distance between the chains equals one, and [21] shows how to construct such monitors. In contrast, our negative results shows that it is not possible to decide with high probability whether the total variation distance equals one when the two Markov are black-box.

Linear distances have been proposed also for quantitative transition systems, e.g. [11]. Moreover, there are other useful distances based on different fundamentals; for instance, the Skorokhod distance [7, 23, 14] measures the discrete differences between systems while allowing for timing distortion; Kullback-Leibler divergence [20] is useful from the information-theoretic point of view. Finally, distances have been also studied with respect to applications in linear-time model checking [24, 5].

## 1.2 Outline

After recalling the basic notions in Section 2, we introduce our framework and illustrate it with examples in Section 3. We define our problem formally in Section 4. In Sections 5 and 6 we provide the proofs of our technically principal negative and positive result, respectively. Section 7 extends the results in the settings of topology, logics and automata, and discusses general conditions for estimability. We conclude in Section 8.

Proofs omitted due to space constraints can be found in [9].

## 2 Preliminaries

We consider a finite set  $Ap$  of atomic propositions and denote  $\Sigma = 2^{Ap}$ .

► **Definition 1** (Markov chain). A (labelled) Markov chain (MC) is a tuple  $\mathcal{M} = (S, \mathbf{P}, \mu, L)$ , where

- $S$  is a finite set of states,
- $\mathbf{P} : S \times S \rightarrow [0, 1]$  is a transition probability matrix, such that for every  $s \in S$  it holds  $\sum_{s' \in S} \mathbf{P}(s, s') = 1$ ,
- $\mu$  is an initial probability distribution over  $S$ ,
- $L : S \rightarrow \Sigma$  is a labelling function.

A run of  $\mathcal{M}$  is an infinite sequence  $\rho = s_1 s_2 \dots$  of states, such that  $\mu(s_1) > 0$  and  $\mathbf{P}(s_i, s_{i+1}) > 0$  for all  $i \geq 1$ ; we let  $\rho[i]$  denote the state  $s_i$ . A path in  $\mathcal{M}$  is a finite prefix of a run of  $\mathcal{M}$ . An  $\omega$ -word is an infinite sequence  $a_1 a_2 \dots \in \Sigma^\omega$  of symbols from  $\Sigma$ ; a word is a finite prefix  $w \in \Sigma^*$  of an  $\omega$ -word. We extend the labelling notation so that for a path  $\pi \in S^k$ , the projected sequence  $L(\pi)$  is the word  $w \in \Sigma^k$ , where  $w[i] = L(\pi[i])$ , and the inverse map is  $L^{-1}(w) = \{\pi \in S^k \mid L(\pi) = w\}$ . Given a path  $\pi = s_1 \dots s_n$ , we denote the  $k$ -prefix of  $\pi$  by  $\pi \downarrow k = s_1 \dots s_k$ , and similarly for prefixes of words.

Each path  $\pi$  in  $\mathcal{M}$  determines the set of runs  $\text{Cone}(\pi)$  consisting of all runs that start with  $\pi$ . To  $\mathcal{M}$  we assign the probability space  $(\text{Runs}, \mathcal{F}, \mathbb{P}_{\mathcal{M}})$ , where  $\text{Runs}$  is the set of all runs in  $\mathcal{M}$ ,  $\mathcal{F}$  is the  $\sigma$ -algebra generated by all  $\text{Cone}(\pi)$ , and  $\mathbb{P}_{\mathcal{M}}$  is the unique probability measure such that  $\mathbb{P}_{\mathcal{M}}(\text{Cone}(s_1 \dots s_n)) = \mu(s_1) \cdot \prod_{i=1}^{n-1} \mathbf{P}(s_i, s_{i+1})$ , where the empty product

equals 1. We will omit the subscript in  $\mathbb{P}_{\mathcal{M}}$  if the Markov chain is clear from the context. Further, we write  $\mathbb{P}_{\mathcal{M}}^s$  for the probability measure, where  $\mu(s) = 1$  and  $\mu(s') = 0$  for  $s' \neq s$ . Finally, we overload the notation and for a path  $\pi$  write  $\mathbb{P}(\pi)$  meaning  $\mathbb{P}(\text{Cone}(\pi))$ , and for a  $(\omega)$ -word  $w$ , we write  $\mathbb{P}(w)$  meaning  $\mathbb{P}(L^{-1}(w))$ .

### 3 Framework for Linear Distances

In this section we introduce our framework for linear distances. For  $i \in \{1, 2\}$ , let  $\mathcal{M}_i = (S, \mathbf{P}_i, \mu_i, L)$  denote a Markov chain<sup>1</sup> and  $(\text{Runs}, \mathcal{F}, \mathbb{P}_i)$  the induced probability space. Since single runs of Markov chains typically have measure 0, we introduce linear distances using measurable sets of runs:

► **Definition 2** ( $\mathcal{L}$ -distance). For a class  $\mathcal{L} \subseteq \mathcal{F}$  of measurable  $\omega$ -languages<sup>2</sup>, the  $\mathcal{L}$ -distance  $D_{\mathcal{L}}$  is defined by

$$D_{\mathcal{L}}(\mathcal{M}_1, \mathcal{M}_2) = \sup_{X \in \mathcal{L}} |\mathbb{P}_1(X) - \mathbb{P}_2(X)|.$$

Note that every  $D_{\mathcal{L}}$  is a pseudo-metric<sup>3</sup>. However, two different MCs can have distance 0, for instance, when they induce the same probability space.

The definition of  $\mathcal{L}$ -distances can be instantiated either (i) by a direct topological description of  $\mathcal{L}$ , or indirectly (ii) by a class  $\mathcal{A}$  of automata inducing the class of recognized languages  $\mathcal{L} = \{L(A) \mid A \in \mathcal{A}\}$ , or (iii) by a set of formulae  $\mathfrak{L}$  of a linear-time logic inducing the languages of models  $\mathcal{L} = \{L(\varphi) \mid \varphi \in \mathfrak{L}\}$  where  $L(\varphi)$  denotes the language of  $\omega$ -words satisfying the formula  $\varphi$ .

We now discuss several particularly interesting instantiations:

► **Example 3** (Total variation). One extreme choice is to consider all measurable languages, resulting in the *total variation distance*  $D_{\text{TV}}(\mathcal{M}_1, \mathcal{M}_2) = \sup_{X \in \mathcal{F}(\Sigma)} |\mathbb{P}_1(X) - \mathbb{P}_2(X)|$ .

► **Example 4** (Trace distances). The other extreme choices are to consider (i) only the generators of  $\mathcal{F}(\Sigma)$ , i.e. the cones  $\{w\Sigma^\omega \mid w \in \Sigma^*\}$ , resulting in the *finite-trace distance*  $D_{\text{FT}}(\mathcal{M}_1, \mathcal{M}_2) = \sup_{w \in \Sigma^+} |\mathbb{P}_1(w) - \mathbb{P}_2(w)|$ ; or (ii) only the elementary events, i.e.  $\Sigma^\omega$ , resulting in the *infinite-trace distance*  $D_{\text{IT}}(\mathcal{M}_1, \mathcal{M}_2) = \sup_{w \in \Sigma^\omega} |\mathbb{P}_1(w) - \mathbb{P}_2(w)|$ .

► **Example 5** (Topological distances). There are many possible choices for  $\mathcal{L}$  between the two extremes above, such as *clopen sets*  $\Delta_1$ , which are finite unions of cones (being both closed and open), *open sets*  $\Sigma_1$ , which are infinite unions of cones, *closed sets*  $\Pi_1$ , or classes higher in the *Borel hierarchy* such as the class of  $\omega$ -regular languages (within  $\Delta_3$ ), or languages given by thresholds for a *long-run average reward* (within  $\Sigma_3$ ).

► **Example 6** (Automata distances). The class of  $\omega$ -regular languages can also be given in terms of automata, for instance by the class of all deterministic *Rabin automata* (DRA). Similarly, the closed sets  $\Pi_1$  correspond to the class of deterministic Büchi automata with all states final. Further, we can restrict the class of all DRA to those of *size at most k* for a fixed  $k \in \mathbb{N}$ , denoting the resulting distance by  $D_{\text{DRA} \leq k}$ .

<sup>1</sup> To avoid clutter, the chains are defined over the same state space with the same labelling, which can be w.l.o.g. achieved by their disjoint union.

<sup>2</sup> Formally, the measurable space of  $\omega$ -languages is given by the set  $\Sigma^\omega$  equipped with a  $\sigma$ -algebra  $\mathcal{F}(\Sigma)$  generated by the set of cones  $\{w\Sigma^\omega \mid w \in \Sigma^*\}$ . This ensures, for every measurable  $\omega$ -language  $X$ , that  $L^{-1}(X)$  is measurable in every MC.

<sup>3</sup> It is symmetric, it satisfies the triangle inequality, and the distance between identical MCs is 0.

► **Example 7** (Logical distances). The class of  $\omega$ -regular languages can also be given in terms of logic, by the monadic second-order logic (with order). Further useful choices include the *first-order logic with order*, corresponding to the star-free languages and to the *linear temporal logic* (LTL), or its fragments such as LTL with only **X** or only **F** and **G** operators etc.

► **Remark.** The introduced distances can also be considered in the discrete setting, resulting in various notions of equivalence. For instance, the *finite-trace equivalence*  $E_{FT}$  can be derived from the finite-trace distance by the following discretization:

$$E_{FT}(\mathcal{M}_1, \mathcal{M}_2) = \begin{cases} 0 & \text{if } D_{FT}(\mathcal{M}_1, \mathcal{M}_2) = 0 \\ 1 & \text{otherwise, i.e., } D_{FT}(\mathcal{M}_1, \mathcal{M}_2) > 0. \end{cases}$$

## 4 Problem Statement

Linear distances can be very useful when we want to compare a black-box system with another system, e.g. a white-box specification or a black-box previous version of the system. Indeed, in such a setting we can typically obtain simulation runs of the system and we must establish a relation between the systems based on these runs only. This is in contrast with branching distances where either both systems are assumed white-box or there are strong requirements on the testing abilities, such as dumping the current state of the system, arbitrary many restarts from there, and nesting this branching arbitrarily. Therefore, we focus on the setting where we can obtain only finite prefixes of runs and we use statistics to (i) deduce information on the whole infinite runs, and (ii) estimate the distance of the systems.

For a given distance function  $D_{\mathcal{L}}$ , the goal is to construct an almost-surely terminating algorithm that given

- any desired finite number of sampled simulation run from Markov chains  $\mathcal{M}_1$  and  $\mathcal{M}_2$  of any desired finite length,
- lower bound  $p_{\min} > 0$  on the minimum (non-zero) transition probability,
- confidence  $\alpha \in (0, 1)$ ,
- interval width  $\delta \in (0, 1)$ ,

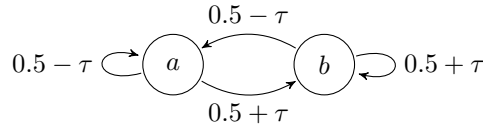
computes an interval  $I$  such that  $|I| \leq \delta$  and  $\Pr[D_{\mathcal{L}}(\mathcal{M}_1, \mathcal{M}_2) \in I] \geq 1 - \alpha$ .

A distance function is called *estimable*, if there exists an algorithm in the above sense, and *inestimable* otherwise.

## 5 Inestimability: Total variation distance

We show that for the total variation distance  $D_{TV}$  there exists no “statistical” algorithm (in the above sense) which is correct for all inputs  $(\mathcal{M}_1, \mathcal{M}_2, \alpha, \delta)$ . Our argument consists of two steps:

1. We construct two chains such that  $D_{TV}(\mathcal{M}_1, \mathcal{M}_2) = 1$ , namely the two MCs shown in Figure 1 (similar to [20]): one with  $\tau = 0$  and the other with small  $\tau > 0$ .
2. We show that any potentially correct algorithm will give with high probability an incorrect output for some choice of  $\tau, \alpha, \delta$ .



■ **Figure 1** A Markov chain with labelling displayed in states.

**Maximizing event.** We start by showing that even an arbitrarily small difference in transition probabilities between two Markov chains may result in total variation distance of 1. Consider two Markov chains as in Figure 1, where  $\mathcal{M}_1$  has  $\tau = 0$ , and  $\mathcal{M}_2$  has  $\tau > 0$ . We assume that the initial distribution for each chain is its stationary distribution. In this setting, every simulation step is like an independent trial with probability  $0.5 - \tau$  (resp.  $0.5 + \tau$ ) of seeing  $a$  (resp.  $b$ ).

Let  $X_n$  (resp.  $Y_n$ ) denote the number of  $b$  symbols in a random path of length  $n$  sampled from  $\mathcal{M}_1$  (resp.  $\mathcal{M}_2$ ). By the central limit theorem the distributions of  $X_n$  and  $Y_n$  are converging to the normal distribution when  $n \rightarrow \infty$ :

$$X_n \approx \mathcal{N}(0.5n, 0.5^2n) \quad Y_n \approx \mathcal{N}((0.5 + \tau)n, n(0.25 - \tau^2)).$$

For  $n \in \mathbb{N}$  let the event  $E_n$  mean “there is at most  $c_n = (0.5 + \tau/2)n$  symbols  $b$  in the path prefix of length  $n$ .” The probabilities of event  $E_n$  in the two Markov chains are:

$$\mathbb{P}_{\mathcal{M}_1}(E_n) = \mathbb{P}_{\mathcal{M}_1}(X_n \leq c_n) = \Phi(\tau\sqrt{n}) \quad \mathbb{P}_{\mathcal{M}_2}(E_n) = \mathbb{P}_{\mathcal{M}_2}(Y_n \leq c_n) = \Phi\left(\frac{-0.5\tau\sqrt{n}}{\sqrt{0.25 - \tau^2}}\right),$$

where  $\Phi$  is the CDF of the standard normal distribution. For  $n \rightarrow \infty$  the probability of  $E_n$  in  $\mathcal{M}_1$  and  $\mathcal{M}_2$  converges to 1 and 0, respectively, so the total variation distance converges to 1.

**Negative result for total variation distance.** Now we show that there is no statistical procedure for estimating total variation distance that would almost-surely terminate.

► **Theorem 8.** *For any  $\delta < 1$  and  $\alpha < \frac{1}{2}$ , there is no algorithm for computing a  $1 - \alpha$  confidence interval of size  $\delta$  for the total variation distance that almost-surely terminates.*

**Proof.** Let us write  $\mathcal{M}(\tau)$  for a Markov chain in Figure 1 with the parameter  $\tau$  and the initial distribution being stationary.

For  $\alpha < \frac{1}{2}$  we define the following decision problem  $B_\alpha$ :

- The input to  $B_\alpha$  is a single path from  $\mathcal{M}(\tau)$  of arbitrary length, where  $\tau$  is unknown,
- The task of  $B_\alpha$  is to output answer **Yes** with probability  $\geq 1 - \alpha$  if  $D_{TV}(\mathcal{M}(0), \mathcal{M}(\tau)) = 1$ , output answer **No** with probability  $\geq 1 - \alpha$  if  $D_{TV}(\mathcal{M}(0), \mathcal{M}(\tau)) = 0$ . Note that  $D_{TV}(\mathcal{M}(0), \mathcal{M}(\tau))$  can equal only 0 or 1.

The remaining part of proof is done in two parts. In the first part, we show that there is no algorithm that solves  $B_\alpha$  and almost-surely terminates. In the second part we reduce the problem  $B_\alpha$  to computing a confidence interval for the total variation distance.

**Part I.** Suppose the opposite of the claim: that for some  $\alpha < \frac{1}{2}$  there is an algorithm which solves  $B_\alpha$  and almost-surely terminates. We represent the algorithm for solving  $B_\alpha$  as a deterministic Turing machine TM, which works as follows:

1. The input tape of TM contains a (single) randomly sampled run of  $\mathcal{M}(\tau)$ ,
2. TM reads a part of the run from the tape and eventually returns **Yes/No** answer.

The input to the TM is random, therefore we can assign a probability distribution to the computations of TM. To this end, we represent the answer of TM by the random variable  $X : \text{Runs} \mapsto \{\mathbf{Yes}, \mathbf{No}\}$ , and we use the random variable  $Y : \text{Runs} \mapsto \mathbb{N} \cup \{\infty\}$  to represent the number of path symbols TM reads before terminating, where  $\infty$  means that TM does not terminate.

Suppose we run TM on the Markov chain  $\mathcal{M}(0)$ . We write  $\mathbb{P}_1$  for the probability measure of TM on this input. The total variation distance between the two Markov chains  $\mathcal{M}(0)$  is 0, so with probability  $\geq 1 - \alpha$  TM returns answer **No**, i.e.  $\mathbb{P}_1(X = \mathbf{No}) \geq 1 - \alpha$ .

By assumption TM almost-surely terminates on every input, so  $\mathbb{P}_1(Y \in \mathbb{N}) = 1$ . Let  $q$  be the following quantile:

$$q = \min\{c \in \mathbb{N} : \mathbb{P}_1(Y \leq c) \geq 0.5 + \alpha\}.$$

► **Claim.**  $q \in \mathbb{N}$ .

It follows that:

$$\mathbb{P}_1(X = \mathbf{No} \wedge Y \leq q) = 1 - \mathbb{P}_1(X = \mathbf{Yes} \vee Y > q) \geq 1 - \mathbb{P}_1(X = \mathbf{Yes}) - \mathbb{P}_1(Y > q) \geq 0.5. \quad (1)$$

Turing machine TM is deterministic, so if it terminates after reading prefix  $\pi$  of some run  $\rho$ , then it terminates after reading prefix  $\pi$  of any run. As a consequence, the event  $Y \leq q$  can be represented as a union of  $\ell$  cones where  $\ell \leq |\Sigma|^q = 2^q$  since  $\Sigma = \{a, b\}$  in  $\mathcal{M}$ :

$$\{\rho : Y(\rho) \leq q\} = \bigcup_{i=1}^{\ell} \text{Cone}(\pi_i),$$

where all  $\pi_i \in \Sigma^q$  are distinct. The event  $X = \mathbf{No} \wedge Y \leq q$  is a refinement of the event  $Y \leq q$ , so it may also be represented as

$$\{\rho : X = \mathbf{No} \wedge Y(\rho) \leq q\} = \bigcup_{i=1}^m \text{Cone}(\pi_i), \quad (2)$$

where  $m \leq \ell \leq 2^q$ . Since every path in  $\mathcal{M}(0)$  of length  $q$  has probability  $0.5^q$ , we get by (2)

$$\mathbb{P}_1(X = \mathbf{No} \wedge Y(\rho) \leq q) = \mathbb{P}_1\left(\bigcup_{i=1}^m \text{Cone}(\pi_i)\right) = \sum_{i=1}^m \mathbb{P}_1(\pi_i) = m0.5^q.$$

Then by (1) it follows that  $m \geq 2^{q-1}$ .

Now, we run TM on the Markov chain  $\mathcal{M}(\epsilon)$  where  $\epsilon = 0.5 - \alpha^{\frac{1}{q}} 2^{\frac{1-q}{q}}$  if  $q > 0$  and  $\epsilon = 0.25$  in the degenerated case of  $q = 0$ .

► **Claim.**  $\epsilon > 0$ .

Let us write  $\mathbb{P}_2$  for the probability measure of TM on the input  $\mathcal{M}(\epsilon)$ . The distance between  $\mathcal{M}(0)$  and  $\mathcal{M}(\epsilon)$  is 1, since  $\epsilon > 0$ . As a consequence, TM should return answer **Yes** on this input with probability  $\geq 1 - \alpha$ , or equivalently answer **No** with probability  $< \alpha$ . We show, however, that the probability of **No** is  $\geq \alpha$ :

$$\mathbb{P}_2(X = \mathbf{No} \wedge Y \leq q) = \sum_{i=1}^m \mathbb{P}_2(\pi_i) \quad \text{by (2)}$$

$$\begin{aligned}
 &= \sum_{i=1}^m (0.5 + \epsilon)^{u_i} (0.5 - \epsilon)^{q-u_i} && u_i \text{ is number of } b\text{'s in } \pi_i \\
 &\geq \sum_{i=1}^m (0.5 - \epsilon)^q = m(0.5 - \epsilon)^q \\
 &\geq 2^{q-1} (0.5 - \epsilon)^q = \alpha. && \text{by } m \geq 2^{q-1}..
 \end{aligned}$$

We obtain a contradiction, thus the assumed machine TM does not exist.

**Part II.** Suppose for a contradiction that for some  $\alpha < \frac{1}{2}, \delta < 1$  there exists an algorithm  $\text{Alg}_{\alpha,\delta}$  that solves the problem defined in the theorem and almost-surely terminates. Then then this algorithm can solve the problem  $B_\alpha$  in the following way:

1. Use  $\text{Alg}_{\alpha,\delta}$  to compute a confidence interval  $I$  for the total variation distance between  $\mathcal{M}(0)$  and  $\mathcal{M}(\tau)$ . Algorithm  $\text{Alg}_{\alpha,\delta}$  can sample any number of paths from  $\mathcal{M}(0)$ . Observe that in  $\mathcal{M}(\tau)$  probability of seeing states  $a$  and  $b$  remains constant over time. Thus, sampling multiple paths from  $\mathcal{M}(\tau)$  by  $\text{Alg}_{\alpha,\delta}$  can be replaced by sampling a single path from  $\mathcal{M}(\tau)$ .
2. Output **Yes** if  $1 \in I$ , **No** if  $0 \in I$ .

We have shown that for any  $\alpha < \frac{1}{2}$  the problem  $B_\alpha$  cannot be solved by an algorithm that almost-surely terminates. As a consequence, the algorithm  $\text{Alg}_{\alpha,\delta}$  cannot exist. ◀

From Part II, it follows that there is no statistical algorithm even for fixed  $\alpha$  and  $\delta$ .

## 6 Estimability: Finite-trace distance

In Section 6.1 we show how to estimate the distance given by traces of a fixed length. In Section 6.2 we show how to reduce the problem of computing the finite-trace distance  $D_{\text{FT}}$  (where traces of arbitrary lengths are considered) to computing a constant number of fixed-length distances.

### 6.1 Estimates for fixed length

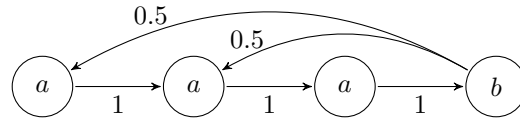
Given two Markov chains  $\mathcal{M}_1$  and  $\mathcal{M}_2$  we wish to estimate the finite-trace distance for fixed length  $k \in \mathbb{N}$

$$D_{\text{FT}}^k = \sup_{w \in \Sigma^k} |\mathbb{P}_1(w) - \mathbb{P}_2(w)|.$$

There is  $m = |\Sigma|^k$  words in  $\Sigma^k$  (we enumerate them as  $w_1, \dots, w_m$ ), so the traces of length  $k$  follow a multinomial distribution, i.e. for  $i = 1, 2$   $\sum_{j=1}^m \mathbb{P}_i(w_j) = 1$ .

We present a statistical procedure that estimates  $D_{\text{FT}}^k$  with arbitrary precision. For  $j \leq |\Sigma|^k$  we call a *contrast*  $\Delta_j$  the difference in probabilities of trace  $w_j$  between  $\mathcal{M}_1$  and  $\mathcal{M}_2$ :  $\Delta_j = |\mathbb{P}_1(w_j) - \mathbb{P}_2(w_j)|$ . The distance  $D_{\text{FT}}^k$  is the maximum over all such contrasts  $D_{\text{FT}}^k = \max_{j \leq m} \Delta_j$ . We use the statistical procedure of [19] to simultaneously estimate all contrasts. We sample random paths from both Markov chains, and let  $n_i^j$  denote the number of observations of trace  $w_j$  in a Markov chain  $\mathcal{M}_i$ . We write  $n_i = \sum_{j \leq m} n_i^j$  for the sum of all observations in  $\mathcal{M}_i$ . The estimator of  $\mathbb{P}_i(w_j)$  is  $\tilde{p}_i^j = \frac{n_i^j}{n_i}$ , and the estimator of  $\Delta_j$  is  $\tilde{\Delta}_j = |\tilde{p}_1^j - \tilde{p}_2^j|$ .





■ **Figure 2** Markov chain with 4 states. The leftmost state is 6-deterministic, but not deterministic.

► **Theorem 9** ([19]). *As  $n_1, n_2 \rightarrow \infty$  the probability approaches  $1 - \alpha$  that simultaneously for all contrasts*

$$|\Delta_j - \tilde{\Delta}_j| \leq S_j M \quad \text{where} \quad S_j = \sqrt{\frac{\tilde{p}_1^j - (\tilde{p}_1^j)^2}{n_1} + \frac{\tilde{p}_2^j - (\tilde{p}_2^j)^2}{n_2}},$$

and  $M$  is the square root of the  $\frac{1-\alpha}{100}$  percentile of the  $\chi^2$  distribution with  $|\Sigma|^k$  degrees of freedom.

The procedure for estimating  $D_{\text{FT}}^k$  works as follows. For  $\epsilon, \alpha > 0$  we sample paths from  $\mathcal{M}_1$  and  $\mathcal{M}_2$  until, by Theorem 9, with probability  $1 - \alpha$  for all contrasts  $|\Delta_j - \tilde{\Delta}_j| \leq \epsilon$ . Then with probability  $1 - \alpha$  it holds that  $|D_{\text{FT}}^k - \max_{j \leq m} \tilde{\Delta}_j| \leq \epsilon$ .

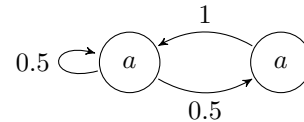
## 6.2 Estimates for unbounded length

Intuitively, the longer the path, the less probable it is, and the less distance it can cause. However, this is only true if along the path probabilistic choices are made repeatedly.

► **Definition 10.** In a Markov chain  $\mathcal{M}$ , a state  $s \in S$  is  $k$ -deterministic, if there exists a word  $w$  of length  $k$ , such that  $\mathbb{P}^s(w) = 1$ . Otherwise,  $s$  is  $k$ -branching. A state  $s \in S$  is deterministic, if it is  $k$ -deterministic for all  $k \in \mathbb{N}$ .

► **Lemma 11.** *If  $s \in S$  is  $k$ -branching, it is also  $(k + 1)$ -branching. Dually, if it is  $k$ -deterministic, it is also  $(k - 1)$ -deterministic.*

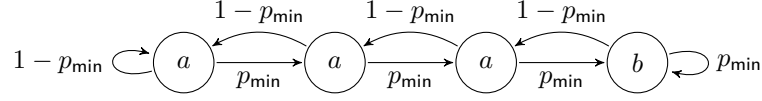
► **Example 12.** Every state is trivially 1-deterministic. In Figure 3, the leftmost state is 3-deterministic and 4-branching. The states of the MC on the right are deterministic.



► **Lemma 13.** *Consider a state  $s$  in a Markov chain  $\mathcal{M}$  with  $n$  states. If state  $s$  is  $n^2$ -deterministic, then it is deterministic.*

Before proceeding to the proof, notice that even though it may seem that every branching state must be  $n + 1$  branching, this is not the case in general. Observe the counterexample in Figure 2. The leftmost state is 6-deterministic (only the word  $aaabaa$  can be generated), while  $n = 4$ .

**Proof.** Consider state  $s$  that is  $n^2$ -deterministic and assume for contradiction that  $s$  is not deterministic. Let  $N > n^2$  be the smallest number such that  $s$  is  $N$ -branching, and thus not  $(N - 1)$ -branching. Then there exist two paths  $\pi = s_1, s_2, \dots, s_N$  and  $\pi' = s_1, s'_2, \dots, s'_N$  such that  $s_1 = s$  and for  $i = 1, 2, \dots, N - 1$ , we have  $L(s_i) = L(s'_i)$  and  $L(s_N) \neq L(s'_N)$ . Looking at a sequence of pairs  $(s_1, s_1), (s_2, s'_2), \dots, (s_{N-1}, s'_{N-1})$ , since there are at most  $n^2$  possible pairs of states over  $S$ , by the pigeon-hole principle at least two pairs will be



■ **Figure 3** Markov chain, s.t.  $\mathbb{P}(a) = \mathbb{P}(aa) = \mathbb{P}(aaa) = 1$ ,  $\mathbb{P}(aaab) = p_{\min}^3$ ,  $\mathbb{P}(aaaa) = 1 - p_{\min}^3$ .

repeating in the observed sequence, say  $(s_i, s'_i) = (s_j, s'_j)$ , where  $i < j$ . But then the paths  $\pi'' = s_1, s_2, \dots, s_i, s_{j+1}, \dots, s_N$  and  $\pi''' = s_1, s_2, \dots, s_i, s_{j+1}, \dots, s_N$  have  $M < N$  states and they witness that  $s_1$  is  $M$ -branching, which by Lemma 11 is in contradiction with  $s$  being  $(N - 1)$ -deterministic. ◀

► **Lemma 14.** *If a state  $s \in S$  is  $k$ -branching, then any word of length  $k$  starting from  $s$  has probability at most  $(1 - p_{\min}^{k-1})$ , i.e.,  $\forall w \in \Sigma^k : \mathbb{P}^s(w) \leq 1 - p_{\min}^{k-1}$ .*

To illustrate this, observe the Markov chain in Figure 3 with leftmost initial state.

**Proof.** Let  $w \in \Sigma^k$ . Since  $s$  is  $k$ -branching, there exists a word  $w' \in \Sigma^k$  such that  $w' \neq w$  and  $\mathbb{P}^s(w') > 0$ . Hence there exists at least one path with  $k - 1$  transitions, producing the trace  $w'$ , and thus  $\mathbb{P}^s(w') \geq p_{\min}^{k-1}$ . Finally,  $\mathbb{P}^s(w) \leq 1 - \mathbb{P}^s(w') \leq 1 - p_{\min}^{k-1}$ . ◀

We show that, for estimating the finite trace distance with the required precision  $\epsilon$ , it suffices to infer probabilities of the words up to some finite length  $k$ , which depends on  $\epsilon$ . The idea is that paths that become deterministic before step  $k$  do not change their probability afterwards, while all other paths together have the probability bounded by  $\epsilon$ .

► **Lemma 15.** *Let  $s$  be a  $n^2$ -deterministic state in a Markov chain  $\mathcal{M}$  with  $n$  states. Then there are words  $u, z$ , such that  $|z| + |u| \leq n$ ,  $|u| \geq 1$ , and  $\mathbb{P}^s(zu^\omega) = 1$ .*

This motivates the following definition, where  $\text{pref}(w)$  denotes the set of all prefixes of the  $(\omega)$ -word  $w$ .

► **Definition 16.** A word  $w \in \Sigma^+$  is called  $k$ -ultimately periodic in a Markov chain  $\mathcal{M}$  if  $\mathbb{P}(w) > 0$  and there exists a word  $u$  such that  $w \in \text{pref}(\Sigma^k u^\omega)$  and  $1 \leq |u| \leq n$ , where  $n$  is the number of states in  $\mathcal{M}$ . ◀

Intuitively, for sufficiently long word  $w$  and large  $\epsilon$ , if  $\mathbb{P}(w) > \epsilon$  and  $w$  is  $k$ -ultimately periodic, then it enters within  $k$  steps a BSCC, which is bisimilar to a cycle (all transition probabilities are 1). One can also prove that this is the only way for a  $\omega$ -word to achieve a probability greater than  $\epsilon$ .

For a word  $w$  we write  $B^k(w)$  for the set of paths that are labelled by  $w$ , have a positive probability and where all states up to step  $k$  are  $n^2$ -branching:

$$B^k(w) = \{ \pi = s_1 \cdots s_{|w|} \in L^{-1}(w) \mid \mathbb{P}(\pi) > 0 \wedge \forall i \leq \min(k, |w|). s_i \text{ is } n^2\text{-branching} \} .$$

In a similar way, we write  $D^k(w)$  for the set of paths that enter a  $(n^2)$ -deterministic state before step  $k$

$$D^k(w) = \{ \pi = s_1 \cdots s_{|w|} \in L^{-1}(w) \mid \mathbb{P}(\pi) > 0 \wedge \exists i \leq \min(k, |w|). s_i \text{ is } n^2\text{-deterministic} \} .$$

For any  $k$ , we can partition paths labeled by  $w$  into  $B^k$ -paths and  $D^k$ -paths:

$$\mathbb{P}(w) = \sum_{\pi \in L^{-1}(w)} \mathbb{P}(\pi) = \sum_{\pi \in B^k(w)} \mathbb{P}(\pi) + \sum_{\pi \in D^k(w)} \mathbb{P}(\pi) . \quad (3)$$

Now we show that the probability of  $B^k$ -paths diminishes exponentially with length  $k$ :

► **Lemma 17.** Consider a Markov chain  $\mathcal{M}$  with  $n$  states. For every  $k \in \mathbb{N}$  and word  $w$ , if  $|w| > k$  then

$$\sum_{\pi \in B^k(w)} \mathbb{P}(\pi) \leq (1 - p_{\min}^{n^2})^{\lfloor \frac{k}{n^2} \rfloor}.$$

► **Lemma 18.** Let  $w$  be a word in a Markov chain  $\mathcal{M}$  with  $n$  states. For every  $\epsilon > 0$ , if  $\mathbb{P}(w) > \epsilon$  and  $|w| > k$  then  $w$  is  $k$ -ultimately periodic in  $\mathcal{M}$ , where  $k = n^2 \lceil \frac{\log \epsilon}{\log(1 - p_{\min}^{n^2})} \rceil + n$ .

**Proof.** Assume that  $|w| > k$ . We split paths labelled by  $w$  into  $B^{k-n}(w)$  and  $D^{k-n}(w)$  as in (3):

$$\mathbb{P}(w) = \sum_{s_1 \cdots s_{|w|} \in L^{-1}(w)} \mathbb{P}(s_1 \cdots s_{|w|}) = \sum_{\substack{s_1 \cdots s_{|w|} \in \\ B^{k-n}(w)}} \mathbb{P}(s_1 \cdots s_{|w|}) + \sum_{\substack{s_1 \cdots s_{|w|} \in \\ D^{k-n}(w)}} \mathbb{P}(s_1 \cdots s_{|w|}). \quad (4)$$

By Lemma 17 we get

$$\sum_{s_1 \cdots s_{|w|} \in B^{k-n}(w)} \mathbb{P}(s_1 \cdots s_{|w|}) \leq \epsilon. \quad (5)$$

Now, from the assumption  $\mathbb{P}(w) > \epsilon$ , (4) and (5), it follows that

$$\sum_{s_1 \cdots s_{|w|} \in D^{k-n}(w)} \mathbb{P}(s_1 \cdots s_{|w|}) > 0.$$

This implies that there is a path  $\pi = s_1 \cdots s_{|w|} \in D^{k-n}(w)$ . By definition of  $D^{k-n}(w)$ ,  $\pi$  has a  $n^2$ -deterministic state before step  $k - n$ , and w.l.o.g. let  $s_{k-n}$  be that state. By Lemma 15, every positive word from state  $s_{k-n}$  is a prefix of  $zu^\omega$  for some words  $z, u$  such that  $|z| + |u| \leq n$ . Therefore  $w \in \text{pref}(yzu^\omega)$ , where  $y = L(s_1 \cdots s_{k-n})$ , i.e.  $w$  is  $|k|$ -ultimately periodic. ◀

► **Lemma 19.** Consider a Markov chain  $\mathcal{M}$  with  $n$  states. Let  $w$  be a  $k$ -ultimately periodic word in  $\mathcal{M}$ , and  $x$  be a prefix of  $w$  such that  $|x| > k + n$ . Then

$$\mathbb{P}(x) - \mathbb{P}(w) \leq (1 - p_{\min}^{n^2})^{\lfloor \frac{k-n}{n^2} \rfloor}.$$

► **Theorem 20.** Consider Markov chains  $\mathcal{M}_1$  and  $\mathcal{M}_2$  that have at most  $n$  states. For  $\epsilon > 0$  it holds that

$$|\text{D}_{\text{FT}}(\mathcal{M}_1, \mathcal{M}_2) - \max_{i \leq k} \text{D}_{\text{FT}}^i(\mathcal{M}_1, \mathcal{M}_2)| \leq \epsilon, \quad \text{where } k = n^2 \lceil \frac{\log \epsilon}{\log(1 - p_{\min}^{n^2})} \rceil + 2n.$$

**Proof.** We show that for any word  $w \in \Sigma^+$ :

$$\left| |\mathbb{P}_1(w) - \mathbb{P}_2(w)| - |\mathbb{P}_1(w \downarrow k) - \mathbb{P}_2(w \downarrow k)| \right| \leq \epsilon. \quad (6)$$

For  $|w| \leq k$  (6) holds trivially. Suppose that  $|w| \geq k$  and consider two cases.

1. If  $\mathbb{P}_i(w \downarrow k) > \epsilon$ , then by Lemma 18  $w \downarrow k$  is  $(k - n)$ -ultimately periodic. Then by Lemma 19  $\mathbb{P}_i(w \downarrow k) \leq \mathbb{P}_i(w) + \epsilon$ .

## 20:12 Linear Distances between Markov Chains

2. If  $\mathbb{P}_i(w \downarrow k) \leq \epsilon$ , then clearly  $\mathbb{P}_i(w \downarrow k) \leq \mathbb{P}_i(w) + \epsilon$ . Both cases can be summarised by

$$\mathbb{P}_i(w) \leq \mathbb{P}_i(w \downarrow k) \leq \mathbb{P}_i(w) + \epsilon. \quad (7)$$

W.l.o.g assume that  $\mathbb{P}_1(w) \geq \mathbb{P}_2(w)$ . Then by (7)

$$\mathbb{P}_1(w \downarrow k) - \mathbb{P}_2(w \downarrow k) \geq \mathbb{P}_1(w) - \mathbb{P}_2(w) - \epsilon,$$

which implies (6). ◀

## 7 Consequences and Discussion

We now discuss the consequences of the (in)estimability results for several specific subclasses of  $\omega$ -regular languages, captured topologically, logically, or by automata. We also remark on the estimability in case when the transition probabilities have finite precision.

### 7.1 Topology

**Negative result for clopen sets.** Note that the proof of inestimability was based on the ability to express the events  $E_n$  for any  $n \in \mathbb{N}$ :

$E_n =$  “there is at most  $c_n = (0.5 + \tau/2)n$  symbols  $b$  in the prefix path of length  $n$ .”

Observe that each  $E_n$  can be expressed as finite union of cones, each expressing exact positions of  $a$ 's and  $b$ 's in the first  $n$  steps. For instance, for  $\tau = 0.2$ , the event  $E_2$ , “there is at most 1 symbol  $b$  in the first 2 steps,” can be described by the union  $\text{Cone}(aa) \cup \text{Cone}(ab) \cup \text{Cone}(ba)$ .

Since finite unions of cones form exactly the clopen sets, the lowest class  $\Delta_1$  in the Borel hierarchy, it follows that distances based on any class in the hierarchy are inestimable.

**Positive result for the infinite-trace distance.** Using the result on finite-trace distance, we can prove that the infinite-trace distance  $D_{\Gamma}$  of Example 4 is also estimable. Indeed, the distance is non-zero only due to  $k$ -ultimately periodic  $\omega$ -words with positive probability. By Lemma 19 we can provide confidence intervals for these probabilities through the  $k$ -prefixes using the fixed-length distance  $D_{\Gamma}^k$ .

### 7.2 Logic

**Negative result for LTL.** The LTL distance as in Example 7 is again inestimable since we can express the event  $E_n$  in LTL by a finite composition of operators  $\mathbf{X}, \wedge, \vee$  (notably this fragment induces the same distance as LTL [5]). Indeed, for instance, for  $\tau = 0.2$ , the event  $E_{10}$ , “there is at most 6 symbols  $b$  in the path prefix of length 10,” is equivalent to “at least 4 symbols  $a$  in the path prefix of length  $n$ ,” and it can be described by a disjunction of  $\binom{10}{4}$  formulae, each determining the possible position of symbols  $a$ , resulting in a formula  $(a \wedge \mathbf{X}a \wedge \mathbf{X}^2a \wedge \mathbf{X}^3a) \vee (a \wedge \mathbf{X}a \wedge \mathbf{X}^2a \wedge \mathbf{X}^4a) \vee \dots \vee (\mathbf{X}^7a \wedge \mathbf{X}^8a \wedge \mathbf{X}^9a \wedge \mathbf{X}^{10}a)$ .

**Positive result for LTL(FG,GF).** The distance generated by the fragment of LTL described by combining operators **FG** and **GF** and Boolean operators is estimable. Notice that the probability of the property  $\varphi \equiv \mathbf{FG}\varphi'$  equals the probability of reaching a BSCC such that  $\varphi'$  holds in all of its states, while the probability of property  $\varphi \equiv \mathbf{GF}\varphi'$  equals the probability that every BSCC contains a state which satisfies  $\varphi'$ . Hence, properties expressed in this fragment of LTL can be checked by inferring all BSCCs of a chain and a simple analysis of them. The statistical estimation of all BSCCs for labelled Markov chains where only the minimal transition probability is known is possible and is shown in [10].

### 7.3 Automata

**Negative result for automata distances.** For the class of all deterministic Rabin automata (DRA), the distance (as in Example 6) is inestimable. This is implied by the inestimability for clopen sets or for LTL. Further, we can also directly encode the event  $E_n$  that “at least  $k$  symbols  $a$  are observed in the path of length  $n$ ” by an automaton: the DRA counts how many symbols  $a$  are seen in the prefix up to length  $n$ ; this can be done with  $k \cdot n$  states where the automaton is in a state  $s_{k',n'}$  if and only if in the  $n' \leq n$  prefix of the input word, there are  $k' \leq k$  symbols  $a$ .

**Positive result for fixed-size automata.** When restricting to the class of DRA of size at most  $k \in \mathbb{N}$ , the distance  $D_{DRA \leq k}$  can be estimated. A naive algorithm amounts to enumerating all automata up to given size  $k$ , then applying statistical model checking to infer the probability of satisfying the automata in each of the Markov chains, and checking for which automaton the probability difference in the two chains is maximized. Statistically inferring the probability of whether a (black-box) Markov chain satisfies a property given by a DRA is a subroutine of the procedure for statistical model checking Markov chains for LTL, described in [10].

### 7.4 Finite Precision

When the transition probabilities have finite precision, e.g. are given by at most two decimal digits, several negative results turn positive. Finite precision allows us to learn the MCs exactly with high probability, by rounding the learnt transition probabilities to the closest multiple of the precision. Subsequently, we can approximate the distance by the algorithms applicable in the white-box setting. In case of the total variation distance, one can apply the approximation algorithm of [8]; for trace distances, the approximation algorithm of [5] is also available. In particular, for the special case of the trace equivalence  $E_{FT}$  we can leverage the fact that Markov chains are equivalent when all their traces up to length  $|\mathcal{M}_1| + |\mathcal{M}_2| - 1$  have equal probability. With the assumption of finite precision one can get by sampling the exact distribution of such traces with high confidence. Note that the same algorithm can not be applied without assuming finite precision, since arbitrarily small difference in chains cannot be detected.

## 8 Conclusions and Future Work

We have introduced a linear-distance framework for Markov chains and considered estimating the distances in the black-box setting from simulation runs. We investigated several distances, delimiting the (in)estimability boarder for distances given topologically, logically, and by automata. As the next step, it is desirable to look for practical algorithms that would

converge fast on practical benchmarks. Another direction is to characterize the largest language for which the distance can be estimated, and, dually, the smallest language that cannot be estimated.

---

### References

- 1 Alessandro Abate. Approximation metrics based on probabilistic bisimulations for general state-space Markov processes: A survey. *Electr. Notes Theor. Comput. Sci.*, 297:3–25, 2013.
- 2 Giorgio Bacci, Giovanni Bacci, Kim G. Larsen, and Radu Mardare. The BisimDist library: Efficient computation of bisimilarity distances for Markovian models. In *QEST*, pages 278–281, 2013.
- 3 Giorgio Bacci, Giovanni Bacci, Kim G. Larsen, and Radu Mardare. Computing behavioral distances, compositionally. In *MFCS*, pages 74–85, 2013.
- 4 Giorgio Bacci, Giovanni Bacci, Kim G. Larsen, and Radu Mardare. On-the-fly exact computation of bisimilarity distances. In *TACAS*, pages 1–15, 2013.
- 5 Giorgio Bacci, Giovanni Bacci, Kim G. Larsen, and Radu Mardare. Converging from branching to linear metrics on Markov chains. In *ICTAC*, pages 349–367, 2015.
- 6 Giorgio Bacci, Giovanni Bacci, Kim G. Larsen, and Radu Mardare. On the total variation distance of semi-Markov chains. In *FoSSaCS*, pages 185–199, 2015.
- 7 Paul Caspi and Albert Benveniste. Toward an approximation theory for computerised control. In *EMSOFT*, pages 294–304, 2002.
- 8 Taolue Chen and Stefan Kiefer. On the total variation distance of labelled Markov chains. In *CSL-LICS*, pages 33:1–33:10, 2014.
- 9 Przemyslaw Daca, Thomas A. Henzinger, Jan Křetínský, and Tatjana Petrov. Linear distances between Markov chains. Technical Report abs/1605.00186, arXiv.org, 2014.
- 10 Przemyslaw Daca, Thomas A. Henzinger, Jan Křetínský, and Tatjana Petrov. Faster statistical model checking for unbounded temporal properties. *TACAS*, pages 112–129, 2016.
- 11 Luca de Alfaro, Marco Faella, and Mariëlle Stoelinga. Linear and branching metrics for quantitative transition systems. In *ICALP*, pages 97–109, 2004.
- 12 Luca de Alfaro, Rupak Majumdar, Vishwanath Raman, and Mariëlle Stoelinga. Game relations and metrics. In *LICS*, pages 99–108, 2007.
- 13 Josee Desharnais, Vineet Gupta, Radha Jagadeesan, and Prakash Panangaden. Metrics for labeled Markov systems. In *CONCUR*, pages 258–273, 1999.
- 14 Jyotirmoy V. Deshmukh, Rupak Majumdar, and Vinayak S. Prabhu. Quantifying conformance using the Skorokhod metric. In *CAV*, pages 234–250, 2015.
- 15 Laurent Doyen, Thomas A. Henzinger, and Jean-François Raskin. Equivalence of labeled Markov chains. *Int. J. Found. Comput. Sci.*, 19(3):549–563, 2008.
- 16 Norm Ferns, Prakash Panangaden, and Doina Precup. Metrics for finite Markov decision processes. In *AAAI*, pages 950–951, 2004.
- 17 Nathanaël Fijalkow, Stefan Kiefer, and Mahsa Shirmohammadi. Trace refinement in labelled Markov decision processes. In *FOSSACS*, pages 303–318, 2016.
- 18 Antoine Girard and George J. Pappas. Approximate bisimulation: A bridge between computer science and control theory. *Eur. J. Control*, 17(5-6):568–578, 2011.
- 19 Leo A Goodman. Simultaneous confidence intervals for contrasts among multinomial populations. *The Annals of Mathematical Statistics*, pages 716–725, 1964.
- 20 Manfred Jaeger, Hua Mao, Kim G. Larsen, and Radu Mardare. Continuity properties of distances for Markov processes. In *QEST*, pages 297–312, 2014.
- 21 Stefan Kiefer and A. Prasad Sistla. Distinguishing hidden markov chains. In *31th Annual ACM/IEEE Symposium on Logic in Computer Science, LICS 2016*, 2016.
- 22 Kim G. Larsen and Arne Skou. Bisimulation through probabilistic testing. In *POPL*, pages 344–352, 1989.

- 23 Rupak Majumdar and Vinayak S. Prabhu. Computing the Skorokhod distance between polygonal traces. In *HSCC*, pages 199–208, 2015.
- 24 Ilya Tkachev and Alessandro Abate. On approximation metrics for linear temporal model-checking of stochastic systems. In *HSCC*, pages 193–202, 2014.
- 25 Franck van Breugel, Babita Sharma, and James Worrell. Approximating a behavioural pseudometric without discount for probabilistic systems. In *FLOSSACS*, pages 123–137, 2007.
- 26 Franck van Breugel and James Worrell. Approximating and computing behavioural distances in probabilistic transition systems. *Theor. Comput. Sci.*, 360(1-3):373–385, 2006.