

Accurate Approximate Diagnosability of Stochastic Systems

Nathalie Bertrand¹, Serge Haddad^{2,*}, and Engel Lefaucheux^{1,2}

¹ Inria Rennes, France `nathalie.bertrand@inria.fr`

² LSV, ENS Cachan & CNRS & Inria, Université Paris-Saclay, France
`serge.haddad@lsv.fr` `engel.lefaucheux@irisa.fr`

Abstract. Diagnosis of partially observable stochastic systems prone to faults was introduced in the late nineties. Diagnosability, *i.e.* the existence of a diagnoser, may be specified in different ways: (1) exact diagnosability (called A-diagnosability) requires that almost surely a fault is detected and that no fault is erroneously claimed while (2) approximate diagnosability (called ε -diagnosability) allows a small probability of error when claiming a fault and (3) accurate approximate diagnosability (called AA-diagnosability) requires that this error threshold may be chosen arbitrarily small. Here we mainly focus on approximate diagnoses. We first refine the almost sure requirement about finite delay introducing a uniform version and showing that while it does not discriminate between the two versions of exact diagnosability this is no more the case in approximate diagnosis. Then we establish a complete picture for the decidability status of the diagnosability problems: (uniform) ε -diagnosability and uniform AA-diagnosability are undecidable while AA-diagnosability is decidable in PTIME, answering a longstanding open question.

Keywords: automata for system analysis and programme verification

1 Introduction

Diagnosis and diagnosability. The increasing use of software systems for critical operations motivates the design of fast automatic detection of malfunctions. In general, diagnosis raises two important issues: deciding whether the system is *diagnosable* and, in the positive case, synthesizing a *diagnoser* possibly satisfying additional requirements about memory size, implementability, etc. One of the proposed approaches consists in modelling these systems by partially observable labelled transition systems (LTS) [11]. In such a framework, diagnosability requires that the occurrence of unobservable faults can be deduced from the previous and subsequent observable events. Formally, an LTS is diagnosable if there exists a diagnoser that satisfies *reactivity* and *correctness* constraints. Reactivity requires that if a fault occurred, the diagnoser eventually detects it. Correctness asks that the diagnoser only claims the existence of a fault when there actually was

* This author was partly supported by ERC project EQualIS (FP7-308087).

one. Diagnosability for LTS was shown to be decidable in PTIME [7] while the diagnoser itself could be of size exponential w.r.t. the size of the LTS. Diagnosis has been extended to numerous models (Petri nets [3], pushdown systems [8], etc.) and settings (centralized, decentralized, distributed), and have had an impact on important application areas, *e.g.* for telecommunication network failure diagnosis. Also, several contributions, gathered under the generic name of active diagnosis, focus on enforcing the diagnosability of a system [4, 5, 10, 13].

Diagnosis of stochastic systems. Diagnosis was also considered in a quantitative setting, and namely for probabilistic labelled transition systems (pLTS) [1, 12], that can be seen as Markov chains in which the transitions are labelled with events. Therefore, one can define a probability measure over infinite runs. In that context, the specification of reactivity and correctness can be relaxed. Here, reactivity only asks to detect faults almost surely (*i.e.* with probability 1). This weaker reactivity constraint takes advantage of probabilities to rule out negligible behaviours. For what concerns correctness, three natural variants can be considered. *A-diagnosability* sticks to strong correctness and therefore asks the diagnoser to only claim fault occurrences when a fault is certain. *ε -diagnosability* tolerates small errors, allowing to claim a fault if the conditional probability that no fault occurred does not exceed ε . *AA-diagnosability* requires the pLTS to be ε -diagnosable for all positive ε , allowing the designer to select a threshold according to the criticality of the system. A-diagnosability and AA-diagnosability were introduced in [12]. Recently, we focused on semantical and algorithmic issues related to A-diagnosability, and in particular we established that A-diagnosability is PSPACE-complete [1]. When it comes to approximate diagnosability (*i.e.* ε and AA-diagnosability), up to our knowledge, a (PTIME-checkable) sufficient condition for AA-diagnosability [12] has been given, but no decidability result is known.

Contributions. Our contributions are twofold. From a semantical point of view, we investigate the specification of reactivity, introducing *uniform reactivity* which requires that once a fault occurs, the probability of detection when time elapses converges to 1 uniformly w.r.t. faulty runs. Uniformity provides the user with a stronger guarantee about the delay before detection. We show that uniform A-diagnosability and A-diagnosability coincide while this is no longer the case for approximate diagnosability. From an algorithmic point of view, we first show that ε -diagnosability and its uniform version are undecidable. Then we characterize AA-diagnosability as a separation property between labelled Markov chains (LMC), precisely a *distance* 1 between appropriate pairs of LMCs built from the pLTS. Thanks to [6], this yields a polynomial time algorithm for AA-diagnosability. AA-diagnosability can thus be checked more efficiently than A-diagnosability (PTIME vs PSPACE), yet, surprisingly, contrary to A-diagnosers, AA-diagnosers may require infinite memory. Finally, we show that uniform AA-diagnosability is undecidable.

Organization. In Section 2, we introduce the different variants of diagnosability and establish the full hierarchy between these specifications. In Section 3, we address the decidability and complexity issues related to approximate diagnosis. Full proofs can be found in the companion research report [2].

2 Specification of Diagnosability

2.1 Probabilistic Labelled Transition Systems

To represent stochastic discrete event systems, we use transition systems labelled with events and in which the transition function is probabilistic.

Definition 1. A probabilistic labelled transition system (*pLTS*) is a tuple $\mathcal{A} = \langle Q, q_0, \Sigma, T, \mathbf{P} \rangle$ where:

- Q is a finite set of states with $q_0 \in Q$ the initial state;
- Σ is a finite set of events;
- $T \subseteq Q \times \Sigma \times Q$ is a set of transitions;
- $\mathbf{P} : T \rightarrow \mathbb{Q}_{>0}$ is the probability function fulfilling for every $q \in Q$:

$$\sum_{(q,a,q') \in T} \mathbf{P}[q, a, q'] = 1.$$

Observe that a pLTS is a labelled transition system (LTS) equipped with transition probabilities. The transition relation of the underlying LTS is defined by: $q \xrightarrow{a} q'$ for $(q, a, q') \in T$; this transition is then said to be *enabled* in q .

Let us now introduce some important notions and notations that will be used throughout the paper. A *run* ρ of a pLTS \mathcal{A} is a (finite or infinite) sequence $\rho = q_0 a_0 q_1 \dots$ such that for all i , $q_i \in Q$, $a_i \in \Sigma$ and when q_{i+1} is defined, $q_i \xrightarrow{a_i} q_{i+1}$. The notion of run can be generalized, starting from an arbitrary state q . We write Ω for the set of all infinite runs of \mathcal{A} starting from q_0 , assuming the pLTS is clear from context. When it is finite, ρ ends in a state q and its *length*, denoted $|\rho|$, is the number of actions occurring in it. Given a finite run $\rho = q_0 a_0 q_1 \dots q_n$ and a (finite or infinite) run $\rho' = q_n a_n q_{n+1} \dots$, we call concatenation of ρ and ρ' and we write $\rho\rho'$ for the run $q_0 a_0 q_1 \dots q_n a_n q_{n+1} \dots$; the run ρ is then a *prefix* of $\rho\rho'$, which we denote $\rho \preceq \rho\rho'$. The *cylinder* generated by a finite run ρ consists of all infinite runs that extend ρ : $\text{Cyl}(\rho) = \{\rho' \in \Omega \mid \rho \preceq \rho'\}$. The sequence associated with $\rho = q a_0 q_1 \dots$ is the word $\sigma_\rho = a_0 a_1 \dots$, and we write indifferently $q \xrightarrow{\rho} q'$ or $q \xrightarrow{\sigma_\rho} q'$ (resp. $q \xrightarrow{\rho} q'$ or $q \xrightarrow{\sigma_\rho} q'$) for an infinite (resp. finite) run ρ . A state q is *reachable* (from q_0) if there exists a run such that $q_0 \xrightarrow{\rho} q$, which we alternatively write $q_0 \Rightarrow q$. The language of pLTS \mathcal{A} consists of all infinite words that label runs of \mathcal{A} and is formally defined as $\mathcal{L}^\omega(\mathcal{A}) = \{\sigma \in \Sigma^\omega \mid q_0 \xrightarrow{\sigma} \}$.

Forgetting the labels and merging (and summing the probabilities of) the transitions with same source and target, a pLTS yields a discrete time Markov chain (DTMC). As usual for DTMC, the set of infinite runs of \mathcal{A} is the support of a probability measure defined by Caratheodory's extension theorem from the probabilities of the cylinders:

$$\mathbb{P}_{\mathcal{A}}(\text{Cyl}(q_0 a_0 q_1 \dots q_n)) = \mathbf{P}[q_0, a_1, q_1] \cdots \mathbf{P}[q_{n-1}, a_n, q_n] .$$

When \mathcal{A} is fixed, we may omit the subscript. To simplify, for ρ a finite run, we will sometimes abuse notation and write $\mathbb{P}(\rho)$ for $\mathbb{P}(\text{Cyl}(\rho))$. If R is a (denumerable) set of finite runs (such that no run is a prefix of another one), we write $\mathbb{P}(R)$ for $\sum_{\rho \in R} \mathbb{P}(\rho)$.

2.2 Partial Observation and Ambiguity

Beyond the pLTS model for stochastic discrete event systems, in order to formalize problems related to fault diagnosis, we partition Σ into two disjoint sets Σ_o and Σ_u , the sets of *observable* and of *unobservable events*, respectively. Moreover, we distinguish a special *fault event* $\mathbf{f} \in \Sigma_u$. Let σ be a finite word over Σ ; its length is denoted $|\sigma|$. The projection of words onto Σ_o is defined inductively by: $\pi(\varepsilon) = \varepsilon$; for $a \in \Sigma_o$, $\pi(\sigma a) = \pi(\sigma)a$; and $\pi(\sigma a) = \pi(\sigma)$ for $a \notin \Sigma_o$. We write $|\sigma|_o$ for $|\pi(\sigma)|$. When σ is an infinite word, its projection is the limit of the projections of its finite prefixes. As usual the projection mapping is extended to languages: for $L \subseteq \Sigma^*$, $\pi(L) = \{\pi(\sigma) \mid \sigma \in L\}$. With respect to the partition of $\Sigma = \Sigma_o \uplus \Sigma_u$, a pLTS \mathcal{A} is *convergent* if, from any reachable state, there is no infinite sequence of unobservable events: $\mathcal{L}^\omega(\mathcal{A}) \cap \Sigma^* \Sigma_u^\omega = \emptyset$. When \mathcal{A} is convergent, for every $\sigma \in \mathcal{L}^\omega(\mathcal{A})$, $\pi(\sigma) \in \Sigma_o^\omega$. In the rest of the paper we assume that pLTS are convergent. We will use the terminology *sequence* for a word $\sigma \in \Sigma^* \cup \Sigma^\omega$, and an *observed sequence* for a word $\sigma \in \Sigma_o^* \cup \Sigma_o^\omega$. The projection of a sequence to Σ_o is thus an observed sequence.

The *observable length* of a run ρ denoted $|\rho|_o \in \mathbb{N} \cup \{\infty\}$, is the number of observable events that occur in it: $|\rho_o| = |\sigma_\rho|_o$. A *signalling run* is a finite run ending with an observable event. Signalling runs are precisely the relevant runs w.r.t. partial observation issues since each observable event provides an external observer additional information about the execution. In the sequel, SR denotes the set of signalling runs, and SR_n the set of signalling runs of observable length n . Since we assume pLTS to be convergent, for every $n > 0$, SR_n is equipped with a probability distribution defined by assigning measure $\mathbb{P}(\rho)$ to each $\rho \in \text{SR}_n$. Given ρ a finite or infinite run, and $n \leq |\rho|_o$, $\rho_{\downarrow n}$ denotes the signalling subrun of ρ of observable length n . For convenience, we consider the empty run q_0 to be the single signalling run, of null length. For an observed sequence $\sigma \in \Sigma_o^*$, we define its cylinder $\text{Cyl}(\sigma) = \sigma \Sigma_o^*$ and the associated probability $\mathbb{P}(\text{Cyl}(\sigma)) = \mathbb{P}(\{\rho \in \text{SR}_{|\sigma|} \mid \pi(\rho) = \sigma\})$, often shortened as $\mathbb{P}(\sigma)$.

Let us now partition runs depending on whether they contain a fault or not. A run ρ is *faulty* if σ_ρ contains \mathbf{f} , otherwise it is *correct*. For $n \in \mathbb{N}$, we write F_n (resp. C_n) for the set of faulty (resp. correct) signalling runs of length n , and further define the set of all faulty and correct signalling runs $\text{F} = \cup_{n \in \mathbb{N}} \text{F}_n$ and $\text{C} = \cup_{n \in \mathbb{N}} \text{C}_n$. W.l.o.g., by considering two copies of each state, we assume that the state space Q is partitioned into correct states and faulty states: $Q = Q_f \uplus Q_c$ such that faulty (resp. correct) states, *i.e.* states in Q_f (resp. Q_c) are only reachable by faulty (resp. correct) runs. An infinite (resp. finite) observed sequence $\sigma \in \Sigma_o^\omega$ (resp. Σ_o^*) is *ambiguous* if there exists a correct infinite (resp. signalling) run ρ and a faulty infinite (resp. signalling) run ρ' such that $\pi(\rho) = \pi(\rho') = \sigma$.

2.3 Fault Diagnosis

Whatever the considered notion of diagnosis in probabilistic systems, *reactivity* requires that when a fault occurs, a diagnoser almost surely will detect it after a finite delay. We refine this requirement by also considering *uniform reactivity* ensuring that given any positive probability threshold α there exists a delay n_α such that the probability to exceed this delay is less or equal than α . Here uniformity means “independently of the faulty run”.

Similarly, *correctness* of the diagnosis may be specified in different ways. Since we focus on approximate diagnosis, a fault can be claimed after an ambiguous observed sequence. This implies that ambiguity should be quantified in order to assess the quality of the diagnosis. To formalise this idea, with every observed sequence $\sigma \in \Sigma_o^*$ we associate a *correctness proportion*

$$\text{CorP}(\sigma) = \frac{\mathbb{P}(\{\rho \in C_{|\sigma|} \mid \pi(\rho) = \sigma\})}{\mathbb{P}(\{\rho \in C_{|\sigma|} \cup F_{|\sigma|} \mid \pi(\rho) = \sigma\})},$$

which is the conditional probability that a signalling run is correct given that its observed sequence is σ . Thus approximate diagnosability also denoted ε -*diagnosability* allows the diagnoser to claim a fault when the correctness proportion does not exceed ε while accurate approximate diagnosability denoted *AA-diagnosability* ensures that ε can be chosen as small as desired but still positive.

Definition 2 (Diagnosability notions). *Let \mathcal{A} be a pLTS and $\varepsilon \geq 0$.*

- \mathcal{A} is ε -diagnosable if for all faulty run $\rho \in F$ and all $\alpha > 0$ there exists $n_{\rho,\alpha}$ such that for all $n \geq n_{\rho,\alpha}$:

$$\mathbb{P}(\{\rho' \in \text{SR}_{n+|\rho|_o} \mid \rho \preceq \rho' \wedge \text{CorP}(\pi(\rho')) > \varepsilon\}) \leq \alpha \mathbb{P}(\rho).$$

\mathcal{A} is uniformly ε -diagnosable if $n_{\rho,\alpha}$ does not depend on ρ .

- \mathcal{A} is (uniformly) AA-diagnosable if it is (uniformly) ε -diagnosable for all $\varepsilon > 0$.

Two variants of diagnosability for stochastic systems were introduced in [12]: AA-diagnosability and A-diagnosability. *A-diagnosability*, which corresponds to exact diagnosis, is nothing else but 0-diagnosability in Definition 2 wording. By definition, A-diagnosability implies AA-diagnosability which implies ε -diagnosability for all $\varepsilon > 0$. Observe also that since the faulty run ρ (and so $\mathbb{P}(\rho)$) is fixed, ε -diagnosability can be rewritten:

$$\lim_{n \rightarrow \infty} \mathbb{P}(\{\rho' \in \text{SR}_{n+|\rho|_o} \mid \rho \preceq \rho' \wedge \text{CorP}(\pi(\rho')) > \varepsilon\}) = 0.$$

We now provide examples that illustrate these notions. Consider \mathcal{A}_1 , the pLTS represented on Figure 1. We claim that \mathcal{A}_1 is AA-diagnosable but neither A-diagnosable, nor uniformly AA-diagnosable. We only give here intuitions on these claims, and refer the reader to the proof of Proposition 3 in [2]. First an

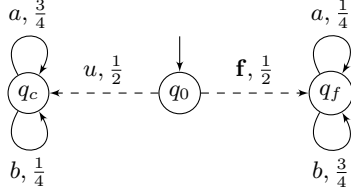


Fig. 1. An AA-diagnosable pLTS \mathcal{A}_1 , that is neither A-diagnosable, nor uniformly AA-diagnosable.

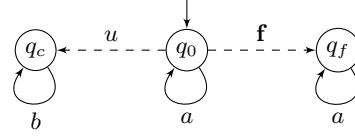


Fig. 2. An uniformly AA-diagnosable pLTS \mathcal{A}_2 , that is not A-diagnosable.

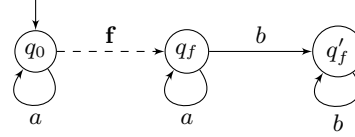


Fig. 3. An A-diagnosable pLTS \mathcal{A}_3 .

ε -diagnoser will look at the proportion of b occurrences and if the sequence is “long” enough and the proportion is “close” to $\frac{3}{4}$, it will claim a fault. However, the delay $n_{\alpha, \rho}$ before claiming a fault cannot be selected independently of the faulty run. Indeed, given the faulty run $\rho_n = q_0 \mathbf{f} q_f (a q_f)^n$, we let $p_{n, m}$ for the probability of extensions of ρ_n by m observable events and with correctness proportion below ε . In order for $p_{n, m}$ to exceed $1 - \alpha$, m must depend on n . So \mathcal{A}_1 is not uniformly AA-diagnosable. \mathcal{A}_1 is neither A-diagnosable since all observed sequences of faulty runs are ambiguous.

Consider now the pLTS \mathcal{A}_2 depicted in Figure 2, for which we consider a uniform distribution on the outgoing edges from q_0 . First note that every faulty run $(q_0 a)^i q_0 \mathbf{f} (q_f a)^j q_f$ has a correct run, namely $q_0 (a q_0)^{i+j}$ with the same observed sequence. So \mathcal{A}_2 is not A-diagnosable. Yet, we argue that it is uniformly AA-diagnosable. The correctness proportion of a faulty run (exponentially) decreases with respect to its length. So the worst run to be considered for the diagnoser is $q_0 \mathbf{f} q_f a q_f$ implying uniformity.

Consider the pLTS \mathcal{A}_3 from Figure 3, with uniform distributions in q_0 and q_f . Viewed as an LTS, it is not diagnosable, since the observed sequence a^ω is ambiguous and forbids the diagnosis of faulty runs without any occurrence of b . On the contrary, let $\rho = q_0 (a q_0)^x \mathbf{f} q_f (a q_f)^y (b q'_f)^z$ be an arbitrary faulty run. If $z > 0$ then $\text{CorP}(\pi(\rho)) = 0$. Otherwise $\mathbb{P}(\{\rho' \in \text{SR}_{n+|\rho|_o} \mid \rho \preceq \rho' \wedge \text{CorP}(\pi(\rho')) > 0\}) = \frac{1}{2^n} \mathbb{P}(\rho)$ and so \mathcal{A}_3 is A-diagnosable.

Proposition 3 establishes the exact relations between the different specifications. Observe that uniform AA-diagnosability is strictly stronger than AA-diagnosability while A-diagnosability and uniform A-diagnosability are equivalent.

Proposition 3. – A pLTS is A-diagnosable if and only if it is uniformly A-diagnosable.

– There exists an AA-diagnosable pLTS, not uniformly $\frac{1}{2}$ -diagnosable and so not uniformly AA-diagnosable.

- There exists a uniformly AA-diagnosable pLTS, not A-diagnosable.

Although we have not explicitly defined diagnosers for diagnosable pLTS, given a fixed threshold $\varepsilon > 0$, a simple diagnoser would monitor the sequence of observed events σ , compute the current correctness proportion, and output “faulty” if $\text{CorP}(\sigma)$ is below ε . However such an ε -diagnoser may need an infinite memory. This contrasts with the case of A-diagnosability for which finite-memory diagnosers suffice [12].

Proposition 4. *There exists an AA-diagnosable pLTS, thus ε -diagnosable for every $\varepsilon > 0$, that admits no finite-memory diagnoser when $0 < \varepsilon \leq \frac{1}{2}$.*

Proof. Consider \mathcal{A}_1 the AA-diagnosable pLTS of Figure 1 and assume there exists a diagnoser with m states for some threshold $0 < \varepsilon \leq \frac{1}{2}$. After any sequence a^n , it cannot claim a fault. So there exist $1 \leq i < j \leq m + 1$ such that the diagnoser is in the same state after observing a^i and a^j .

Consider the faulty run $\rho = q_0 f q_f (a q_f)^i$. Due to the reactivity requirement, there must be a run $\rho \rho'$ for which the diagnoser claims a fault. This implies that for all n , the diagnoser claims a fault after $\rho_n = \rho (a q_f)^{n(j-i)} \rho'$ but $\lim_{n \rightarrow \infty} \text{CorP}(\pi(\rho_n)) = 1$, which contradicts the correctness requirement. \square

3 Analysis of Approximate Diagnosability

A-diagnosability was proved to be a PSPACE-complete problem [1]. We now focus on the other notions of approximate diagnosability introduced in Definition 2, and study their decidability and complexity.

Reducing the emptiness problem for probabilistic automata [9] (PA), we obtain the following first result:

Theorem 5. *For any rational $0 < \varepsilon < 1$, the ε -diagnosability and uniform ε -diagnosability problems are undecidable for pLTS.*

We now turn to the decidability status of AA-diagnosability and uniform AA-diagnosability. We prove that AA-diagnosability can be solved in polynomial time by establishing a characterization in terms of distance on labelled Markov chains; this constitutes the most technical contribution of this section.

A *labelled Markov chain* (LMC) is a pLTS where every event is observable: $\Sigma = \Sigma_o$. In order to exploit results of [6] on LMC in our context of pLTS, we introduce the mapping \mathcal{M} that performs *in polynomial time* the probabilistic closure of a pLTS w.r.t. the unobservable events and produces an LMC. For sake of simplicity, we denote by \mathcal{A}_q , the pLTS \mathcal{A} where the initial state has been substituted by q .

Definition 6. *Given a pLTS $\mathcal{A} = \langle Q, q_0, \Sigma, T, \mathbf{P} \rangle$ with $\Sigma = \Sigma_o \uplus \Sigma_u$, the labelled Markov chain $\mathcal{M}(\mathcal{A}) = \langle Q, q_0, \Sigma_o, T', \mathbf{P}' \rangle$ is defined by:*

- $T' = \{(q, a, q') \mid \exists \rho \in \text{SR}_1(\mathcal{A}_q) \rho = q \cdots a q'\}$ (and so a is observable).
- For all $(q, a, q') \in T'$, $\mathbf{P}'(q, a, q') = \mathbb{P}\{\rho \in \text{SR}_1(\mathcal{A}_q) \mid \rho = q \cdots a q'\}$.

Let E be an *event* of Σ^ω (i.e. a measurable subset of Σ^ω for the standard measure), we denote by $\mathbb{P}^{\mathcal{M}}(E)$ the probability that event E occurs in the LMC \mathcal{M} . Given two LMC \mathcal{M}_1 and \mathcal{M}_2 , the (probabilistic) distance between \mathcal{M}_1 and \mathcal{M}_2 generalizes the concept of distance for distributions. Given an event E , $|\mathbb{P}^{\mathcal{M}_1}(E) - \mathbb{P}^{\mathcal{M}_2}(E)|$ expresses the absolute difference between the probabilities that E occurs in \mathcal{M}_1 and in \mathcal{M}_2 . The distance is obtained by getting the supremum over the events.

Definition 7. Let \mathcal{M}_1 and \mathcal{M}_2 be two LMC over the same alphabet Σ . Then $d(\mathcal{M}_1, \mathcal{M}_2)$ the distance between \mathcal{M}_1 and \mathcal{M}_2 is defined by:

$$d(\mathcal{M}_1, \mathcal{M}_2) = \sup(\mathbb{P}^{\mathcal{M}_1}(E) - \mathbb{P}^{\mathcal{M}_2}(E) \mid E \text{ event of } \Sigma^\omega).$$

The *distance 1 problem* asks, given labelled Markov chains \mathcal{M}_1 and \mathcal{M}_2 , whether $d(\mathcal{M}_1, \mathcal{M}_2) = 1$. We summarize in the next proposition, the results by Chen and Kiefer on LMC that we use later.

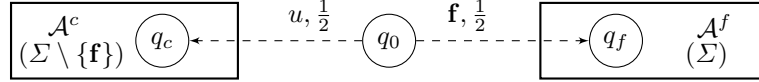
Proposition 8 ([6]).

- Given two LMC $\mathcal{M}_1, \mathcal{M}_2$, there exists an event E such that:

$$d(\mathcal{M}_1, \mathcal{M}_2) = \mathbb{P}^{\mathcal{M}_1}(E) - \mathbb{P}^{\mathcal{M}_2}(E).$$

- The distance 1 problem for LMC is decidable in polynomial time.

Towards the decidability of AA-diagnosability, let us first explain how to solve the problem on a subclass of pLTS called *initial-fault pLTS*. Informally, an initial-fault pLTS \mathcal{A} consists of two disjoint pLTS \mathcal{A}^f and \mathcal{A}^c and an initial state q_0 with an outgoing unobservable correct transition leading to \mathcal{A}^c and a transition labelled by \mathbf{f} leading to \mathcal{A}^f (see the figure below). Moreover no faulty transitions occur in \mathcal{A}^c . We denote such a pLTS by $\mathcal{A} = \langle q_0, \mathcal{A}^f, \mathcal{A}^c \rangle$.



The next lemma establishes a strong connection between distance of LMC and diagnosability of initial-fault pLTS.

Lemma 9. Let $\mathcal{A} = \langle q_0, \mathcal{A}^f, \mathcal{A}^c \rangle$ be an initial-fault pLTS. Then \mathcal{A} is AA-diagnosable if and only if $d(\mathcal{M}(\mathcal{A}^f), \mathcal{M}(\mathcal{A}^c)) = 1$.

Proof. We write \mathbb{P} , \mathbb{P}_f and \mathbb{P}_c for the probability distributions of pLTS \mathcal{A} , \mathcal{A}^f and \mathcal{A}^c . By construction of $\mathcal{M}(\mathcal{A}^f)$ and $\mathcal{M}(\mathcal{A}^c)$, for every observed sequence σ , $\mathbb{P}^{\mathcal{M}(\mathcal{A}^f)}(\sigma) = \mathbb{P}_f(\sigma)$ and similarly $\mathbb{P}^{\mathcal{M}(\mathcal{A}^c)}(\sigma) = \mathbb{P}_c(\sigma)$. In words, the mapping \mathcal{M} leaves unchanged the probability of occurrence of an observed sequence.

- If \mathcal{A} is AA-diagnosable, for every $\varepsilon > 0$ and every faulty run ρ :

$$\lim_{n \rightarrow \infty} \mathbb{P}(\{\rho' \in \text{SR}_{n+|\rho|_o} \mid \rho \preceq \rho' \wedge \text{CorP}(\pi(\rho')) > \varepsilon\}) = 0. \quad (1)$$

Pick some $0 < \varepsilon < 1$. By applying Equation (1) on the faulty run $\rho_f = q_0 \mathbf{f} q_f$ with $|\pi(\rho_f)| = 0$, there exists some $n \in \mathbb{N}$ such that:

$$\mathbb{P}(\{\rho \in \text{SR}_n \mid \rho_f \preceq \rho \wedge \text{CorP}(\pi(\rho)) > \varepsilon\}) \leq \varepsilon.$$

Let \mathfrak{S} be the set of observed sequences of faulty runs with length n and correctness proportion not exceeding threshold ε :

$$\mathfrak{S} = \{\sigma \in \Sigma_o^n \mid \exists \rho \in \text{SR}_n, \pi(\rho) = \sigma \wedge \rho_f \preceq \rho \wedge \text{CorP}(\sigma) \leq \varepsilon\}.$$

$E = \text{Cyl}(\mathfrak{S})$ is the event consisting of the infinite suffixes of those sequences. Let us show that $\mathbb{P}_c(E) \leq \frac{\varepsilon}{1-\varepsilon}$ and $\mathbb{P}_f(E) \geq 1 - 2\varepsilon$.

$$\mathbb{P}_f(E) = 1 - 2\mathbb{P}(\{\rho \in \text{SR}_n \mid \rho_f \preceq \rho \wedge \text{CorP}(\pi(\rho)) > \varepsilon\}) \geq 1 - 2\varepsilon.$$

The factor 2 comes from the probability $\frac{1}{2}$ in \mathcal{A} to enter \mathcal{A}^f that \mathbb{P}_f does not take into account contrary to \mathbb{P} .

Moreover, for every observed sequence $\sigma \in \mathfrak{S}$, there exists a faulty run ρ such that $\pi(\rho) = \sigma$. Thus, $\text{CorP}(\sigma) \leq \varepsilon$. Using the definition of CorP :

$$\text{CorP}(\sigma) = \frac{\mathbb{P}(\{\rho \in \text{C}_n \mid \pi(\rho) = \sigma\})}{\mathbb{P}(\{\rho \in \text{SR}_n \mid \pi(\rho) = \sigma\})} = \frac{\mathbb{P}_c(\sigma)}{\mathbb{P}_c(\sigma) + \mathbb{P}_f(\sigma)} \leq \varepsilon.$$

Thus, $\mathbb{P}_c(\sigma) \leq \frac{\varepsilon}{1-\varepsilon} \mathbb{P}_f(\sigma)$. Hence:

$$\mathbb{P}_c(E) = \sum_{\sigma \in \mathfrak{S}} \mathbb{P}_c(\sigma) \leq \sum_{\sigma \in \mathfrak{S}} \frac{\varepsilon}{1-\varepsilon} \mathbb{P}_f(\sigma) = \frac{\varepsilon}{1-\varepsilon} \mathbb{P}_f(E) \leq \frac{\varepsilon}{1-\varepsilon}.$$

Therefore $d(\mathcal{M}(\mathcal{A}^c), \mathcal{M}(\mathcal{A}^f)) \geq \mathbb{P}_f(E) - \mathbb{P}_c(E) \geq 1 - \varepsilon(2 + \frac{1}{1-\varepsilon})$. Letting ε go to 0, we obtain $d(\mathcal{M}(\mathcal{A}^c), \mathcal{M}(\mathcal{A}^f)) = 1$.

• Conversely assume that $d(\mathcal{M}(\mathcal{A}^f), \mathcal{M}(\mathcal{A}^c)) = 1$. Due to Proposition 8, there exists an event $E \subseteq \Sigma_o^\omega$ such that $\mathbb{P}_f(E) = 1$ and $\mathbb{P}_c(E) = 0$.

For all $n \in \mathbb{N}$, let \mathfrak{S}_n be the set of prefixes of length n of the observed sequences of E : $\mathfrak{S}_n = \{\sigma \in \Sigma_o^n \mid \exists \sigma' \in E, \sigma \preceq \sigma'\}$.

For all $\varepsilon > 0$, let $\mathfrak{S}_n^\varepsilon$ be the subset of sequences of \mathfrak{S}_n whose correctness proportion exceeds threshold ε : $\mathfrak{S}_n^\varepsilon = \{\sigma \in \mathfrak{S}_n \mid \text{CorP}(\sigma) > \varepsilon\}$.

As $\bigcap_{n \in \mathbb{N}} \text{Cyl}(\mathfrak{S}_n) = E$, $\lim_{n \rightarrow \infty} \mathbb{P}_c(\mathfrak{S}_n) = \mathbb{P}_c(E) = 0$.

So $\lim_{n \rightarrow \infty} \mathbb{P}_c(\mathfrak{S}_n^\varepsilon) = 0$.

On the other hand for all $n \in \mathbb{N}$,

$$\mathbb{P}_c(\mathfrak{S}_n^\varepsilon) = \sum_{\sigma \in \mathfrak{S}_n^\varepsilon} \mathbb{P}_c(\sigma) > \sum_{\sigma \in \mathfrak{S}_n^\varepsilon} \frac{\varepsilon}{1-\varepsilon} \mathbb{P}_f(\sigma) = \frac{\varepsilon}{1-\varepsilon} \mathbb{P}_f(\mathfrak{S}_n^\varepsilon).$$

Therefore we have $\lim_{n \rightarrow \infty} \mathbb{P}_f(\mathfrak{S}_n^\varepsilon) = 0$.

Let ρ be a faulty run and $\alpha > 0$. There exists $n_\alpha \geq |\rho|_o$ such that for all $n \geq n_\alpha$, $\mathbb{P}_f(\mathfrak{S}_n^\varepsilon) \leq \alpha$. Let $n \geq n_\alpha$, and $\tilde{\mathfrak{S}}_n$ be the set of observed sequences of length n triggered by a run with prefix ρ and whose correctness proportion exceeds ε :

$$\tilde{\mathfrak{S}}_n = \{\sigma \in \Sigma_o^n \mid \exists \rho' \in \text{SR}_n, \rho \preceq \rho' \wedge \pi(\rho') = \sigma \wedge \text{CorP}(\sigma) > \varepsilon\}.$$

Let us prove that $\mathbb{P}(\tilde{\mathfrak{S}}_n) \leq \alpha$. On the one hand, since $\mathbb{P}_f(\mathfrak{S}_n) \geq \mathbb{P}_f(E) = 1$, $\mathbb{P}_f(\tilde{\mathfrak{S}}_n \cap (\Sigma_o^n \setminus \mathfrak{S}_n)) = 0$. On the other hand, since $\mathbb{P}_f(\mathfrak{S}_n^\varepsilon) < \alpha$, $\mathbb{P}_f(\tilde{\mathfrak{S}}_n \cap \mathfrak{S}_n) \leq$

$\mathbb{P}_f(\mathfrak{S}_n^\varepsilon) \leq \alpha$. Thus $\mathbb{P}_f(\tilde{\mathfrak{S}}_n) = \mathbb{P}_f(\tilde{\mathfrak{S}}_n \cap \mathfrak{S}_n) + \mathbb{P}_f(\tilde{\mathfrak{S}}_n \cap (\Sigma_o^n \setminus \mathfrak{S}_n)) \leq \alpha$. Because α was taken arbitrary, we obtain that $\lim_{n \rightarrow \infty} \mathbb{P}_f(\tilde{\mathfrak{S}}_n) = 0$.

Observe now that $\mathbb{P}(\{\rho' \in \text{SR}_n \mid \rho \preceq \rho' \wedge \text{CorP}(\pi(\rho')) > \varepsilon\}) = \frac{1}{2} \mathbb{P}_f(\tilde{\mathfrak{S}}_n)$. Therefore, $\lim_{n \rightarrow \infty} \mathbb{P}(\{\rho' \in \text{SR}_n \mid \rho \preceq \rho' \wedge \text{CorP}(\pi(\rho')) > \varepsilon\}) = 0$. So \mathcal{A} is AA-diagnosable. \square

In order to understand why characterizing AA-diagnosability for general pLTS is more involved, let us study the pLTS \mathcal{A}_2 presented in Figure 2 where outgoing transitions of any state are equidistributed. Recall that \mathcal{A}_2 is AA-diagnosable (and even uniformly AA-diagnosable).

Let us look at the distance between pairs of a correct and a faulty states of \mathcal{A} that can be reached by runs with the same observed sequence. On the one hand, $d(\mathcal{M}(\mathcal{A}_{q_0}), \mathcal{M}(\mathcal{A}_{q_f})) \leq \frac{1}{2}$ since for any event E either (1) $a^\omega \in E$ implying $\mathbb{P}^{\mathcal{M}(\mathcal{A}_{q_f})}(E) = 1$ and $\mathbb{P}^{\mathcal{M}(\mathcal{A}_{q_0})}(E) \geq \frac{1}{2}$ or (2) $a^\omega \notin E$ implying $\mathbb{P}^{\mathcal{M}(\mathcal{A}_{q_f})}(E) = 0$ and $\mathbb{P}^{\mathcal{M}(\mathcal{A}_{q_0})}(E) \leq \frac{1}{2}$. On the other hand, $d(\mathcal{M}(\mathcal{A}_{q_c}), \mathcal{M}(\mathcal{A}_{q_f})) = 1$ since $\mathbb{P}^{\mathcal{M}(\mathcal{A}_{q_f})}(a^\omega) = 1$ and $\mathbb{P}^{\mathcal{M}(\mathcal{A}_{q_c})}(a^\omega) = 0$.

We claim that the pair (q_0, q_f) is irrelevant, since the correct state q_0 does not belong to a bottom strongly connected component (BSCC) of the pLTS, while (q_c, q_f) is relevant since q_c belongs to a BSCC triggering a “recurrent” ambiguity.

The next theorem characterizes AA-diagnosability, establishing the soundness of this intuition. Moreover, it states the complexity of deciding AA-diagnosability.

Theorem 10. *Let \mathcal{A} be a pLTS. Then, \mathcal{A} is AA-diagnosable if and only if for every correct state q_c belonging to a BSCC and every faulty state q_f reachable by runs with same observed sequence, $d(\mathcal{M}(\mathcal{A}_{q_c}), \mathcal{M}(\mathcal{A}_{q_f})) = 1$.*

The AA-diagnosability problem is decidable in polynomial time for pLTS.

The full proof of Theorem 10 is given [2]. Let us sketch the key ideas to establish the characterization of AA-diagnosability in terms of the distance 1 problem.

The left-to-right implication is the easiest one, and is proved by contraposition. Assume there exist two states in \mathcal{A} , $q_c \in Q_c$ belonging to a BSCC and $q_f \in Q_f$ reachable resp. by ρ_c and ρ_f with $\pi(\rho_c) = \pi(\rho_f)$, and with $d(\mathcal{M}(\mathcal{A}_{q_c}), \mathcal{M}(\mathcal{A}_{q_f})) < 1$. Applying Lemma 9 to the initial-fault pLTS $\mathcal{A}' = \langle q'_0, \mathcal{A}_{q_f}, \mathcal{A}_{q_c} \rangle$, one deduces that \mathcal{A}' is not AA-diagnosable. First we relate the probabilities of runs in \mathcal{A} and \mathcal{A}' . Then we show that considering the additional faulty runs with same observed sequence as ρ_f does not make \mathcal{A} AA-diagnosable.

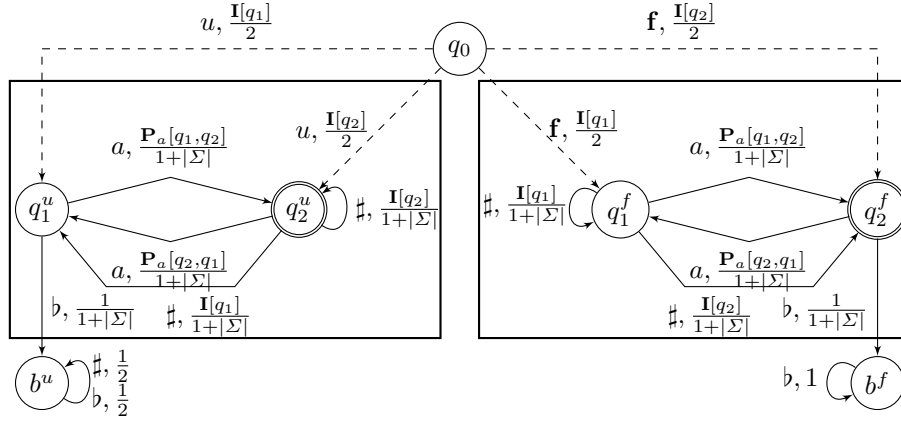
The right-to-left implication is harder to establish. For ρ_0 a faulty run, $\alpha > 0$, $\varepsilon > 0$, $\sigma_0 = \pi(\rho_0)$ and $n_0 = |\sigma_0|$, we start by extending the runs with observed sequences σ_0 by n_b observable events where n_b is chosen in order to get a high probability that the runs end in a BSCC. For such an observed sequence $\sigma \in \Sigma_o^{m_b}$, we partition the possible runs with observed sequence $\sigma_0\sigma$ into three sets: \mathfrak{R}_σ^F is the subset of faulty runs; \mathfrak{R}_σ^C (resp. \mathfrak{R}_σ^T) is the set of correct runs ending (resp. not ending) in a BSCC. At first, we do not take into account the “transient” runs in \mathfrak{R}_σ^T . We apply Lemma 9 to obtain an integer n_σ such that from \mathfrak{R}_σ^F and \mathfrak{R}_σ^C we can diagnose with (appropriate) high probability and low correctness proportion after

n_σ observations. Among the runs that trigger diagnosable observed sequences, some exceed the correctness proportion ε , when taking into account the runs from \mathfrak{R}_σ^T . Yet, we show that the probability of such runs is small, when cumulated over all extensions σ , leading to the required upper bound α .

Using the characterization, one can easily establish the complexity of AA-diagnosability. Indeed, reachability of a pair of states with the same observed sequence is decidable in polynomial time by an appropriate “self-synchronized product” of the pLTS. Since there are at most a quadratic number of pairs to check, and given that the distance 1 problem can be decided in polynomial time, the PTIME upper-bound follows.

In contrast, uniform AA-diagnosability is shown to be undecidable by a reduction from the emptiness problem for probabilistic automata, that is more involved than the one for Theorem 5.

Theorem 11. *The uniform AA-diagnosability problem is undecidable for pLTS.*



The reduction is illustrated above, and we sketch here the undecidability proof. Assuming there exists a word $w \in \Sigma^*$ accepted with probability greater than $\frac{1}{2}$ in the probabilistic automaton. We pick arbitrary $\alpha < 1$ and n_α . Then, one can exhibit a faulty signalling run ρ_n with $\pi(\rho_n) = (w\sharp)^n$ for some appropriate n , such that for every extension $\rho_n \preceq \rho$ with $|\rho| = |\rho_n| + n_\alpha$, one has $\text{CorP}(\rho) > \frac{1}{2}$. This shows that the constructed pLTS is not uniformly $\frac{1}{2}$ -diagnosable.

Assuming now that all words are accepted with probability less than $\frac{1}{2}$. Then for any observed sequence $\sigma \in (\Sigma \cup \{\sharp\})^*$, $\text{CorP}(\sigma) \leq \frac{1}{2}$. After reaching a BSCC, the correctness proportion decreases uniformly, due to the \sharp -loop on b^u . Given positive α and ε , one can thus find integers n_0 and n_1 such that a BSCC is reached after n_0 observable events with probability at least $1 - \alpha$, and after n_1 more the correctness proportion, which was at most $\frac{1}{2}$, decreases below ε . This shows the uniform AA-diagnosability.

4 Conclusion

This paper completes our previous work [1] on diagnosability of stochastic systems, by giving here a full picture on approximate diagnosis. On the one hand, we performed a semantical study: we have refined the reactivity specification by introducing a uniform requirement about detection delay w.r.t. faults and studied its impact on both the exact and approximate case. On the other hand, we established decidability and complexity of all notions of approximate diagnosis: we have shown that (uniform) ε -diagnosability and uniform AA-diagnosability are undecidable while AA-diagnosability can be solved in polynomial time.

There are still interesting issues to be tackled, to continue our work on monitoring of stochastic systems. For example, prediction and prediagnosis, which are closely related to diagnosis and were analyzed in the exact case in [1], should be studied in the approximate framework.

References

1. Bertrand, N., Haddad, S., Lefaucheu, E.: Foundation of diagnosis and predictability in probabilistic systems. In: Proceedings of FSTTCS 2014. LIPIcs, vol. 29, pp. 417–429. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik (2014)
2. Bertrand, N., Haddad, S., Lefaucheu, E.: Accurate approximate diagnosability of stochastic systems. <https://hal.inria.fr/hal-01220954> (2015)
3. Cabasino, M., Giua, A., Lafortune, S., Seatzu, C.: Diagnosability analysis of unbounded Petri nets. In: Proceedings of CDC 2009. pp. 1267–1272. IEEE (2009)
4. Cassez, F., Tripakis, S.: Fault diagnosis with static and dynamic observers. *Fundamenta Informaticae* 88, 497–540 (2008)
5. Chanthery, E., Pencolé, Y.: Monitoring and active diagnosis for discrete-event systems. In: Proceedings of SP 2009. pp. 1545–1550. Elsevier (2009)
6. Chen, T., Kiefer, S.: On the total variation distance of labelled Markov chains. In: Proceedings of CSL-LICS 2014. pp. 33:1–33:10. ACM (2014)
7. Jiang, S., Huang, Z., Chandra, V., Kumar, R.: A polynomial algorithm for testing diagnosability of discrete-event systems. *IEEE Transactions on Automatic Control* 46(8), 1318–1321 (2001)
8. Morvan, C., Pinchinat, S.: Diagnosability of pushdown systems. In: Proceedings of HVC 2009. LNCS, vol. 6405, pp. 21–33. Springer (2009)
9. Paz, A.: *Introduction to Probabilistic Automata*. Academic Press (1971)
10. Sampath, M., Lafortune, S., Teneketzis, D.: Active diagnosis of discrete-event systems. *IEEE Transactions on Automatic Control* 43(7), 908–929 (1998)
11. Sampath, M., Sengupta, R., Lafortune, S., Sinnamohideen, K., Teneketzis, D.: Diagnosability of discrete-event systems. *IEEE Transactions on Automatic Control* 40(9), 1555–1575 (1995)
12. Thorsley, D., Teneketzis, D.: Diagnosability of stochastic discrete-event systems. *IEEE Transactions on Automatic Control* 50(4), 476–492 (2005)
13. Thorsley, D., Teneketzis, D.: Active acquisition of information for diagnosis and supervisory control of discrete-event systems. *Journal of Discrete Event Dynamic Systems* 17, 531–583 (2007)