



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

Fachbereich Elektrotechnik und Informationstechnik  
Bioinspired Communication Systems

# **Bayesian Inference of Information Transfer in Graph-Based Continuous-Time Multi-Agent Systems**

**Master- Thesis**

Elektro- und Informationstechnik

Eingereicht von

Gizem Ekinici

am

21.07.2020

1. Gutachten: Prof. Dr. techn. Heinz Koeppel
2. Gutachten: Dominik Linzner



## **Erklärung zur Abschlussarbeit gemäß §22 Abs. 7 und §23 Abs. 7 APB TU Darmstadt**

Hiermit versichere ich, Gizem Ekinici, die vorliegende Arbeit gemäß §22 Abs. 7 APB der TU Darmstadt ohne Hilfe Dritter und nur mit den angegebenen Quellen und Hilfsmitteln angefertigt zu haben. Alle Stellen, die Quellen entnommen wurden, sind als solche kenntlich gemacht worden. Diese Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen. Mir ist bekannt, dass im Falle eines Plagiats (§38 Abs.2 APB) ein Täuschungsversuch vorliegt, der dazu führt, dass die Arbeit mit 5,0 bewertet und damit ein Prüfungsversuch verbraucht wird. Abschlussarbeiten dürfen nur einmal wiederholt werden. Bei der abgegebenen Arbeit stimmen die schriftliche und die zur Archivierung eingereichte elektronische Fassung gemäß §23 Abs. 7 APB überein.

English translation for information purposes only:

Thesis statement pursuant to §22 paragraph 7 and §23 paragraph 7 of APB TU Darmstadt: I herewith formally declare that I, Gizem Ekinici, have written the submitted thesis independently pursuant to §22 paragraph 7 of APB TU Darmstadt. I did not use any outside support except for the quoted literature and other sources mentioned in the paper. I clearly marked and separately listed all of the literature and all of the other sources which I employed when producing this academic work, either literally or in content. This thesis has not been handed in or published before in the same or similar form. I am aware, that in case of an attempt at deception based on plagiarism (§38 Abs. 2 APB), the thesis would be graded with 5,0 and counted as one failed examination attempt. The thesis may only be repeated once. In the submitted thesis the written copies and the electronic version for archiving are pursuant to § 23 paragraph 7 of APB identical in content.

Darmstadt, den 21.07.2020

---

(Gizem Ekinici)



## Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.



# Contents

<b>List of Symbols</b>	<b>i</b>
<b>List of Figures</b>	<b>ii</b>
<b>1. Introduction</b>	<b>1</b>
1.1. Motivation . . . . .	2
1.2. Structure of the Thesis . . . . .	2
<b>2. Foundations</b>	<b>3</b>
2.1. Problem Formulation . . . . .	3
2.2. Continuous-Time Bayesian Networks . . . . .	4
2.2.1. Continuous-Time Markov Processes . . . . .	4
2.2.1.1. Homogenous Continuous-Time Markov Processes . . . . .	5
2.2.1.2. Conditional Markov Processes . . . . .	6
2.2.2. The CTBN Model . . . . .	7
2.3. Partially Observable Markov Decision Processes . . . . .	8
2.3.1. Exact Belief State Update . . . . .	10
2.3.2. Exact Belief State Update for CTMP by Filtering . . . . .	10
2.3.3. Belief State Update using Particle Filter . . . . .	11
2.3.3.1. Particle Filtering . . . . .	11
2.3.3.2. Marginalized Continuous-Time Markov Processes . . . . .	12
2.4. Sampling Algorithms . . . . .	13
2.4.1. Gillespie Algorithm for Generative CTBN . . . . .	13
2.4.2. Thinning Algorithm . . . . .	14
<b>3. Methodology</b>	<b>15</b>
3.1. The Model . . . . .	15
3.1.1. CTBN Model . . . . .	15
3.1.2. POMDP Model . . . . .	16
3.1.2.1. Observation Model . . . . .	17
3.1.2.2. Belief State . . . . .	17
3.1.2.3. Optimal Policy . . . . .	19
3.2. Inference of Observation Model . . . . .	19
3.2.1. Likelihood Model . . . . .	20
3.2.2. Inference under Noisy Observation Model . . . . .	20
3.3. Data Generation . . . . .	21
3.3.1. Sampling Algorithm . . . . .	21

<b>4. Results</b>	<b>23</b>
4.1. Configurations . . . . .	23
4.2. Simulation . . . . .	24
4.3. Inference of Observation Model . . . . .	27
4.3.1. Equivalence Classes . . . . .	27
4.3.2. Learning Observation Model . . . . .	29
4.3.3. Inference with Non-informative Priors on Parent Parameters . . . . .	33
4.3.4. Robustness Test under Channel Noise . . . . .	33
<b>5. Discussion</b>	<b>34</b>
<b>6. Outlook</b>	<b>35</b>
<b>Bibliography</b>	<b>37</b>
<b>A. Amalgamation Operation</b>	<b>38</b>
A.1. Amalgamation of Independent Processes . . . . .	38
<b>B. Marginalized Likelihood Function for Homogenous Continuous Time Markov Processes</b>	<b>39</b>
<b>C. Equivalence Classes of Observation Models</b>	<b>40</b>
C.1. Observation Models in Experiments . . . . .	44
<b>D. Additional Results</b>	<b>45</b>



# List of Symbols

$\mathcal{X}$	state space of random variable $X$
$X(t)$	value of random variable $X$ at time $t$
$X^{[0,T]}$	discrete valued trajectory of random variable $X$ in time interval $[0, T]$
$\mathbf{X}^T$	transpose of matrix/vector $\mathbf{X}$

# List of Figures

2.1. Communication model. . . . .	3
2.2. Two component CTBN . . . . .	4
2.3. A policy tree . . . . .	8
3.1. Hierarchical model. . . . .	15
3.2. Closer look to agent-environment interaction from the perspective of POMDP framework. . . . .	16
4.1. Parent trajectories and observation . . . . .	24
4.2. Belief state trajectories . . . . .	25
4.3. $Q_3$ and $X_3$ trajectories . . . . .	26
4.4. Degenerate marginal particle filter . . . . .	27
4.5. Equivalence classes in the case of exact belief update . . . . .	28
4.6. An equivalence class in the case of marginal particle filtering . . . . .	29
4.7. Average log-likelihood in the case of exact belief update . . . . .	30
4.8. Average log-likelihood in the case of marginal particle filtering . . . . .	30
4.9. AUROC results over increasing number of samples . . . . .	32
4.10. AUPR results over increasing number of samples . . . . .	32
C.1. Parent trajectories for the models leading to the same belief state . . . . .	41
C.2. Observation, belief state and $Q_3$ trajectories derived by $\psi_1$ in Equation C.1 corresponding to parent trajectories in Figure C.1 . . . . .	41
C.3. Observation, belief state and $Q_3$ trajectories derived by $\psi_2$ in Equation C.2 corresponding to parent trajectories in Figure C.1 . . . . .	42
C.4. Parent trajectories for the models leading to the same belief state . . . . .	43
C.5. Observation, belief state and $Q_3$ trajectories derived by $\psi_1$ in Equation C.3 corresponding to parent trajectories in Figure C.4 . . . . .	43
C.6. Observation, belief state and $Q_3$ trajectories derived by $\psi_2$ in Equation C.4 corresponding to parent trajectories in Figure C.4 . . . . .	44
D.1. ROC curves with n number of trajectories for dataset generated using exact belief state update, $\psi_0$ -vs-rest . . . . .	45
D.2. ROC curves with n number of trajectories for dataset generated using particle filtering, $\psi_0$ -vs-rest . . . . .	46

# 1. Introduction

In multi-agent systems, the well-being of the population strongly relies on the cooperation between the agents. Such systems can be found in nature from cellular level to swarms [15, 22]. They consist of relatively simple individuals and as limited abilities as these individuals might have, they give rise to a complex behaviour of the population through cooperation. These complex behaviors might be crucial for the survival of the population, e.g. running away from predators, partitioning of vital substances.

The cooperation of the agents is optimized through the exchange of information between the agents. The individual interests might conflict, as in many cases it would. However, the blabla. Cells exist in stochastic environments. In order to maintain life, they are required to process noisy and fluctuating information coming from environment and response accordingly. In addition to the extracellular environment, the internal dynamics of cells are also found to be stochastic, such as their gene expressions [19].

Perkins *et al.* [15] argue that, as the extracellular environment is a stochastic process, the intracellular processing of these stochastic signals and choosing an appropriate response can only be probabilistic and consists of three main steps. First step is to infer from the noisy signals and make a prediction of the current or future state. Second is to make a proper choice of action considering the advantages and disadvantages, given the predictions. Final step is to take this action in a way that will contribute to the survival of the cell population.

Many studies took probabilistic approaches to explain the behaviour of cellular networks quantitatively and statistical inference is presented as a possible framework to explain the mechanisms that a cell may use to interpret the state of the environment from noisy signals. Libby *et al* [11] used Bayesian inference approach to model the gene expression of a bacterium in an environment with high and low level of sugar, and showed this model is consistent with the measurements. This approach is later extended for the situations where the signal fluctuates over time, e.g. non steady-state sugar concentration. Andrews *et al* [?] proposed that the cell adopts to such environment by updating its beliefs in real-time. They modelled this decision-making mechanism as a sequential application of Bayesian inference, where the posterior probability that is inferred in the current step is used as a prior probability in the next step.

Bosher *et al.* [1] took a similar approach to explain cellular decision making, focusing on information theory. They argued that mutual information between the signal and the response is shown to be a suitable measure to quantify the cell's ability to infer. It is argued that if the mutual information between the signal to be inferred and the output of a signal transduction mechanism is high, only then the cell could be able to perform a high quality inference. For example, Bialek *et al*, in their study of development of fruit fly, showed that the mutual information between gap gene expression and the position of nucleus is the information needed for each nucleus to determine their position in a cell.

The importance of information sharing in multi-agent systems are shown by blabla. This

paper tries to explain the effect of information sharing on the cooperative optimal strategy in a stochastically fluctuating environment. A coordination game of two players is presented in a two-state environment, where the players should coordinate on the two possible actions accordingly. The optimal strategy problem in this game is solved in two different communication scenarios. In the first game scenario  $G_p$ , the players only get individual signals from the environment, and do not communicate with each other. The communication with the environment is modelled as binary noisy channel with probability  $p$  of getting the true signal. In the second scenario  $G_{pq}$ , the players also share information about the cues they get from the environment, through another binary noisy channel with probability  $q$  of sending it correctly. With game-theoretic approach, the average long-term fitness of the players is defined as the fraction of time that the players cooperatively choose the correct action weighted by the game payoffs. The strategies maximizing the long-term fitness are considered optimal, and used to investigate the influence of information sharing on this cooperative game.

## 1.1. Motivation

## 1.2. Structure of the Thesis

**Chapter 2** presents the theoretical background of this work. It reviews the details of continuous-time Bayesian networks and partially observable Markov decision processes, and introduces the sampling algorithms used in this work.

**Chapter 3** is dedicated to the details of the problem. It explains how the frameworks are utilised, and presents the algorithms used for data generation and inference.

**Chapter 4** presents the experimental results of the simulation and the inference.

**Chapter 5** discusses the results, highlights the conclusions and reviews the limitations.

**Chapter 6** gives suggestions for future directions and extensions to the research.

## 2. Foundations

This chapter presents the theory applied in this thesis. First, the details of the communication problem are described briefly to put the theory into perspective, and then the mathematical background of the frameworks used to model this problem is introduced.

### 2.1. Problem Formulation

The communication model considered in this work is given in Figure 2.1. The parent nodes,  $X_1$  and  $X_2$ , emit messages which contain information about their states. These messages are translated by an observation model,  $\psi$ , and an agent node,  $X_3$  makes a decision based on this translated message,  $y$ . The main objective is to infer the observation model, given a set of trajectories of nodes.

The transition models of the nodes and the dependencies between them are modelled as a continuous-time Bayesian network (CTBN), denoted by  $S$ .

The messages that are emitted by the parent nodes  $X_1$  and  $X_2$  are modelled as independent homogeneous continuous-time Markov processes  $X_i(t)$ , with state space  $\mathcal{X}_i = \{x_1, x_2, \dots, x_m\}$  for  $i \in \{1, 2\}$ .

The agent node  $X_3$  does not have direct access to the messages but observes a translation of them. The observation model is defined as the likelihood of a translation given the parent messages.

$$\psi(x_1, x_2) := p(y(t) \mid X_1(t) = x_1, X_2(t) = x_2) \quad (2.1)$$

The agent  $X_3$  is modelled as inhomogeneous continuous-time Markov process with state space  $\mathcal{X}_3 = \{x_1, x_2, \dots, x_m\}$  and set of actions  $a \in \{a_0, a_1, \dots, a_k\}$  to choose from.

Given the observation, the agent forms a belief over the parent states,  $b(x_1, x_2; t)$ , that summarizes the past observations. The policy of the agent,  $\pi(a \mid b)$ , is assumed to be shaped by evolution (close) to optimality. Based on the belief state, the agent takes an action, which in the setting described above corresponds to changing its internal dynamics.

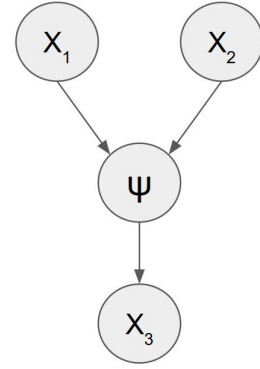


Figure 2.1.: Communication model.

## 2.2. Continuous-Time Bayesian Networks

Consider a directed acyclic graph denoted by  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V}$  is a set of nodes and  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$  is a set of edges such that  $\mathcal{E} = \{(m, n) : m, n \in \mathcal{V}\}$ . In this graph, the parent nodes of node  $n$  are defined as the set of nodes that feed into it and denoted by  $\text{Par}_{\mathcal{G}}(n) = \{m \in \mathcal{V} : (m, n) \in \mathcal{E}\}$ . A directed acyclic graph is characterized as a Bayesian network where each node represents a random variable such that  $\mathcal{V} = \{X_1, X_2, \dots, X_N\}$  and the joint distribution  $p(X_1, X_2, \dots, X_N)$  factors as

$$p(X_1, X_2, \dots, X_N) = \prod_{i=1}^N p(X_i \mid \text{Par}_{\mathcal{G}}(X_i)). \quad (2.2)$$

A continuous-time Bayesian network (CTBN) is a graphical model with graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  that represents a collection of random variables whose values evolve continuously over time. In the CTBN framework, the dependencies of a set of Markov processes (MPs) can be modelled efficiently through a directed graph, relying on two assumptions. The first assumption is that only one node can transition at a time, and the second is that the instantaneous dynamics of each node depends only on its parent nodes [3, 13]. A two component CTBN is illustrated in Figure 2.2.

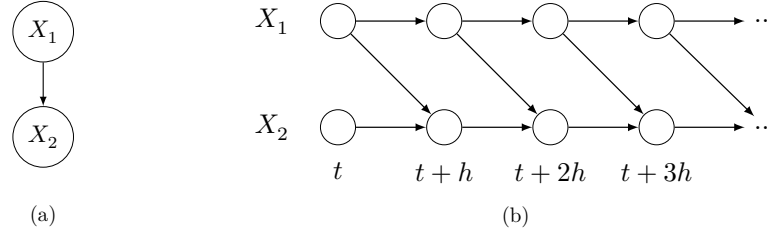


Figure 2.2.: Two component CTBN. (a) Network representation (b) CTBN in (a) unrolled in time as  $h \rightarrow 0$

### 2.2.1. Continuous-Time Markov Processes

A continuous-time Markov process (CTMP) is a continuous-time stochastic process which satisfies the Markov property, namely, the probability distribution over the states at a later time is conditionally independent of the past states given the current state [3].

Consider a CTMP  $X(t)$  over a single variable with a countable state space  $\mathcal{X}$ . Then the Markov property can be written as

$$\Pr(X(t_k) = x_{t_k} \mid X(t_{k-1}) = x_{t_{k-1}}, \dots, X(t_0) = x_{t_0}) = \Pr(X(t_k) = x_{t_k} \mid X(t_{k-1}) = x_{t_{k-1}}) \quad (2.3)$$

where  $X(t)$  denotes the state of the variable at time  $t$  such that  $X(t) = x_t \in \mathcal{X}, t \geq 0$ , and  $t_0 < t_1 < \dots < t_k$ .

A CTMP is represented by its transition intensity matrix,  $\mathbf{Q} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ . In this matrix, the intensity  $q_i$  represents the instantaneous probability of leaving state  $x_i$  and  $q_{i,j}$  represents the

instantaneous probability of switching from state  $x_i$  to  $x_j$ , where  $x_i, x_j \in \mathcal{X}$ .

$$\mathbf{Q} = \begin{bmatrix} -q_1 & q_{1,2} & \cdots & q_{1,n} \\ q_{2,1} & -q_2 & \cdots & q_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ q_{n,1} & q_{n,2} & \cdots & -q_n \end{bmatrix} \quad (2.4)$$

where  $q_i = \sum_{j \neq i} q_{i,j}$  [13].

### 2.2.1.1. Homogenous Continuous-Time Markov Processes

A continuous-time Markov process is time-homogenous when the transition intensities do not depend on time. Let  $X(t)$  be a homogenous CTMP, with a countable state space  $\mathcal{X}$  and transition intensity matrix  $\mathbf{Q}$ . Infinitesimal transition probability from state  $x_i$  to  $x_j$  in terms of the transition intensities  $q_{i,j}$  can be written as [3]:

$$p_{i,j}(h) = \delta_{ij} + q_{i,j}h + o(h) \quad (2.5)$$

where  $p_{i,j}(h) \equiv \Pr(X(t+h) = x_j \mid X(t) = x_i)$  are Markov transition functions,  $\delta_{i,j} = \delta(x_i, x_j)$  is Kronecker delta and  $o(\cdot)$  is a function decaying to zero faster than its argument such that  $\lim_{h \rightarrow 0} \frac{o(h)}{h} = 0$ .

The *Chapman-Kolmogorov- or master-equation* is then derived as follows:

$$\begin{aligned} p_j(t) &= \Pr(X(t) = x_j) \\ &= \sum_{\forall i} p_{i,j}(h) p_i(t-h) \\ \lim_{h \rightarrow 0} p_j(t) &= \lim_{h \rightarrow 0} \sum_{\forall i} [\delta_{ij} + q_{i,j}h + o(h)] p_i(t-h) \\ &= \lim_{h \rightarrow 0} p_j(t-h) + \lim_{h \rightarrow 0} h \sum_{\forall i} q_{i,j} p_i(t-h) \\ \lim_{h \rightarrow 0} \frac{p_j(t) - p_j(t-h)}{h} &= \lim_{h \rightarrow 0} \sum_{\forall i} q_{i,j} p_i(t-h) \\ \frac{d}{dt} p_j(t) &= \sum_{\forall i} q_{i,j} p_i(t) \end{aligned} \quad (2.6)$$

Equation 2.6 can be written in matrix form:

$$\frac{d}{dt} p(t) = p(t) \mathbf{Q} \quad (2.7)$$

where the time-dependent probability distribution  $p(t)$  is a row vector with entries  $\{p_i(t)\}_{x_i \in \mathcal{X}}$ . The solution of the system of ordinary differential equations (ODEs) is

$$p(t) = p(0) \exp(t\mathbf{Q}) \quad (2.8)$$

with initial distribution  $p(0)$ .

The amount of time staying in a state  $x_i$  is exponentially distributed with parameter  $q_i$ . The

probability density function  $f$  and cumulative distribution function  $F$  for staying in the state  $x_i$  can be written as [13]

$$f(t) = q_i \exp(-q_i t), t \geq 0 \quad (2.9)$$

$$F(t) = 1 - \exp(-q_i t), t \geq 0. \quad (2.10)$$

Given the transitioning from state  $x_i$ , the probability of landing on state  $x_j$  is  $q_{i,j}/q_i$ .

**Likelihood Function** Consider a single transition denoted as  $d = \langle x_i, x_j, t \rangle$ , where the transition occurs from state  $x_i$  to  $x_j$  after spending  $t$  amount of time at state  $x_i$ . The likelihood of this transition is the product of the probability of having remained at state  $x_i$  for duration  $t$  from Equation 2.9, and the probability of transitioning to  $x_j$ .

$$p(d \mid \mathbf{Q}) = (q_i \exp(-q_i t)) \left( \frac{q_{i,j}}{q_i} \right) \quad (2.11)$$

The likelihood of a trajectory sampled from a homogenous CTMC, denoted by  $X^{[0,T]}$ , can be decomposed as the product of the likelihood of single transitions. The sufficient statistics summarizing this trajectory can be written as  $T[x_i]$ , the total amount of time spent in state  $x_i$  and  $M[x_i, x_j]$  the total number of transitions from state  $x_i$  to  $x_j$ .

$$M[x_i, x_j] = \sum_{d \in X^{[0,T]}} \mathbb{1}(X(t) = x_i) \mathbb{1}(X(t+h) = x_j) \quad (2.12)$$

$$T[x_i] = \sum_{d \in X^{[0,T]}} \mathbb{1}(X(t) = x_i) \quad (2.13)$$

where  $\mathbb{1}(\cdot)$  is the indicator function. Then the likelihood of a trajectory  $X^{[0,T]}$  can be written as:

$$\begin{aligned} p(X^{[0,T]} \mid \mathbf{Q}) &= \prod_{d \in X^{[0,T]}} p(d \mid \mathbf{Q}) \\ &= \left( \prod_i q_i^{M[x_i]} \exp(-q_i T[x_i]) \right) \left( \prod_i \prod_{j \neq i} \left( \frac{q_{i,j}}{q_i} \right)^{M[x_i, x_j]} \right) \\ &= \prod_{j \neq i} \exp(-q_{i,j} T[x_i]) q_{i,j}^{M[x_i, x_j]} \end{aligned} \quad (2.14)$$

where  $M[x_i] = \sum_{j \neq i} M[x_i, x_j]$  is the total number transitions leaving state  $x_i$ .

### 2.2.1.2. Conditional Markov Processes

A continuous-time Markov process is *time-inhomogenous* when the transition intensities change over time. In a CTBN, while every node is a Markov process, the leaf nodes are characterized as *conditional* Markov processes, a type of inhomogeneous MP, where the intensities change over time, but not as a function of time rather as a function of parent states [13].

Let  $X$  be a conditional Markov process in a graph  $\mathcal{G}$ , with a set of parents  $U = \text{Par}_{\mathcal{G}}(X)$ .



Its *conditional intensity matrix* (CIM)  $\mathbf{Q}_{X|U}$  can be viewed as a set of homogenous intensity matrices  $\mathbf{Q}_{X|u}$ , with entries  $q_{i,j|u}$  (similar to Equation 2.4), for each instantiation of parent nodes  $U(t) = u$  such that  $u \in \mathcal{U} = \bigtimes_{X_m \in \text{Par}_{\mathcal{G}}(X)} \chi_m$ , where  $\bigtimes$  denotes Cartesian product [13]. As a result, given a trajectory of parent nodes,  $\mathbf{X}$  has a trajectory of intensity matrix as

$$\mathbf{Q}^{[0,T]} = [\mathbf{Q}_{X|U(t_0)}, \mathbf{Q}_{X|U(t_1)}, \dots, \mathbf{Q}_{X|U(t_N)}], \quad 0 < t_0 < \dots < t_N \leq T. \quad (2.15)$$

Markov transition function for a conditional MP can be written as

$$\Pr(X(t+h) = x_j \mid X(t) = x_i, U(t) = u, \mathbf{Q}_{X|u}) = \delta(i, j) + q_{i,j|u}h + o(h). \quad (2.16)$$

**Likelihood Function** Given the instantiation of its parents, the complete information on the dynamics of  $\mathbf{X}$  is obtained. Then the likelihood of a trajectory drawn from a conditional MP  $\mathbf{X}$  can be written similar to Equation 2.14,

$$\begin{aligned} p(X^{[0,T]} \mid \mathbf{Q}_{X|U}) &= \left( \prod_u \prod_i q_{i|u}^{M[x_i|u]} \exp(-q_{i|u}T[x_i \mid u]) \right) \left( \prod_u \prod_i \prod_{j \neq i} \left( \frac{q_{i,j|u}}{q_{i|u}} \right)^{M[x_i, x_j|u]} \right) \\ &= \prod_u \prod_i \prod_{j \neq i} \exp(-q_{i,j|u}T[x_i \mid u]) q_{i,j|u}^{M[x_i, x_j|u]} \end{aligned} \quad (2.17)$$

with the sufficient statistics,  $T[\cdot]$  and  $M[\cdot]$  introduced in Section 2.2.1.1, are also conditioned on parent nodes.

## 2.2.2. The CTBN Model

Evidently, a homogeneous CTMP can be considered as a conditional MP whose set of parents is empty. Thus, a CTBN can be formed as a set of conditional Markov processes.

Let  $S$  be a CTBN with graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  and local variables  $\mathcal{V} = \{X_1, X_2, \dots, X_N\}$ , each with a state space  $\chi_n$ . This results in factorizing state spaces for  $S$  such that  $\mathcal{S} = \chi_1 \times \chi_2 \times \dots \times \chi_N$ . The joint states of the variables are denoted by  $s = (x_1, x_2, \dots, x_N) \in \mathcal{S}$  where  $x_1 \in \chi_1, \dots, x_N \in \chi_N$ . The dependencies of each variable are defined as a set of its parents  $U_n = \text{Par}_{\mathcal{G}}(X_n)$  with values  $U_n(t) = u_n$  such that  $u_n \in \mathcal{U}_n = \bigtimes_{X_m \in \text{Par}_{\mathcal{G}}(X_n)} \chi_m$ . In the following, the set of all

conditional transition intensity matrices are denoted as  $\mathbf{Q} = \{Q_{X_1|U_1}, \dots, Q_{X_N|U_N}\}$ .

Consider a trajectory drawn from  $S$ , such that  $S^{[0,T]} = \{X_1^{[0,T]}, X_2^{[0,T]}, \dots, X_N^{[0,T]}\}$ . Following Equation 2.17, the likelihood of this trajectory can be written as

$$p(S \mid \mathbf{Q}) = \prod_{n=1}^N \prod_{u \in \mathcal{U}_n} \prod_{x_i \in \chi_n} \prod_{x_j \in \chi_n \setminus x_i} \exp(q_{i,j|u}^n T_n[x_i \mid u]) (q_{i,j|u}^n)^{M_n[x_i, x_j|u]} \quad (2.18)$$

where  $T_n[\cdot]$  and  $M_n[\cdot]$  indicates the sufficient statistics for  $X_n$ .

It should be noted that a CTBN can also be represented by a single conditional intensity matrix, through *amalgamation* operation [13].

## 2.3. Partially Observable Markov Decision Processes

The partially observable Markov decision process (POMDP) framework provides a model of an agent which interacts with its environment but is unable to obtain certain information about its state. Instead, the agent gets an observation which is a function of the true state, e.g. noisy observations, translation. The main goal, similar to Markov decision processes (MDPs), is to learn a policy solving a task by optimizing a reward function.

The problem of decision making under uncertainty can be considered in two parts for the agent. The first is to keep a belief state which summarizes past experiences, and the second is to optimize a policy  $\pi$  which will give an action based on the belief state [9, 12].

Consider a POMDP represented as a tuple  $(S, A, T, R, \mathcal{Y}, \psi)$ , where  $S$  is a countable set of states of the world,  $A$  is a set of actions,  $T : S \times A \rightarrow \Pi(S)$  is the state-transition function,  $R : S \times A \rightarrow \mathbb{R}$  is reward function,  $\mathcal{Y}$  is set of observations and  $\psi : S \times A \rightarrow \Pi(\mathcal{Y})$  is the observation function. The transition function  $T$  gives a probability distribution over states, given a state and action, such that  $T(s, a, s') = \Pr(s' \mid s, a)$ . The observation function  $\psi$  gives a probability distribution over observations given a state and action, such that  $\psi(s', a, y) = \Pr(y \mid s', a)$ . The reward function  $R$  gives the expected immediate reward for each action and state  $R(s, a)$ . The belief state, if represented as a probability distribution over states, provides a summary over the agent's past experiences. This representation also allows the agent to account for its uncertainty while making decisions. The optimal policy of a POMDP agent leads to optimal behaviour as a function of the agent's belief state. The expected amount of future rewards upon executing a policy is defined as *value function* and used to evaluate a given policy.  $V_\pi(s)$  denotes the expected sum of reward obtained by following policy  $\pi$ , starting from state  $s$ .

Finite-horizon model represents a problem where the agent has  $t$  steps to take, that is given the current state of the world, the agent will receive an observation and make a decision  $t$  times. In this model, it is more likely that the agent has a nonstationary policy, a policy which is a sequence of policies indexed by time. That corresponds to a time-dependent behaviour. Consider the famous and well-studied tiger problem as an example. In this problem, the agent is in front of two doors, one of which has a tiger behind, the other has a reward. The agent has three actions to choose from, it can open the left door or the right door, or it can listen to obtain information. If the agent has only one step to take, then even though it is not completely certain about the location of the tiger, it might go for opening the door. However, when it has two steps left, it might be wiser to listen and obtain more information about the state of the world.

A nonstationary  $t$ -step policy can be represented as a *policy tree* shown in Figure 2.3. The policy tree describes the optimal behaviour for  $t$  steps conditioned on the observation. The top node shows the first action to be taken, then depending on the observation, a different branch is followed until the end of the steps.

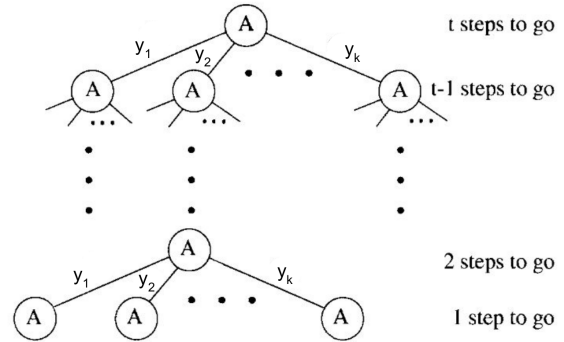


Figure 2.3.: A policy tree of depth  $t$  representing  $t$ -step nonstationary policy. Source: [9]

The value of a nonstationary t-step tree is given in Equation 2.19-Equation 2.21, where  $s$  is starting state,  $a(p)$  is the first action at the top of the tree,  $s'$  is the next state, and  $y_i(p)$  is the (t-1)-step policy that has been chosen after taking the action  $a(p)$  and getting observation  $y_i$ .  $\gamma$  is the *discount factor*, such that  $0 \leq \gamma \leq 1$ , used to regulate the contribution of the future rewards to value function [21].

$$V_p(s) = R(s, a(p)) + \gamma \cdot (\text{Expected value of the future}) \quad (2.19)$$

$$= R(s, a(p)) + \gamma \sum_{s' \in \mathcal{S}} \Pr(s' \mid s, a(p)) \sum_{y_i \in \mathcal{Y}} \Pr(y_i \mid s', a(p)) V_{y_i(p)}(s') \quad (2.20)$$

$$= R(s, a(p)) + \gamma \sum_{s' \in \mathcal{S}} T(s, a(p), s') \sum_{y_i \in \mathcal{Y}} \psi(s', a(p), y_i) V_{y_i(p)}(s') \quad (2.21)$$

The value of a policy tree  $p$  starting from state  $s$  has two components, as can be seen in Equation 2.19. The first component is the immediate reward, that is, the reward that the agent will get for taking action  $a(p)$  while at state  $s$ . The second component is the expected future reward, discounted by  $\gamma$ . This value is computed first by taking the expectation over next possible states  $s' \in \mathcal{S}$ , and their values. The value of a state depends on the observation that the agent receives, which will determine the (t-1)-step policy to be executed. Therefore, a second expectation is performed over the observations.

As mentioned before, the agent is never certain about the state of the world. Therefore, the relevant value function is that of a policy tree starting from a belief state  $b$ , and calculated as the expected value over the states.

$$V_p(b) = \sum_{s \in \mathcal{S}} b(s) V_p(s) \quad (2.22)$$

Equation 2.22 gives the value of a policy tree starting from belief state  $b$ . The optimal policy, then, is chosen as the one which has the maximum value. The optimal t-step value starting from belief state  $b$  is executing the best policy tree for that belief state.

$$V_t(b) = \max_{p \in \mathcal{P}} V_p(b) \quad (2.23)$$

where  $\mathcal{P}$  denotes a finite set of policy trees. It is noteworthy that the value function of every policy tree  $V_p$  is linear in  $b$ . As can be seen from Equation 2.23,  $V_t$  is defined as the maximum of all  $V_p$  over  $b$ . It is the envelope of these value functions; therefore, it is piecewise-linear and convex.

Infinite-horizon discounted model consider the value function over an infinitely long trajectory of the agent. In a POMDP problem, the infinite horizon discounted value function is still convex [24], but it may not always be piecewise-linear. However, it can be approximated by a finite horizon value function for sufficiently many steps [18, 5].

In the problem considered in this thesis, the agent node  $X_3$  cannot observe the incoming messages directly, rather a summary of them. This setting presents a POMDP problem. However, since the optimal policy of the agent is assumed to be given, the main focus in the POMDP framework is belief state estimation.

In the following, update methods for the belief state are introduced, where belief state refers to the posterior probability distribution over the environment states.

### 2.3.1. Exact Belief State Update

In a scenario where compact representations of the *transition model*,  $T(s, a, s')$ , and *observation model*,  $\psi(s', a, y)$ , are available, the belief state update can be obtained as [9]

$$b'(s') = \Pr(s' | y, a, b) \quad (2.24)$$

$$= \frac{\Pr(y | s', a, b) \Pr(s' | a, b)}{\Pr(y | a, b)} \quad (2.25)$$

$$= \frac{\Pr(y | s', a) \sum_{s \in \mathcal{S}} \Pr(s' | a, b, s) \Pr(s | a, b)}{\Pr(y | a, b)} \quad (2.26)$$

$$= \frac{\psi(s', a, y) \sum_{s \in \mathcal{S}} T(s, a, s') b(s)}{\Pr(y | a, b)}. \quad (2.27)$$

The relation between transition model  $T$  and transition intensity matrix  $\mathbf{Q}$  can be written as  $T = \exp(t\mathbf{Q})$  from Equation 2.8. For the derivation above, from Equation 2.24 to Equation 2.25, Bayes' theorem, and from Equation 2.25 to Equation 2.26, law of total probability is applied. It is also noteworthy that the denominator in Equation 2.27 is in the following form,

$$\Pr(y | a, b) = \sum_{\forall s' \in \mathcal{S}} \psi(y | s', a) \sum_{\forall s \in \mathcal{S}} T(s' | s, a) b(s) \quad (2.28)$$

which is computationally expensive in the case of continuous state space.

### 2.3.2. Exact Belief State Update for CTMP by Filtering

Equation 2.27 is the discrete-time solution of belief state. However, since in the model described in Section 2.1, the parent nodes are modelled as CTMPs, thus the environment state for the agent is the state of a CTMP, the belief state should be solved in continuous-time. This is achieved by the inference of the posterior probability of CTMP [8].

A *filtering problem* in statistical context refers to the inference of the conditional probability of the true state of the system at a point in time, given the history of observations [7].

Let  $X$  be a CTMP with transition intensity matrix  $\mathbf{Q}$ . Assume discrete-time observations denoted by  $y_1 = y(t_1), \dots, y_N = y(t_N)$ . The belief state can be written as:

$$b(x_i; t_N) = \Pr(X(t_N) = x_i | y_1, \dots, y_N) \quad (2.29)$$

From the master equation given in Equation 2.6, it follows that:

$$\frac{d}{dt} b(x_j; t) = \sum_{\forall i} q_{i,j} \cdot b(x_i; t) \quad (2.30)$$

The time-dependent belief state  $b(t)$  is a row vector with  $\{b(x_i; t)\}_{x_i \in \mathcal{X}}$ . This posterior probability can be described by a system of ODEs:

$$\frac{db(t)}{dt} = b(t)\mathbf{Q} \quad (2.31)$$

where the initial condition  $b(0)$  is a row vector with  $\{b(x_i; t)\}_{x_i \in \mathcal{X}}$  [8]. The solution to this ODE is

$$b(t) = b(0) \exp(t\mathbf{Q}). \quad (2.32)$$

The belief state update at discrete times of observation  $y_t$  is derived as

$$\begin{aligned} b(x_i; t_N) &= \Pr(X(t_N) = x_i, | y_1, \dots, y_N) \\ &= \frac{\Pr(y_1, \dots, y_N, X(t_N) = x_i)}{\Pr(y_1, \dots, y_N)} \\ &= \frac{\Pr(y_N | y_1, \dots, y_{N-1}, X(t_N) = x_i)}{\Pr(y_N | y_1, \dots, y_{N-1})} \frac{\Pr(y_1, \dots, y_{N-1}, X(t_N) = x_i)}{\Pr(y_1, \dots, y_{N-1})} \\ &= Z_N^{-1} \Pr(y_N | X(t_N) = x_i) \Pr(X(t_N) = x_i | y_1, \dots, y_{N-1}) \\ &= Z_N^{-1} p(y_N | x_i) b(x_i; t_N^-) \end{aligned} \quad (2.33)$$

where  $Z_N = \sum_{x_i \in \mathcal{X}} p(y_N | x_i) b(x_i; t_N^-)$  is the normalization factor [8].

### 2.3.3. Belief State Update using Particle Filter

In a more realistic scenario, the exact update of belief state may not be feasible for several reasons. The computation of exact belief update is expensive for large state spaces, which can be seen from Equation 2.28. Moreover, a problem with continuous state spaces requires a belief state represented as probability distributions over an infinite state space rather than a collection of probabilities as given in Section 2.3.1 [23]. Such representation cannot be obtained using the exact method. Another reason could be the lack of a compact representation of transition or observation models. Under such circumstances, the belief state is obtained using sample-based approximation methods [23].

It should be noted that since the belief state is nothing but the conditional probability of true states given the observations, the problem at hand poses a filtering problem as described in Section 2.3.2.

#### 2.3.3.1. Particle Filtering

Particle filtering is one of the most commonly used Sequential Monte Carlo (SMC) algorithms. The popularity of this method thrives from the fact that, unlike other approximation methods such as Kalman Filter, it does not assume a linear Gaussian model. This advantage offers great flexibility and finds application in a wide range of areas [4].

The key idea in particle filtering is to approximate a target distribution  $p(x)$  by a set of samples, i.e. particles, drawn from that distribution. This is achieved by sequentially updating the particles through two steps. The first step is *importance sampling*. Since the target distribution is not available, the particles are generated from a *proposal distribution*  $q(x)$  and weighted according to the difference between target and proposal distributions. The second step is to resample the particles using these weights with replacement [7].

Consider a problem of deriving the expectation  $\hat{f}(x) = \mathbb{E}[f(x)] = \int f(x)p(x)dx$ , and suppose

$p(x)$  is an intractable density function from which the particles cannot be sampled. Instead they are drawn from a proposal distribution  $q(x)$ , which yields an empirical approximation such that

$$x^{(i)} \sim q(x)$$

$$q(x) \approx \frac{1}{N} \sum_{i=1}^N \delta_{x^{(i)}}(x)$$

where  $\delta_{x^{(i)}}(x)$  is Dirac delta. The expectation in question can be written as

$$\int f(x)p(x)dx = \int f(x)\frac{p(x)}{q(x)}q(x)dx$$

$$\int f(x)\frac{p(x)}{q(x)}\left(\frac{1}{N}\sum_{i=1}^N\delta_{x^{(i)}}(x)\right)dx = \frac{1}{N}\sum_{i=1}^N\frac{p(x^{(i)})}{q(x^{(i)})}f(x^{(i)})$$

where  $w(x^{(i)}) = \frac{p(x^{(i)})}{q(x^{(i)})}$  is defined as *importance weight* of a particle. Then the particles are resampled using the importance weights with replacement, which concludes one iteration of sequential updating [7].

### 2.3.3.2. Marginalized Continuous-Time Markov Processes

In this work, the particles to represent the belief state are drawn from a marginalized CTBN. Consider a CTBN,  $S$ , with local variables  $X_n$ ,  $n \in \{1, \dots, N\}$ , and a set of conditional intensity matrices  $\mathbf{Q}$ . In the following, it is assumed that every non-diagonal entry in  $\mathbf{Q}_{n|u}$  is Gamma distributed with shape and rate parameters,  $\alpha_{i,j|u}^n$  and  $\beta_{i,j|u}^n$ .

The marginal process description of  $S$  considering a single trajectory in the interval  $[0, t)$  can be written as

$$\Pr(X_n(t+h) = x_j \mid X_n(t) = x_i, U_n(t) = u, S^{[0,t]}) \quad (2.34)$$

$$= \int \Pr(X_n(t+h) = x_j \mid X_n(t) = x_i, U_n(t) = u, Q_{n|u}, S^{[0,t]}) p(Q_{n|u}) dQ_{n|u} \quad (2.35)$$

$$= \delta_{i,j} + \mathbb{E}[q_{i,j|u}^n \mid S^{[0,t]} = s^{[0,t]}] h + o(h). \quad (2.36)$$

By integrating out the intensity matrix  $Q_{n|u}$ , the parameter is replaced by its expected value given the history of the process. It should be noted that by doing so, the process becomes parameter-free, and thus self-exciting [20].

The derivation of the conditional expectation for a marginal CTBN follows from the Bayes' rule:

$$p(\mathbf{Q} \mid S^{[0,t]}) = \frac{p(S^{[0,t]} \mid \mathbf{Q}) p(\mathbf{Q})}{p(S^{[0,t]})} \quad (2.37)$$

Equation 2.37, written for single trajectory  $S^{[0,t]}$ , can be extended for multiple trajectories. Consider  $K$  trajectories drawn from  $S$ , denoted by  $\xi_t = \{S^{[0,t],1}, S^{[0,t],2}, \dots, S^{[0,t],K}\}$ . Since the

trajectories are conditionally independent, given  $\mathbf{Q}$ , using Equation 2.18 the likelihood of set  $\xi_t$  is written as,

$$\Pr(\xi_t | \mathbf{Q}) = \prod_{n=1}^N \prod_{u \in \mathcal{U}_n} \prod_{x_i \in \mathcal{X}_n} \prod_{x_j \in \mathcal{X}_n \setminus x_i} \exp(q_{i,j|u}^n T_n[x_i | u]) (q_{i,j|u}^n)^{M_n[x_i, x_j | u]} \quad (2.38)$$

where the joint sufficient statistics of  $X_n$  over all K trajectories are denoted by  $T_n[x_i | u] = \sum_{k=1}^K T_n^k[x_i | u]$  and  $M_n[x_i, x_j | u] = \sum_{k=1}^K M_n^k[x_i, x_j | u]$ . Given independent Gamma-priors on transition intensities, the expectation in Equation 2.36 can be evaluated as follows:

$$\mathbb{E}[q_{i,j|u}^n | \xi_t] = \frac{\alpha_{i,j|u}^n + M_n[x_i, x_j | u]}{\beta_{i,j|u}^n + T_n[x_i | u]} \quad (2.39)$$

The algorithm for belief state update through marginal particle filtering is given in the following chapter.

## 2.4. Sampling Algorithms

### 2.4.1. Gillespie Algorithm for Generative CTBN

Gillespie algorithm is a computer-oriented Monte Carlo simulation procedure that is originally proposed to simulate the reactions of molecules in any spatially homogeneous chemical system. Such systems are regarded as Markov processes and represented via their master equations, which cannot be directly used to obtain realizations of the process. Gillespie algorithm is an efficient tool to overcome this problem [6].

This algorithm can also be applied to sample *events* from a CTBN given the transition intensity matrices, where an event refers to a transition occurring at a specific point in time. This procedure is introduced as *Generative CTBN* in [13].

---

#### Algorithm 1: Generative CTBN

---

**Input** : Structure of the network with N local variables  $X_1, X_2, \dots, X_n$  with state-space  $\mathcal{X}_n = \{x_1, \dots, x_m\}$   
Transition intensity matrices  $\mathbf{Q}_n$  with entries  $q_{i,j}^n$   
 $T_{max}$  to terminate simulation

**Output** : Sample trajectory of the network

**Initialize:** Initialize node values  $X_n(0) = x_i \in \mathcal{X}_n$

- 1: **while**  $t < T_{max}$  **do**
- 2:  $\tau \sim \exp(\sum_{\forall n} \sum_{\forall i \neq j} q_{i,j}^n)$
- 3: transitioning node is randomly drawn with probability  $P(X_n) = \frac{q_i^n}{\sum_{\forall n} q_i^n}$
- 4: next state is randomly drawn with probability  $P(x_j) = \frac{q_{i,j}^n}{q_i^n}$
- 5:  $t \leftarrow t + \tau$
- 6: **end while**

---

### 2.4.2. Thinning Algorithm

Thinning algorithm is a method introduced to simulate nonhomogenous Poisson processes [10]. Later, it is adapted to sample from Hawkes processes, a self-exciting process with time-dependent intensity function [14, 16]. This algorithm is used here to simulate the inhomogeneous Markov process.

---

**Algorithm 2:** Thinning Algorithm

---

**Input** :  $\lambda(t)$  the intensity function of the inhomogenous process  
           $N$  number of events to terminate simulation

**Output** : Sample trajectory of the process

**Initialize:** Time  $t = 0$

```
1: while  $i < N$  do
2:   the upper bound for intensity,  $\lambda^*$ 
3:   transition time  $\tau$  drawn by  $u \sim U(0, 1)$  and  $\tau = \frac{-\ln(u)}{\lambda^*}$ 
4:    $t \leftarrow t + \tau$ 
5:   draw  $s \sim U(0, 1)$ 
6:   if  $s \leq \frac{\lambda(t)}{\lambda^*}$  then
7:     sample accepted and  $t_i = t, i = i + 1$ 
8:   end if
9: end while
```

---



## 3. Methodology

This chapter presents the methodology used in this thesis. First, it is explained how different frameworks, which were introduced in Chapter 2, are put into use. Then, the algorithms used in data generation and inference are given in detail. The results from these experiments are presented in the succeeding chapter.

### 3.1. The Model

A detailed graphical model explored in this thesis is given in the Figure 3.1. This model presents an intersection of continuous-time Bayesian network and partially observable Markov decision process frameworks.

- The transition models of the nodes  $X_1, X_2$  and  $X_3$ , and the dependencies between them are modelled as CTBN.
- The interaction of agent node  $X_3$  and its environment is modelled as POMDP.

#### 3.1.1. CTBN Model

The transition models of the nodes and the dependencies between them are modelled as continuous-time Bayesian network (CTBN), denoted by  $S$  with graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V} = \{X_1, X_2, X_3\}$  and  $\mathcal{E} = \{(X_1, X_3), (X_2, X_3)\}$ . The network  $S$  represents a stochastic process over a structured factorising state space  $\mathcal{S} = \chi_1 \times \chi_2 \times \chi_3$ .

The parent nodes  $X_1$  and  $X_2$  emit their states as messages. The dynamics of these nodes are modelled as independent homogeneous continuous-time Markov processes  $X_i(t)$ , with binary-valued states  $\chi_i = \{0, 1\}$  for  $i \in \{1, 2\}$ . These processes are defined by transition intensity matrices  $\mathbf{Q}_i$ , which are assumed to be Gamma distributed with shape and rate parameters  $\boldsymbol{\alpha} = [\alpha_0, \alpha_1]$  and  $\boldsymbol{\beta} = [\beta_0, \beta_1]$ ,

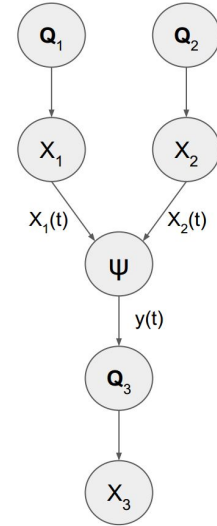


Figure 3.1.: Hierarchical model.

respectively, and are in the following forms.

$$\mathbf{Q}_i = \begin{bmatrix} -q_0^i & q_0^i \\ q_1^i & -q_1^i \end{bmatrix} \quad (3.1)$$

$$\mathbf{Q}_i \sim \text{Gam}(\boldsymbol{\alpha}^i, \boldsymbol{\beta}^i) \text{ for } i \in \{1, 2\} \quad (3.2)$$

It should be noted that in Equation 3.1, the suffixes are simplified using the fact that  $q_i = \sum_{i \neq j} q_{i,j}$ .

The agent  $X_3$  is modelled as inhomogenous continuous-time Markov process with binary states  $\chi_3 = \{0, 1\}$  and set of actions  $a \in \{a_0, a_1\}$ , and set of transition intensity matrices which contains one matrix corresponding to each action,  $\mathbf{Q}_{3|a} = \{\mathbf{Q}_{3|a_0}, \mathbf{Q}_{3|a_1}\}$ .

The dependencies are represented by set of parents for each node  $U_n = \text{Par}_{\mathcal{G}}(X_n)$  and for the model shown in Figure 3.1 can be written as follows:

$$\begin{aligned} U_1, U_2 &= \emptyset \\ U_3 &= \{X_1, X_2\} \end{aligned}$$

In order to have a compact representation of parent messages, a subsystem of  $S$  consisting of only the parent nodes,  $X_1$  and  $X_2$  can be considered as a single system. These two processes can be represented as a *joint* process,  $X_P$ , with factorising state space  $\chi_P = \chi_1 \times \chi_2$ . The transition intensity matrix of the new joint system,  $\mathbf{Q}_P$  is obtained by amalgamation operation, denoted by  $*$ , between  $\mathbf{Q}_1$  and  $\mathbf{Q}_2$  (see Appendix A) [13].

$$\mathbf{Q}_P = \mathbf{Q}_1 * \mathbf{Q}_2 \quad (3.3)$$

### 3.1.2. POMDP Model

In a conventional POMDP scenario, there are two problems to be addressed, one is belief state update and the other is policy optimization. As mentioned in Section 2.3, in the problem at hand, the policy of agent  $X_3$  is assumed to be optimal and given. Thus, the POMDP model of the agent only consists of belief state update. A detailed view of the agents interaction with its environment from POMDP framework perspective is given in the Figure 3.2.

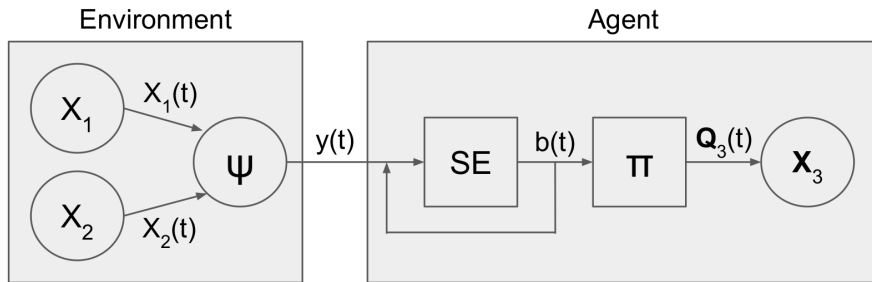


Figure 3.2.: Closer look to agent-environment interaction from the perspective of POMDP framework.

It should be noted that, the interaction in Figure 3.2 is only one-sided, the state or action of the agent does not affect the environment.

### 3.1.2.1. Observation Model

The messages sent by the parent nodes are translated by the observation model. The agent node  $X_3$  does not have a direct access to the messages, but observes a translation of them. The observation is denoted by  $y(t) = y_t$  such that  $y_t \in \mathcal{Y}$  where  $\mathcal{Y}$  is the observation space. The observation model defines a probability distribution over the observation for each combination of parent messages.

$$\psi(x_1, x_2) = p(y(t) \mid X_1(t) = x_1, X_2(t) = x_2) \quad (3.4)$$

where  $x_1 \in \mathcal{X}_1$  and  $x_2 \in \mathcal{X}_2$ . As explained in Section 3.1.1, using the joint process  $X_P$  for the sake of conciseness, Equation 3.4 can be written as

$$\psi(x_P) = p(y(t) \mid X_P(t) = x_P) \quad (3.5)$$

where  $x_P \in \mathcal{X}_P$ .

$\psi(x_P)$  is defined as deterministic categorical distribution over the observation space  $\mathcal{Y}$ . For each state  $x_P$ , there is one possible observation  $y_i \in \mathcal{Y}$ , such that

$$p(y_i \mid x_P) = 1 \wedge p(y_j \mid x_P) = 0 \quad \forall j \neq i. \quad (3.6)$$

$\psi$  denotes the matrix with rows  $\{\psi(x_P)\}_{x_P \in \mathcal{X}_P}$ .

### 3.1.2.2. Belief State

The belief state provides a summary over agents past experiences and allows the agent to take its own uncertainty into account. The belief state is formed by the *state estimator* (labelled as *SE* in Figure 3.2) over the parent states, denoted by  $b(x_P; t)$ .

$$b(x_P; t) = \Pr(X_P(t) = x_P \mid y_1, \dots, y_t) \quad (3.7)$$

**Exact Belief State Update** As discussed in Section 2.3.2, given the transition intensity matrices of parent nodes,  $\mathbf{Q}_1$  and  $\mathbf{Q}_2$ , the continuous-time belief state update poses a filtering problem for CTMPs. This problem can be formulated according to the joint process of parents.

$$b(x_P; t) = \Pr(X_P(t) = x_P \mid y_1, \dots, y_t) \quad (3.8)$$

Consider discrete-time observations from this process, denoted by  $y_1 = y(t_1), \dots, y_l = y(t_l)$  and time-dependent belief state  $b(t)$  as a row vector with  $\{b(x_P; t)\}_{x_P \in \mathcal{X}_P}$ . Following Equation 2.32 and Equation 2.33, the belief state update is evaluated as

$$b(t) = b(0) \exp(t\mathbf{Q}_P) \quad (3.9)$$

with the initial condition  $b(0)$ . The update at discrete times of observation  $y_t$  is

$$b(x_P; t_l) = Z_l^{-1} p(y_l | X_P(t_l) = x_P) b(x_P; t_l^-) \quad (3.10)$$

$$= Z_l^{-1} \psi(x_P) b(x_P; t_l^-) \quad (3.11)$$

where  $Z_l = \sum_{x_P \in \mathcal{X}_P} \psi(x_P) b(x_P; t_l^-)$  is the normalization factor.

**Belief State Update Using Marginalized Particle Filter** The assumption that full information of parent dynamics being available is unrealistic. In an environment as described above, the agent is more likely not to have access to the parameters  $\mathbf{Q}_1$  and  $\mathbf{Q}_2$ , may rather have some prior beliefs over them. Moreover, when the state estimator utilizes exact update method, these parameters are assumed to be available for the inference as well. Thus, in order to simulate a more realistic model and be able to marginalize out these parameters from inference problem, the joint parent process  $X_P$  is replaced with its marginalized counterpart. Using the Gamma-priors over  $\mathbf{Q}_1$  and  $\mathbf{Q}_2$  (Equation 3.2) and sufficient statistics over the particle history, the particles are drawn from this marginalized process as explained in Section 2.3.3.2. With every new observation, the particles are propagated through the marginal process, while the sufficient statistics are updated and the parameters are re-estimated after each particle using the Equation 2.39. The belief state then obtained as the distribution of states over the particles,

$$b(x_P; t) = \frac{1}{N} \sum_{i=1}^N \delta_{k_i(t), x_P} \quad (3.12)$$

where  $N$  is the number of particles,  $k_i \in \mathbf{k}$  is the set of particles, and  $\delta$  is the Kronecker delta.

---

**Algorithm 3:** Marginal particle filter for belief state update [20]

---

**Input** : Observation  $y_l$  at time  $t_l$ , set of particles  $\mathbf{k}^{l-1}$ , estimated  $\hat{Q}$

**Output:** New set of particles  $\mathbf{k}^l$ ,  $\mathbf{b}^{[t_{l-1}, t_l]}$

---

```

1: for  $k_m \in \mathbf{k}^{l-1}$  do
2:    $k_m = \{x_m, \hat{Q}\} \leftarrow \text{Propagate particle through marginal process from } t_{l-1} \text{ to } t_l$ 
3:    $\hat{Q} \leftarrow \text{sufficient statistics added from } k_m[t_{l-1}, t_l]$ 
   // observation likelihood assigned as particle weight
4:    $w_m \leftarrow p(y_l | X_P(t_l) = x_m)$ 
5: end for
   // belief state from  $t_{l-1}$  to  $t_l$ 
6:  $\mathbf{b}^{[t_{l-1}, t_l]} \leftarrow \left\{ \frac{1}{N} \sum_{i=1}^N \delta_{k_i^{[t_{l-1}, t_l]}, x_P} \right\}_{x_P \in \mathcal{X}_P}$ 
   // normalize weights
7:  $w_m \leftarrow \frac{w_m}{\sum_m w_m}$ 
   // resample particles
8: for  $k_m \in \mathbf{k}^l$  do
9:    $k_m \leftarrow \text{Sample from } \mathbf{k}^l \text{ with probabilities } w_m \text{ with replacement}$ 
10: end for
```

---

In Algorithm 3, the weight update for the particles is performed on line 4, based on the

observation model. Given the deterministic nature of the observation model as described in Equation 3.6, in rare cases where the observation  $y_k$  stems from an unlikely transition of parent nodes, the particles fail to simulate this transition and all of them get rejected. One possible solution to this problem would be increasing the number of particle to increase the probability of sampling unlikely transitions. However, in practice, this solution is not feasible due to computational cost. Instead, the degeneration of the particle filter is dealt with by assigning uniform probabilities to the particles, effectively ignoring the unlikely transitions. The situation is illustrated by examples from simulation in Section 4.2.

### 3.1.2.3. Optimal Policy

The optimal policy is defined using a polynomial function of belief state.

$$\pi(b) = \begin{cases} a_0 & \text{if } \mathbf{w}b^\top > 0.5 \\ a_1 & \text{otherwise} \end{cases} \quad (3.13)$$

where  $\mathbf{w}$  is a row vector of weights.

Given the optimal policy,  $\pi(b)$ , the agent takes an action based on the belief state. In the setting described above, taking an action means to change its internal dynamics to the transition intensity matrix corresponding to that action.

$$a(t) = \pi(b(t)) \quad (3.14)$$

$$\mathbf{Q}_3(t) = \begin{cases} \mathbf{Q}_{3|a_0} & \text{if } a(t) = a_0 \\ \mathbf{Q}_{3|a_1} & \text{otherwise} \end{cases} \quad (3.15)$$

## 3.2. Inference of Observation Model

Inference problem is considered for deterministic observation models, such that each state  $x_P \in \mathcal{X}_P$  can only be translated to one observation. Considering the number of states of parents and the observations, this results in a number of possible observation models.

Consider a trajectory in the dataset, denoted by  $S^{[0,T]} = \{X_1^{[0,T]}, X_2^{[0,T]}, X_3^{[0,T]}\}$ . The set of parameters to the system, as introduced before, is written as  $\theta = \{\mathbf{Q}_1, \mathbf{Q}_2, \mathbf{Q}_3, \pi, \psi\}$ . Given the parent trajectories,  $X_1^{[0,T]}$  and  $X_2^{[0,T]}$ , the belief state and the resulting  $\mathbf{Q}_3$  trajectory is computed for each observation model. Then the likelihood of  $S^{[0,T]}$  trajectory given the parameters  $\theta$  are compared for maximum likelihood estimation.

$$\hat{\psi} = \arg \max_{\psi} p(S^{[0,T]} | \theta) \quad (3.16)$$

### 3.2.1. Likelihood Model

The likelihood of a sample trajectory  $S^{[0,T]}$  can be written as:

$$\begin{aligned}
p(S^{[0,T]} | \theta) &= p(X_1^{[0,T]}, X_2^{[0,T]}, X_3^{[0,T]} | \mathbf{Q}_1, \mathbf{Q}_2, \mathbf{Q}_3, \pi, \psi) \\
&= p(X_3^{[0,T]} | X_1^{[0,T]}, X_2^{[0,T]}, \mathbf{Q}_1, \mathbf{Q}_2, \mathbf{Q}_3, \pi, \psi) p(X_1^{[0,T]} | \mathbf{Q}_1) p(X_2^{[0,T]} | \mathbf{Q}_2) \\
&= p(X_3^{[0,T]} | X_1^{[0,T]}, X_2^{[0,T]}, \mathbf{Q}_3, \pi, \psi) p(X_1^{[0,T]} | \mathbf{Q}_1) p(X_2^{[0,T]} | \mathbf{Q}_2) \\
&= p(X_3^{[0,T]} | \mathbf{Q}_3^{[0,T]}) p(X_1^{[0,T]} | \mathbf{Q}_1) p(X_2^{[0,T]} | \mathbf{Q}_2)
\end{aligned} \tag{3.17}$$

As mentioned before, it is plausible to marginalize out the parameters  $\mathbf{Q}_1$  and  $\mathbf{Q}_2$ , for a more realistic model and inference. Noting that in case the belief state is updated using filtering of CTMPs (See Section 3.1.2.2),  $\mathbf{Q}_3^{[0,T]}$  becomes a deterministic function of all the parameters including  $\mathbf{Q}_1$  and  $\mathbf{Q}_2$ , the marginalization cannot be carried out analytically on Equation 3.17. On the other hand, marginal particle filtering removes this dependency on  $\mathbf{Q}_1$  and  $\mathbf{Q}_2$  by using marginalized counterpart of CTMPs (See Section 3.1.2.2), leaving it straightforward to marginalize out the parameters on Equation 3.17.

Marginalizing the likelihood over  $\mathbf{Q}_1$  and  $\mathbf{Q}_2$ :

$$\begin{aligned}
p(S^{[0,T]} | \pi, \psi) &= \int \int p(S^{[0,T]} | \theta) p(\mathbf{Q}_1) p(\mathbf{Q}_2) d\mathbf{Q}_1 d\mathbf{Q}_2 \\
&= \int \int p(X_3^{[0,T]} | \mathbf{Q}_3^{[0,T]}) p(X_1^{[0,T]} | \mathbf{Q}_1) p(X_2^{[0,T]} | \mathbf{Q}_2) p(\mathbf{Q}_1) p(\mathbf{Q}_2) d\mathbf{Q}_1 d\mathbf{Q}_2 \\
&= p(X_3^{[0,T]} | \mathbf{Q}_3^{[0,T]}) \int p(X_1^{[0,T]} | \mathbf{Q}_1) p(\mathbf{Q}_1) d\mathbf{Q}_1 \int p(X_2^{[0,T]} | \mathbf{Q}_2) p(\mathbf{Q}_2) d\mathbf{Q}_2
\end{aligned} \tag{3.18}$$

Marginalized likelihood function for binary-valued homogenous CTMP is derived in Appendix B.

Plugging Equation B.3 in Equation 3.18 for both  $X_1$  and  $X_2$ :

$$\begin{aligned}
p(S^{[0,T]} | \pi, \Phi) &= p(X_3^{[0,T]} | \mathbf{Q}_3^{[0,T]}) \prod_{x_1 \in \{0,1\}} \frac{\beta_{x_1}^{\alpha_{x_1}}}{\Gamma(\alpha_{x_1})} (T_{x_1} + \beta_{x_1})^{M_{x_1} + \alpha_{x_1}} \Gamma(M_{x_1} + \alpha_{x_1}) \\
&\quad \prod_{x_2 \in \{0,1\}} \frac{\beta_{x_2}^{\alpha_{x_2}}}{\Gamma(\alpha_{x_2})} (T_{x_2} + \beta_{x_2})^{M_{x_2} + \alpha_{x_2}} \Gamma(M_{x_2} + \alpha_{x_2})
\end{aligned} \tag{3.19}$$

In order to avoid numerical instability, the log-likelihood is preferred for the calculations instead of likelihood.

### 3.2.2. Inference under Noisy Observation Model

In order to assess the robustness of the inference, we added an error probability  $p_e$  to the observation model. As explained in Section 3.1.2.1, the observation model  $\psi(x_P)$  is assumed to be deterministic. This corresponds to a unique translation, denoted here by  $y_{\text{true}}$  of each

parent state  $x_P$ . In the case of noisy observation model, the resulting observation might differ from the correct translation  $y_{\text{true}}$  with probability  $p_e$ , and the erroneous translation is drawn uniformly from the remaining observation space  $y_i \in \mathcal{Y} \setminus \{y_{\text{true}}\}$ . This noisy observation model is denoted by  $\psi^{p_e}$  can be considered as a noisy communication channel with error probability  $p_e$ .

$$p(y_{\text{true}} | x_P) = 1 - p_e \wedge p(y_j | x_P) = \frac{p_e}{|\mathcal{Y}| - 1} \quad \forall y_j \neq y_{\text{true}} \quad (3.20)$$

For simulation, the noisy observation model is assumed to be available to the agent, which mainly affects the agent's belief state. For exact belief state update, the noisy observation model is employed in Equation 3.11. In the case of particle filtering, the weights are assigned to the particle using the noisy model (see Algorithm 3, line 4). It is noteworthy that by doing so, the degeneration of particle filter as described in Section 3.1.2.2 is prevented.

### 3.3. Data Generation

The dataset contains a number of trajectories drawn from CTBN S. Following the notation in Chapter 2,  $K$  trajectories in time interval  $[0, T]$  are denoted by  $\xi_T = \{S^{[0,T],1}, S^{[0,T],2}, \dots, S^{[0,T],K}\}$ , where  $S^{[0,T],k} = \{X_1^{[0,T],k}, X_2^{[0,T],k}, X_3^{[0,T],k}\}$  denotes a single trajectory for all nodes. Every trajectory comprises of state transitions in the interval, and the times of these transitions.

#### 3.3.1. Sampling Algorithm

In order to sample trajectories from CTBN, two sampling algorithms introduced in Section 2.4 are combined. Gillespie algorithm is used to sample from the parent nodes,  $X_1$  and  $X_2$ , while thinning algorithm is applied to overcome the challenges that come with conditional intensity matrix of the agent,  $X_3$ . It should be noted that Algorithm 1 is applicable to any nodes in a CTBN, both homogenous and conditional MPs. However, since in this setting, the intensity matrix is conditioned on the belief state and the policy, instead of directly on the parent states, a more general algorithm suitable for inhomogenous MPs, thinning algorithm, is preferred. Algorithm 4 describes the procedure to draw samples using marginal particle

filtering.

---

**Algorithm 4:** Sampling trajectories with marginal particle filtering

---

**Input** : Gamma-prior parameters on parents' transition intensity matrices  
 $\alpha^1, \beta^1, \alpha^2, \beta^2$   
Set of agent's transition intensity matrices  $\mathbf{Q}_3$   
 $T_{max}$  to terminate simulation

**Output** : Sample trajectory of the network

**Initialize:** Sample  $\mathbf{Q}_1$  and  $\mathbf{Q}_2$  from their priors  
Initialize nodes uniformly  $X_n(0) = x_i \in \chi_n$   
Initialize particles uniformly  $p^i(0) = x_p \in \chi_P$   
 $t = 0$

- 1: **while**  $t < T_{max}$  **do**
- 2: Draw next transition for  $X_1$  and  $X_2$  ( $\tau_{parent}$ ,  $x_1$  and  $x_2$  using Algorithm 1)
- 3:  $t_{parent} \leftarrow t + \tau_{parent}$  // transition time for parents
- 4:  $y_{t_{parent}} \sim \psi(x_1, x_2)$  // new observation at  $t_{parent}$
- 5: Update particle filter and obtain  $\mathbf{b}^{[t, t_{parent}]}$
- 6:  $a^{[t, t_{parent}]} \leftarrow \pi(\mathbf{b}^{[t, t_{parent}]})$
- 7:  $\mathbf{Q}_3^{[t, t_{parent}]} \leftarrow \mathbf{Q}_3|_{a^{[t, t_{parent}]}}$
- 8:  $t_{agent} \leftarrow t$
- 9: **while**  $t_{agent} < t_{parent}$  **do**
- 10: the upper bound for intensity,  $q_3^*$ <sup>1</sup>
- 11: transition time  $\tau_{agent}$  drawn by  $u \sim U(0, 1)$  and  $\tau_{agent} = \frac{-\ln(u)}{q_3^*}$
- 12:  $t_{agent} \leftarrow t_{agent} + \tau_{agent}$
- 13: draw  $s \sim U(0, 1)$ , accept transition if  $s \leq \frac{q_3(t_{agent})}{q_3^*}$
- 14: **end while**
- 15:  $t \leftarrow t_{parent}$
- 16: **end while**

---



---

<sup>1</sup> $q$  is the transition intensity associated with the current state of the agent.



## 4. Results

The experimental results are presented in this chapter. First, the parameters for the variables introduced in Chapter 3 are given. Then a sample of simulated trajectories are shown as an example. Finally, the inference results are presented.

### 4.1. Configurations

The configurations given below are used for the results presented in the following sections, if not specified otherwise.

- Gamma priors for parent dynamics such that  $\mathbf{Q}_i \sim \text{Gam}(\boldsymbol{\alpha}^i, \boldsymbol{\beta}^i)$  for  $i \in \{1, 2\}$ , and  $\boldsymbol{\alpha} = [\alpha_0, \alpha_1]$  and  $\boldsymbol{\beta} = [\beta_0, \beta_1]$

$$\boldsymbol{\alpha}^1 = [5, 10] \quad \boldsymbol{\beta}^1 = [5, 20] \quad (4.1)$$

$$\boldsymbol{\alpha}^2 = [10, 10] \quad \boldsymbol{\beta}^2 = [10, 5] \quad (4.2)$$

- Transition intensity matrices of  $X_1$  and  $X_2$  sampled from priors given above

$$\mathbf{Q}_1 = \begin{bmatrix} -1.117 & 1.117 \\ 0.836 & -0.836 \end{bmatrix} \quad (4.3)$$

$$\mathbf{Q}_2 = \begin{bmatrix} -1.1 & 1.1 \\ 2.445 & -2.445 \end{bmatrix} \quad (4.4)$$

- Length of trajectories  $T = 5\text{s}$
- State space,  $\chi_P = \chi_1 \times \chi_2 = \{(x_1, x_2)\}_{x_1 \in \chi_1, x_2 \in \chi_2} = \{(0, 0), (0, 1), (1, 0), (1, 1)\}$
- Observation space,  $\mathcal{Y} = \{0, 1, 2\}$
- Action space,  $A = \{a_0, a_1\} = \{0, 1\}$
- The set of transition intensity matrices of  $X_3$

$$\mathbf{Q}_3 = \left\{ \mathbf{Q}_{3|a_0}, \mathbf{Q}_{3|a_1} \right\} = \left\{ \begin{bmatrix} -0.5 & 0.5 \\ 2 & -2 \end{bmatrix}, \begin{bmatrix} -3 & 3 \\ 0.02 & -0.02 \end{bmatrix} \right\} \quad (4.5)$$

- Number of particles,  $N = 200$
- Weights of the policy introduced in Equation 3.13,  $\mathbf{w} = [0.02, 0.833, 0.778, 0.87]$

- Observation model  $\psi_{\text{true}} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

## 4.2. Simulation

The synthetic dataset is generated utilising Algorithm 4.  $K$  trajectories in time interval  $[0, T]$  are denoted by  $\xi_T = \{S^{[0,T],1}, S^{[0,T],2}, \dots, S^{[0,T],K}\}$ , where  $S^{[0,T],k} = \{X_1^{[0,T],k}, X_2^{[0,T],k}, X_3^{[0,T],k}\}$  denotes a single trajectory for all nodes. It is noteworthy that the initial states are drawn from discrete uniform distribution.

$$X_i(0) \sim \mathcal{U}\{0, 1\} \text{ for } i \in \{1, 2, 3\} \quad (4.6)$$

Figure 4.1(a)-(b) shows an example of parent trajectories. In Figure 4.1(c), the resulting trajectory of the joint parent process  $X_P$  is illustrated. As mentioned in Section 3.1.1, this joint process over the parent nodes provides a compact representation.

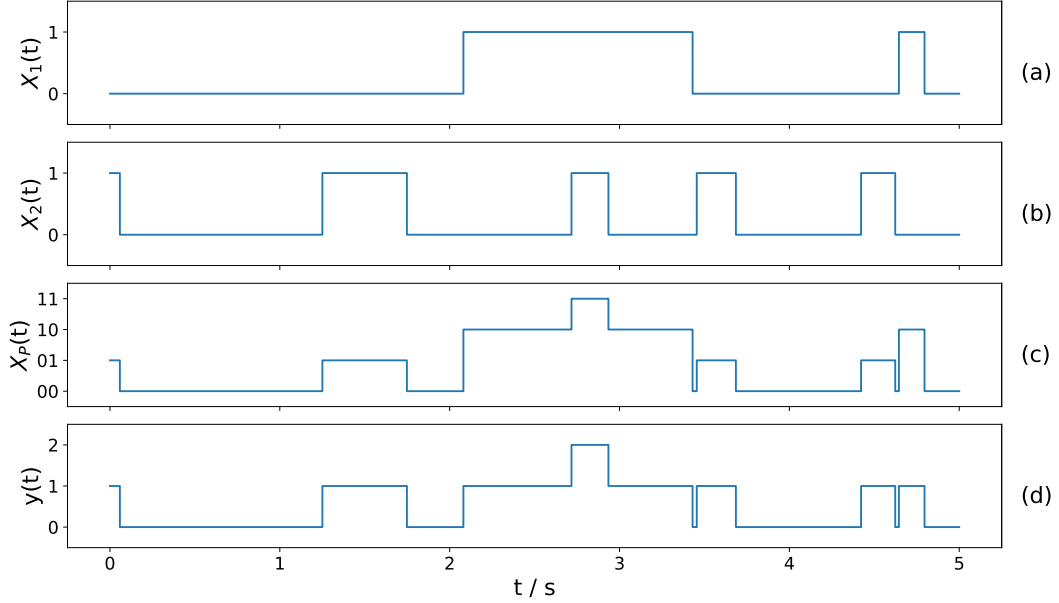


Figure 4.1.: A sample of parent trajectories and observation. (a)-(b) A sample trajectory of parent nodes  $X_1$  and  $X_2$  of length  $T = 5\text{s}$ , (c) The trajectory of the joint parent process  $X_P$ , (d) The observation trajectory resulting from  $X_P$  given in (c) and  $\psi_{\text{true}}$  given in Section 4.1.

In Figure 4.1(c), the states of  $X_P$  taking values in  $\chi_P = \chi_1 \times \chi_2$  is preferred to be represented as a combination of the parent states for readability, so that  $\chi_P = \{00, 01, 10, 11\}$ , where  $x_P \in \chi_P$  simply corresponds to  $x_1 x_2$ ,  $x_1 \in \chi_1$ ,  $x_2 \in \chi_2$ . Figure 4.1(d) shows the observation trajectory resulting from  $X_P(t)$  and the observation model  $\psi_{true}$  given in Section 4.1.

Figure 4.2 illustrates the belief state trajectory given the observations in Figure 4.1(d). For the reference, the belief state update using marginal particle filter and exact update is given together. As can be seen from the figures, the exact update is well approximated by the particle filter.

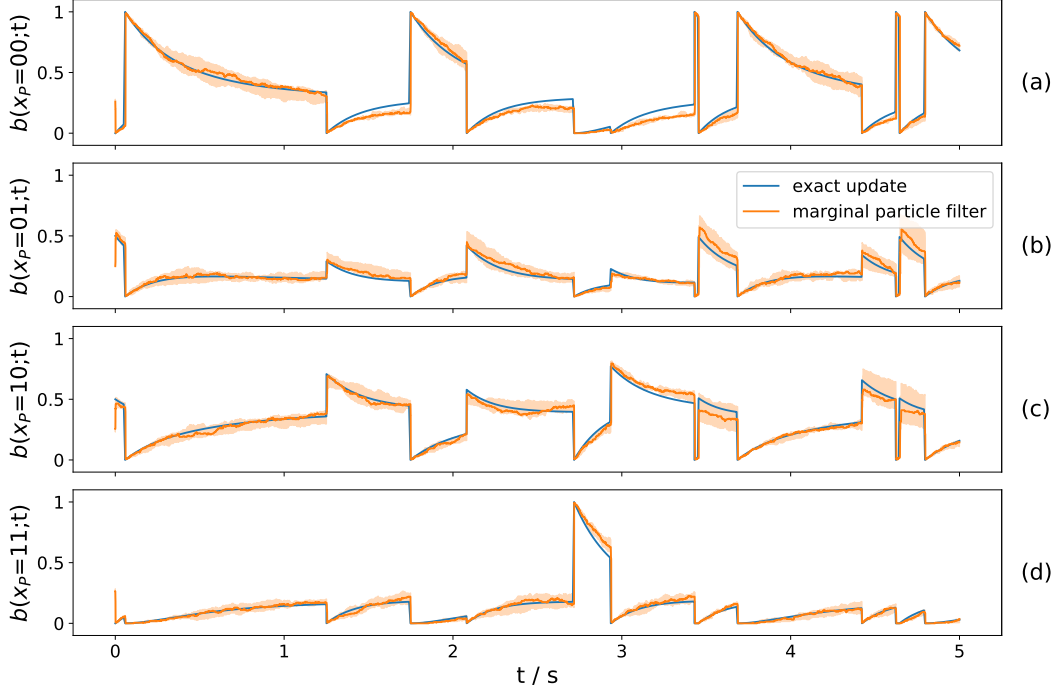


Figure 4.2.: Belief state trajectories corresponding to the observations given in Figure 4.1(d), comparing exact update method and marginal particle filtering.

Finally, the resulting  $Q_3$  and the trajectory of agent are given in Figure 4.3. The  $Q_3$  trajectory shown in Figure 4.3(b), are derived from the belief state update by marginal particle filter which is given in Figure 4.3(a) using Equation 3.13 and Equation 3.15.

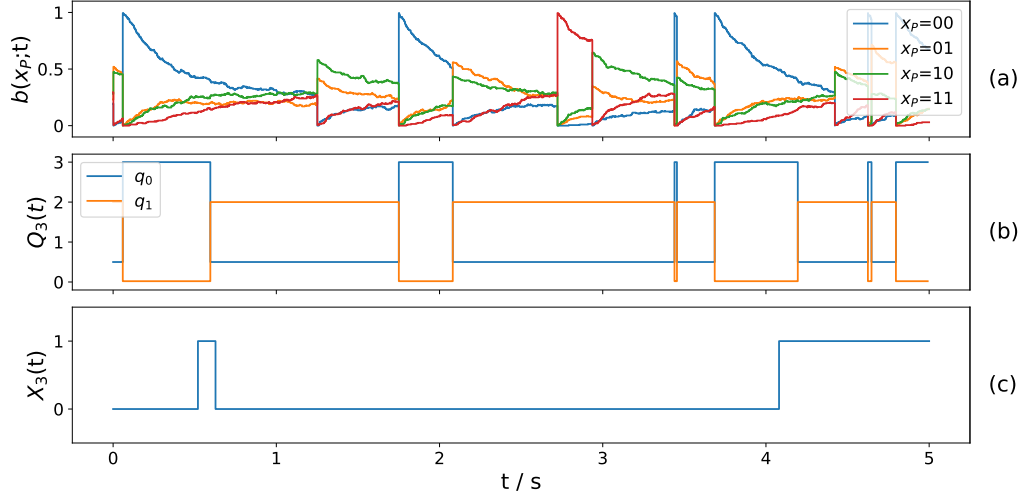


Figure 4.3.: Belief state updated by marginal particle filter and the resulting  $Q_3$  and  $X_3$  trajectories.

As mentioned in Section 3.1.2.2, the degeneration of the particle filter in case of an unlikely changes in observation, i.e. unanticipated transitions of parent nodes, is handled by by assigning uniform probabilities to the particles. It effectively corresponds to ignoring the rapid changes of the observation, which may cause diverging from the exact update, however, it is recovered with the next observation. Figure 4.4 provides an example of the situation. The rapid change in question that caused particle filter method to diverge from exact update method is highlighted in Figure 4.4(a). The particles fail to simulate this observation, and due to uniformization, the transition from 2 to 1 is ignored by the marginal particle filter method. The divergence can be observed in Figure 4.4(d)-(e) clearly.

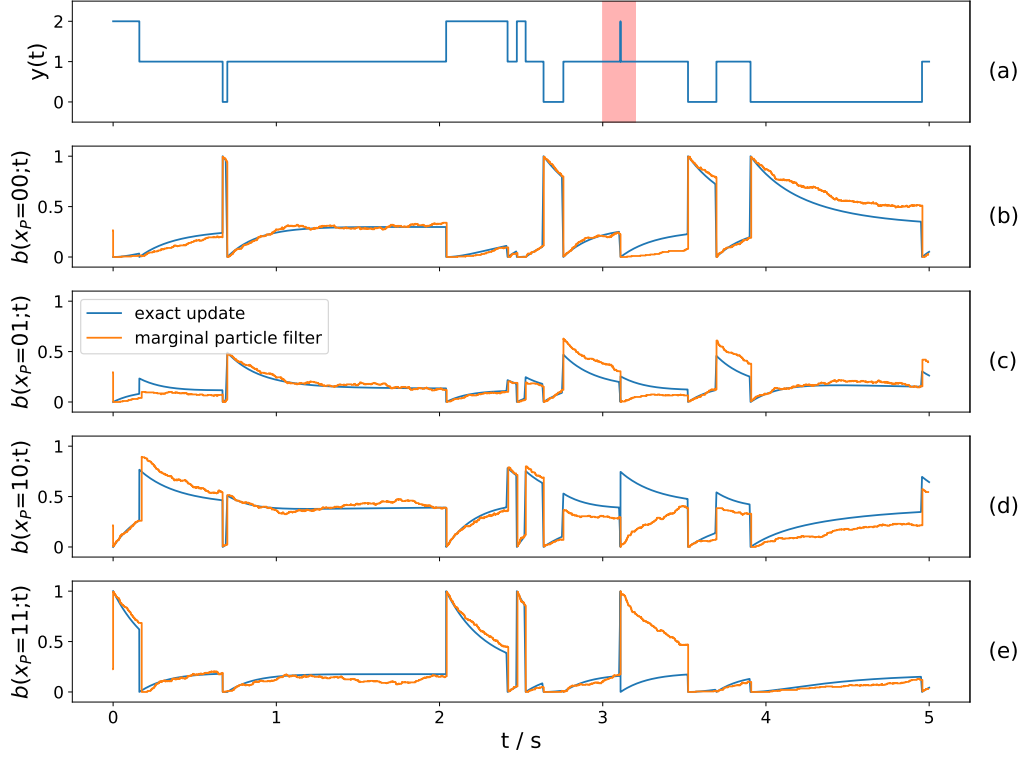


Figure 4.4.: A sample with degenerate marginal particle filter. The unlikely observation which has caused the degeneration is highlighted in (a). It can be seen in (d) and (e) that this observation causes particle filter approximation to diverge from exact update results but it is recovered with the next observation.

### 4.3. Inference of Observation Model

We considered the problem of inferring the observation model as a classification problem. As a measure of the performance of the classifier, we utilised area under the Receiver-Operator-Characteristic curve (AUROC) and Precision-Recall curve (AUPR).

#### 4.3.1. Equivalence Classes

As mentioned in Section 3.2, the deterministic nature of the observation model results in a number of possible observation models. The setting described in Chapter 3 with configurations given in Section 4.1 leads to 81 observation models. However, with this experimental setup and the methods, it is only possible to distinguish these observation models into 10 different classes. Due to this equivalence, the inference problem is considered only for 10 observation models, each one representing one class. The reasons of this phenomena are discussed in detail in

Appendix C, together with the observation models considered in the inference problem. The set of observation model that can be classified is denoted as  $\psi$  in the following.

Figure 4.5 illustrates the equivalence of observation models clearly. The plot depicts the results of an experiment with 200 samples,  $|\xi_T| = 200$ , generated using the observation model  $\psi_{\text{true}}$  given in Section 4.1, and the average log-likelihood of samples computed for all possible observation models. Here, the belief state is updated using exact method as described in Section 3.1.2.2, in order to depict the exact equivalence within one class. As can be seen, the results show the separation of the set of observation models into 10 distinct classes. The legend is removed to avoid clutter.

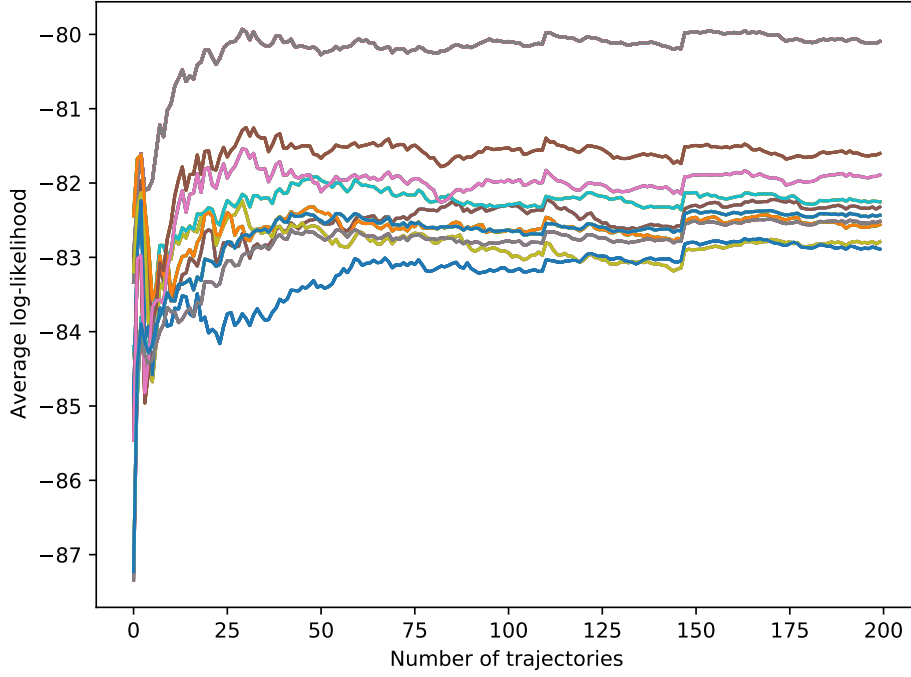


Figure 4.5.: Average log-likelihood  $\log p(S^{[0,T]} | \theta)$  over samples generated using exact belief state update, depicting the equivalence classes in the set of observation models.

In order to show the validity of the equivalence in the case of marginal particle filter, the average log-likelihoods of 200 samples given two observation models in the same class are illustrated in Figure 4.6. The samples are generated with  $\psi_{\text{true}} = \psi_0$ , and the rest of the observation models here fall in the same equivalence class as  $\psi_0$ . As can be seen from the graph, the observation models lead to so similar results to each other that they are assumed to be identical.

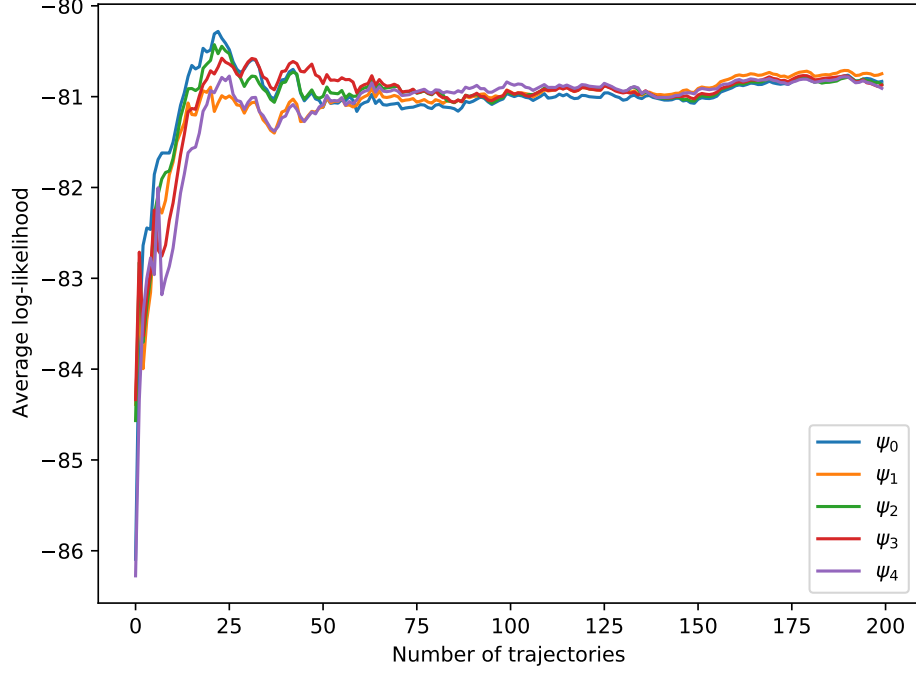


Figure 4.6.: Average log-likelihood  $\log p(S^{[0,T]} | \theta)$  over samples generated using marginal particle filtering, where  $\psi_1$  and  $\psi_2$  belongs in the same class.

#### 4.3.2. Learning Observation Model

Figure 4.7 illustrates the average log-likelihood over 100 samples given the observation models. The samples are generated with  $\psi_{true} = \psi_0$  given in Section 4.1, and exact update method is utilised for belief state update. As can be seen, the curves converge quickly and the true model is well separated from others. Consequently, the maximum likelihood estimation by Equation 3.16 leads to the correct result.

A similar experiment is documented in Figure 4.8, where exact update method is replaced by marginal particle filtering.  $\psi_0$  denotes the observation model that has generated the dataset, and it is correctly estimated as the true model. The jump around 50 trajectories can be explained by an encounter with a highly likely parent trajectories.

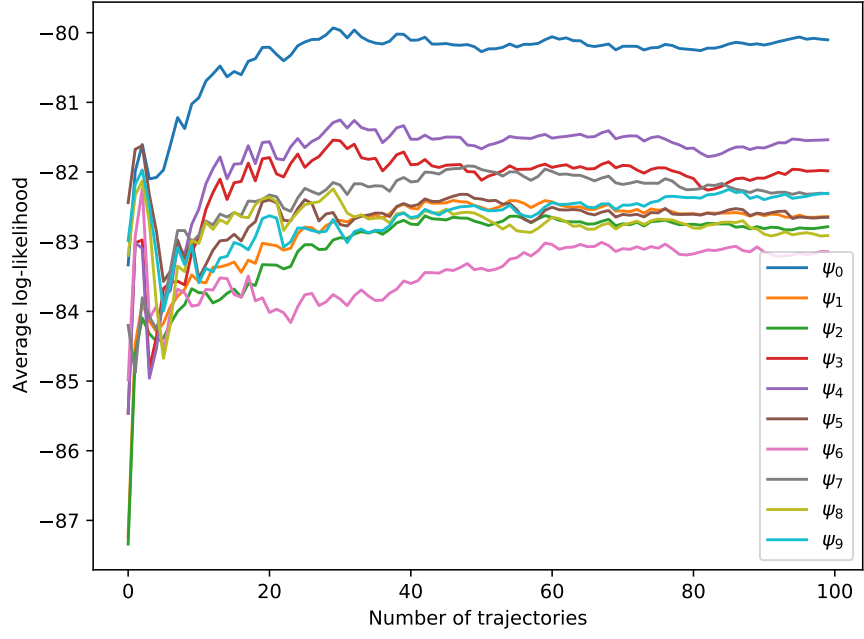


Figure 4.7.: Average log-likelihood  $\log p(S^{[0,T]} | \theta)$  with  $\psi_i \in \psi$  over samples generated using exact belief state update

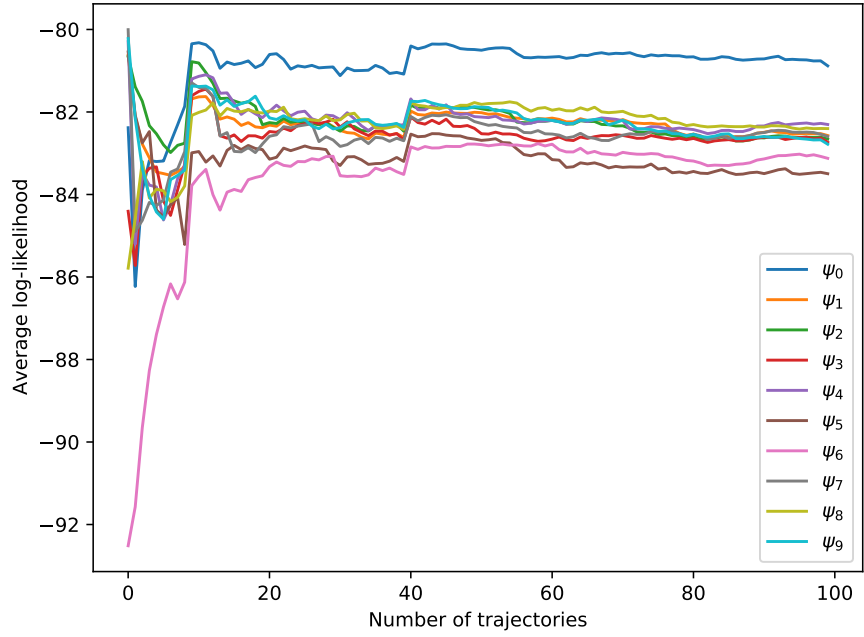


Figure 4.8.: Average log-likelihood  $\log p(S^{[0,T]} | \theta)$  with  $\psi_i \in \psi$  over samples generated using marginal particle filtering



We approach the inference problem as classification between the observation models  $\psi \in \Psi$  given in Appendix C.1. We consider the estimated likelihood values of each sample given an observation model as the score of the sample belonging to the corresponding class. Since this setting represents a multi-class classification problem, the performance metrics AUROC and AUPR are calculated as one-vs-rest.

Consider a binary classification problem. True positives (TP) are the samples predicted as 1 correctly, and false positives (FP) are the samples predicted as 1 while the true label was 0. True negatives (TN) are the samples predicted as 0 correctly, and false negatives (FN) are the samples with true label 1, but predicted 0. True positive rate (TPR), also called *recall* (R) or *sensitivity*, is the ratio of TPs over the number of samples which are labelled as 1. False positive rate (FPR) is the ratio of FPs over the number of samples labelled as 0. The precision (P) is the ratio of TPs over the number of all the samples predicted as 1.

$$\text{TPR} = \text{R} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (4.7)$$

$$\text{FPR} = \frac{\text{FP}}{\text{TN} + \text{FP}} \quad (4.8)$$

$$\text{P} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (4.9)$$

Receiver Operating Characteristics (ROC) curve illustrates the tradeoff between true positives and false positives, over different values of threshold for classification [17]. The area under ROC curve (AUROC) is a performance metric that shows how well the classifier can distinguish the classes. Higher AUROC metric indicates better performance, taking values in the interval  $[0,1]$ . Precision-Recall (PR) curve illustrates the relation between precision and recall, providing a metric to quantify how many of the predictions were correct [2].

We provided the classifier with increasing number of samples for inference. This is achieved through bootstrapping a given number of trajectories, and using the mean likelihood over the bootstrap batch as a new sample. The following AUROC plots shows the results over 50 trajectory generated using each observation model as the true model. According to this, in our dataset, we have 500 trajectories, 50 from each class labelled through a vector with 10 entries, having only 1 for the true observation model and 0 for the rest. When number of trajectories is 1, each sample in the dataset considered individually. When number of trajectories is 2, within each class, 50 sample batches of size 2 are bootstrapped such that none of the batches consist of the same samples. By doing so, we keep the sample size at 50 per each class, regardless of the batch size.

Figure 4.9 shows the AUROC results over 10 runs, comparing the state estimator using exact update and marginal particle filtering. We plotted the median AUROC as a line and 25-75% percentile as the shaded area. Figure 4.10 illustrates the AUPR in the same manner. As expected given the unbiased classifier, both metrics approach to 1 as the number of samples increases. Due to the stochasticity introduced by the marginal particle filtering as the state estimator, the results obtained with this method show slightly lower performance.

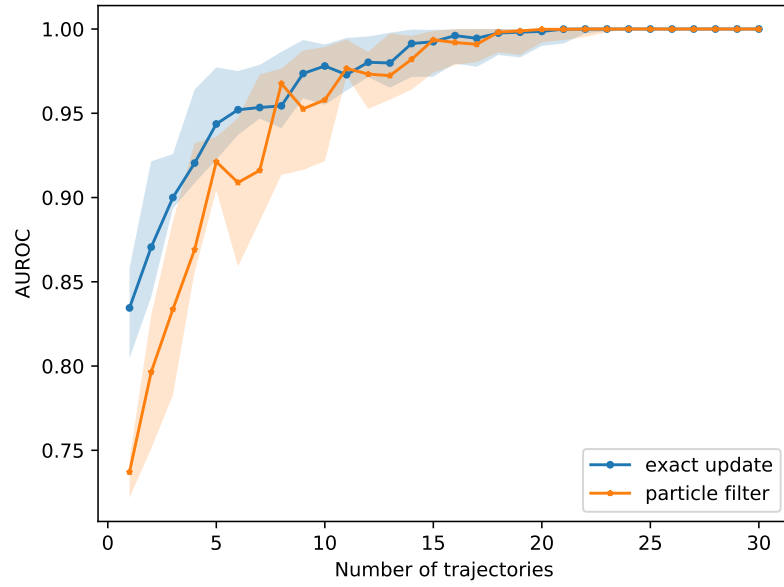


Figure 4.9.: AUROC results over increasing number of samples for  $\psi_0$ -vs-rest. We plotted the median with line and the 25-75% percentile with the shaded area over 10 runs.

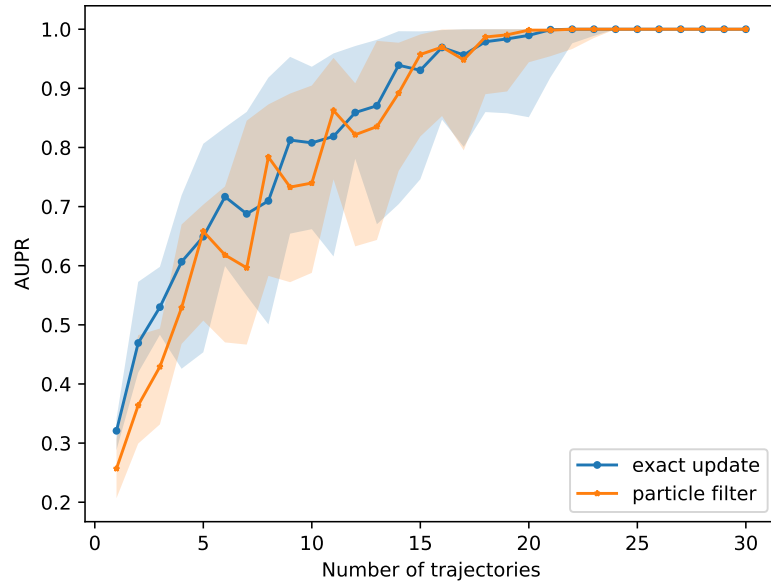


Figure 4.10.: AUPR results over increasing number of samples for  $\psi_0$ -vs-rest. We plotted the median with line and the 25-75% percentile with the shaded area over 10 runs.

#### **4.3.3. Inference with Non-informative Priors on Parent Parameters**

#### **4.3.4. Robustness Test under Channel Noise**

## 5. Discussion

In this work, a communication between two parent nodes and an agent node is modelled combining CTBN and POMDP frameworks. While the parent nodes emit messages containing information about their states, the agent observes a translation of these messages from which it needs to form its belief state and make decisions. The nodes evolve as components of a CTBN, modelled as in Section 3.1.1. Given that the messages of the parent nodes are unavailable to the agent, the interaction between parent nodes and the agent node is modelled as POMDP, as described in Section 3.1.2.

The belief state is updated utilising two methods. The first one is exact update, discussed in Section 3.1.2.2, and assumes that the transition intensities of the parents  $\mathbf{Q}_1$  and  $\mathbf{Q}_2$  are available both for the agent and for the classifier. However, due to the fact that this would not present a realistic system, particle filtering with marginalized CTBN is introduced for as state estimator. Here, both the agent and the classifier was able to perform the belief state update, given Gamma-priors of  $\mathbf{Q}_1$  and  $\mathbf{Q}_2$ .

The performance of the classifier is analysed in terms of the metrics AUROC and AUPR, and the results are given in Figure 4.9 and Figure 4.10, respectively. Using exact method to update the belief state, excellent performance is achieved for the classification task regarding AUROC metric. The stochasticity introduced by the particle filtering results in a slightly lower performance, compared to exact update. Nevertheless, in both cases, as the number of samples increases, the curve approaches 1, which is expected in the case of unbiased classifier.

The equivalence classes are introduced in Section 4.3.1. The set of observation models can be divided into 10 equivalence classes such that the likelihoods of a sample trajectory  $S^{[0,T]}$  given any observation model within one class are equal. This clearly limits the classifiers ability to determine the true model. Due to this limitation, set of observation models is reduced to 10 observation model, each one representing an equivalence class. These observation models are given in Appendix C.1. Consequently, the result which states that the true observation model is  $\psi_i$  is equivalent to that the true observation model belongs to  $i^{\text{th}}$  equivalence class. As shown with examples in Appendix C, there are two reason to this equivalence. The first reason can be specified as the equivalent effect of observation models on the belief state, which is inherent to the observation model structure. The second reason is that the different belief states might lead to the same behaviour. This case intuitively can be explained by the fact that the agent does not need to use all the information it has for decision making.

## 6. Outlook

The first step in the future of this work is to eliminate the equivalence classes to be able to classify every observation model. This problem to a certain extent can be mitigated by joint inference of observation model and policy. The joint inference could be performed as a joint classification problem, where the combination of discrete values of these parameters are treated as classes. This is only feasible by defining appropriate constraints on the policy such that, as for the case of observation models in this work, the policy space is countable.

Another exciting direction is to solve the policy optimization problem, instead of assuming that the optimal policy is given. By doing so, it would be feasible to utilise the classification method for observation model described in this work, in real-world data, providing insights in the interactions of agents and environments.

# Bibliography

- [1] Clive G. Bowsher and Peter S. Swain. Environmental sensing, information transfer, and cellular decision-making. *Current Opinion in Biotechnology*, 28:149–155, 2014.
- [2] Kendrick Boyd, Kevin H. Eng, and C. David Page. Area under the precision-recall curve: Point estimates and confidence intervals. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 8190 LNAI(PART 3):451–466, 2013.
- [3] Ido Cohn, Tal El-Hay, Nir Friedman, and Raz Kupferman. Mean Field Variational Approximation for Continuous-Time Bayesian Networks. *Journal of Machine Learning Research*, 11:2745–2783, 2010.
- [4] Arnaud Doucet and A M Johansen. A tutorial on particle filtering and smoothing: Fifteen years later. *Handbook of Nonlinear Filtering*, (December):4–6, 2009.
- [5] Sondik Edward. The Optimal Control of Partially Observable Markov Processes Over the Infinite Horizon : Discounted Costs Author ( s ): Edward J . Sondik Published by : INFORMS Stable URL : <https://www.jstor.org/stable/169635> REFERENCES Linked references are available on. 26(2):282–304, 2019.
- [6] Daniel T. Gillespie. A general method for numerically simulating the stochastic time evolution of coupled chemical reactions. *Journal of Computational Physics*, 22(4):403–434, 1976.
- [7] Simon Godsill. PARTICLE FILTERING: THE FIRST 25 YEARS AND BEYOND. pages 7760–7764, 2019.
- [8] Lirong Huang, Loïc Paulevé, Christoph Zechner, Michael Unger, Anders Hansen, and Heinz Koepl. Supporting Information for Reconstructing dynamic molecular states from single-cell time series. 2016.
- [9] Leslie Pack Kaelbling, Michael L Littman, and Anthony R Cassandra. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 101(1):99–134, 1998.
- [10] P. A.W. Lewis and G. S. Shedler. Simulation of Nonhomogeneous Poisson Processes By Thinning. *Naval research logistics quarterly*, 26(3):403–413, 1979.

- [11] Eric Libby, Theodore J. Perkins, and Peter S. Swain. Noisy information processing through transcriptional regulation. *Proceedings of the National Academy of Sciences of the United States of America*, 104(17):7151–7156, 2007.
- [12] K P Murphy. A survey of POMDP solution techniques. *Environment*, 2(September):X3, 2000.
- [13] Uri Nodelman, Christian R Shelton, and Daphne Koller. Continuous Time Bayesian Networks. 1995.
- [14] Yosihiko Ogaata. On Lewis’ Simulation Method for Point Processes. *IEEE Transactions on Information Theory*, 27(1):23–31, 1981.
- [15] Theodore J. Perkins and Peter S. Swain. Strategies for cellular decision-making. *Molecular Systems Biology*, 5(326):1–15, 2009.
- [16] Marian-Andrei Rizoïu, Young Lee, and Swapnil Mishra. Hawkes processes for events in social media. *Frontiers of Multimedia Research*, pages 191–218, 2017.
- [17] G. K. Robinson. Confidence Intervals and Regions. *Wiley Encyclopedia of Clinical Trials*, 5(2):251–255, 2008.
- [18] Katsushige Sawaki and Akira Ichikawa. Optimal Control for Partially Observable Markov Decision Processes Over an Infinite Horizon. *Journal of the Operations Research Society of Japan*, 21(1):1–16, 1978.
- [19] Vahid Shahrezaei and Peter S. Swain. The stochastic nature of biochemical networks. *Current Opinion in Biotechnology*, 19(4):369–374, 2008.
- [20] Lukas Studer, Loïc Paulevé, Christoph Zechner, Matthias Reumann, María Rodríguez Martínez, and Heinz Koeppl. Marginalized continuous time Bayesian networks for network reconstruction from incomplete observations. *30th AAAI Conference on Artificial Intelligence, AAAI 2016*, pages 2051–2057, 2016.
- [21] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning, Second Edition: An Introduction - Complete Draft*. 2018.
- [22] Ying Tan and Zhong yang Zheng. Research Advance in Swarm Robotics. *Defence Technology*, 9(1):18–39, 2013.
- [23] Sebastian Thrun. Monte Carlo POMDPs. 40(10):117–151, 1904.
- [24] C. C. White and D. P. Harrington. Application of Jensen’s inequality to adaptive suboptimal design. *Journal of Optimization Theory and Applications*, 32(1):89–99, 1980.

## A. Amalgamation Operation

A CTBN with multiple variables can be represented with a single CIM. This is done by amalgamation operation. Amalgamation defines a combining operation over multiple CIMs and produces a single CIM for the entire system. [13]

### A.1. Amalgamation of Independent Processes

Consider a CTBN with graph  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$  over two variables such that  $\mathcal{V} = \{X_1, X_2\}$ . Assume variables  $X_1$  and  $X_2$  with intensity matrices  $\mathbf{Q}_1$  and  $\mathbf{Q}_2$ , are both parent nodes, i.e.  $\mathcal{E} = \emptyset$  and  $Par_{\mathcal{G}}(X_1) = Par_{\mathcal{G}}(X_2) = \emptyset$ . This CTBN can be identified as a subsystem of the CTBN model described in Section 3.1.1.

Analogous to Equation 2.5, Markov transition function for the joint process can be derived as

$$\begin{aligned} \Pr(X_P(t+h) = x'_p \mid X_P(t) = x_p) &= \Pr(X_1(t+h) = x'_1, X_2(t+h) = x_2 \mid X_1(t) = x_1, X_2(t) = x_2) \\ &= \Pr(X_1(t+h) = x'_1 \mid X_1(t) = x_1, X_2(t) = x_2) \\ &\quad \Pr(X_2(t+h) = x_2 \mid X_1(t) = x_1, X_2(t) = x_2) \\ &= (\delta_{x'_1, x_1} + hq_{x_1, x'_1}^1 + o(h))(1 + hq_{x_2, x_2}^2 + o(h)) \\ &= \delta_{x'_1, x_1} + hq_{x_1, x'_1}^1 + h\delta_{x'_1, x_1} q_{x_2, x_2}^2 + o(h) \end{aligned} \quad (\text{A.1})$$

where  $x_1, x'_1 \in \mathcal{X}_1$ ,  $x_2, x'_2 \in \mathcal{X}_2$ ,  $x_p = (x_1, x_2)$ ,  $x'_p = (x'_1, x_2) \in \mathcal{X}_P$ .

Suppose the intensity matrices of  $X_1$  and  $X_2$  are in the form

$$\mathbf{Q}_i = \begin{bmatrix} -q_0^i & q_0^i \\ q_1^i & -q_1^i \end{bmatrix} \quad \text{for } i \in \{1, 2\} \quad (\text{A.2})$$

Then the intensity matrix for the joint process  $X_P$  with factorising state space  $\mathcal{X}_P = \mathcal{X}_1 \times \mathcal{X}_2$  can be written as

$$\mathbf{Q}_P = \begin{bmatrix} -q_0^2 - q_0^1 & q_0^2 & q_0^1 & 0 \\ q_1^2 & -q_1^2 - q_0^1 & 0 & q_0^1 \\ q_1^1 & 0 & -q_1^1 - q_0^2 & q_0^2 \\ 0 & q_1^1 & q_1^2 & -q_1^1 - q_1^2 \end{bmatrix} \quad \text{for } i \in \{1, 2\} \quad (\text{A.3})$$

As it can be observed from Equation A.3, the transition intensities which corresponds to state transition in both variables, i.e. anti-diagonal entries, are zero, due to the one of the assumptions in CTBN framework that only one variable can transition at a time, as given in Section 2.2.



## B. Marginalized Likelihood Function for Homogenous Continuous Time Markov Processes

Let  $X$  be a homogenous CTMP. For convenience, it is assumed to be binary-valued,  $\chi = \{x_0, x_1\}$ . The transition intensity matrix can be written in the following form:

$$\mathbf{Q} = \begin{bmatrix} -q_0 & q_0 \\ q_1 & -q_1 \end{bmatrix} \quad (\text{B.1})$$

where the transition intensities  $q_0$  and  $q_1$  are gamma-distributed with parameters  $\alpha_0, \beta_0$  and  $\alpha_1, \beta_1$ , respectively. The marginal likelihood of a sample trajectory  $X^{[0,T]}$  can be written as follows:

$$\begin{aligned} P(X^{[0,T]}) &= \int P(X^{[0,T]} | Q) P(Q) dQ \\ &= \int_0^\infty \prod_{j \neq i} \exp(-q_{i,j} T[x_i]) \frac{q_{i,j}^{M[x_i, x_j]} \beta_{i,j}^{\alpha_{i,j}} q_{i,j}^{\alpha_{i,j}-1} \exp(-\beta_{i,j} q_{i,j})}{\Gamma(\alpha_{i,j})} dq_{i,j} \\ &= \prod_{i \in \{0,1\}} \int_0^\infty q_i^{M[x_i]} \exp(-q_i T[x_i]) \frac{\beta_i^{\alpha_i} q_i^{\alpha_i-1} \exp(-\beta_i q_i)}{\Gamma(\alpha_i)} dq_i \\ &= \prod_{i \in \{0,1\}} \frac{\beta_i^{\alpha_i}}{\Gamma(\alpha_i)} \int_0^\infty q_i^{M[x_i] + \alpha_i - 1} \exp(-q_i (T[x_i] + \beta_i)) dq_i \end{aligned} \quad (\text{B.2})$$

$$\begin{aligned} &= \prod_{i \in \{0,1\}} \frac{\beta_i^{\alpha_i}}{\Gamma(\alpha_i)} \left( -(T[x_i] + \beta_i)^{-M[x_i] - \alpha_i} \Gamma(M[x_i] + \alpha_i, q_i (T[x_i] + \beta_i)) \right) \Big|_0^\infty \\ &= \prod_{i \in \{0,1\}} \frac{\beta_i^{\alpha_i}}{\Gamma(\alpha_i)} ((T[x_i] + \beta_i)^{-M[x_i] - \alpha_i} \Gamma(M[x_i] + \alpha_i)) \end{aligned} \quad (\text{B.3})$$

In Equation B.2, the integral is solved using computer algebra system WolframAlpha as follows:

$$\int x^a \exp(-xb) dx = -b^{-a-1} \Gamma(a+1, bx) + C \quad (\text{B.4})$$

## C. Equivalence Classes of Observation Models

The equivalence classes are inherent to the problem setting and caused by two reasons. In this chapter, these reasons are explained and illustrated.

### Identical Effect on Belief State

Some observation models fall into the same class due to their exact same effect on the belief state. An example of this situation is illustrated below. In order to show the exact equivalence, the simulations are performed belief state update using exact update method as described in Section 3.1.2.2. Consider the problem of calculating the likelihood of one sample  $S^{[0,T]}$  given two observation models,  $\psi_1$  and  $\psi_2$  given below. Given a sample of parent trajectories shown in Figure C.1, it is obvious that these two observation models lead to different observation trajectories as shown in Figure C.2(a) and Figure C.3(a). Nonetheless, using Equation 3.11, the resulting belief state is exactly the same. This leads to the exact same trajectories for  $Q_3$  and the likelihood of the sample given these two observation models,  $p(S^{[0,T]} | \psi_1)$  and  $p(S^{[0,T]} | \psi_2)$  end up exactly same.

$$\psi_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad p(S^{[0,T]} | \psi_1) = -83.334 \quad (\text{C.1})$$

$$\psi_2 = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad p(S^{[0,T]} | \psi_2) = -83.334 \quad (\text{C.2})$$

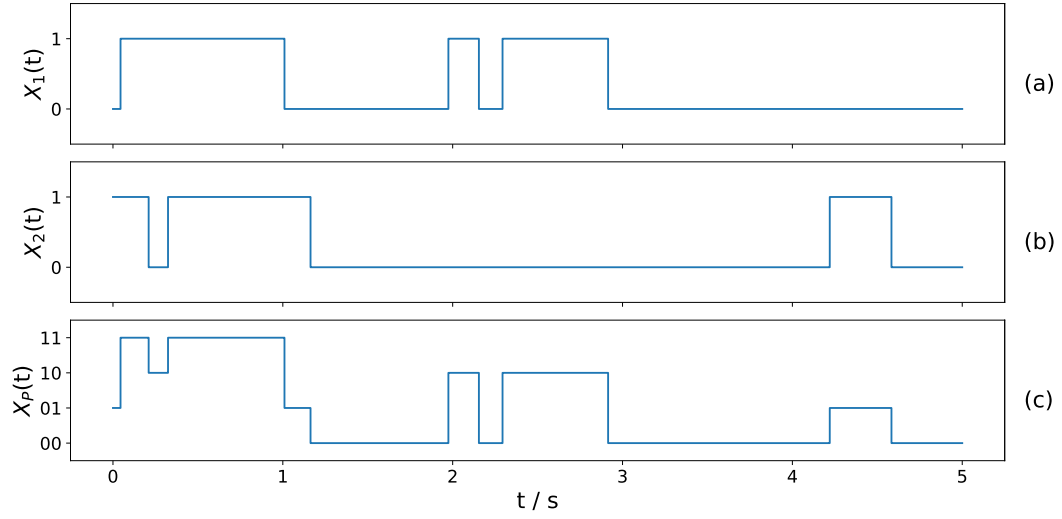


Figure C.1.: Parent trajectories for the models leading to the same belief state

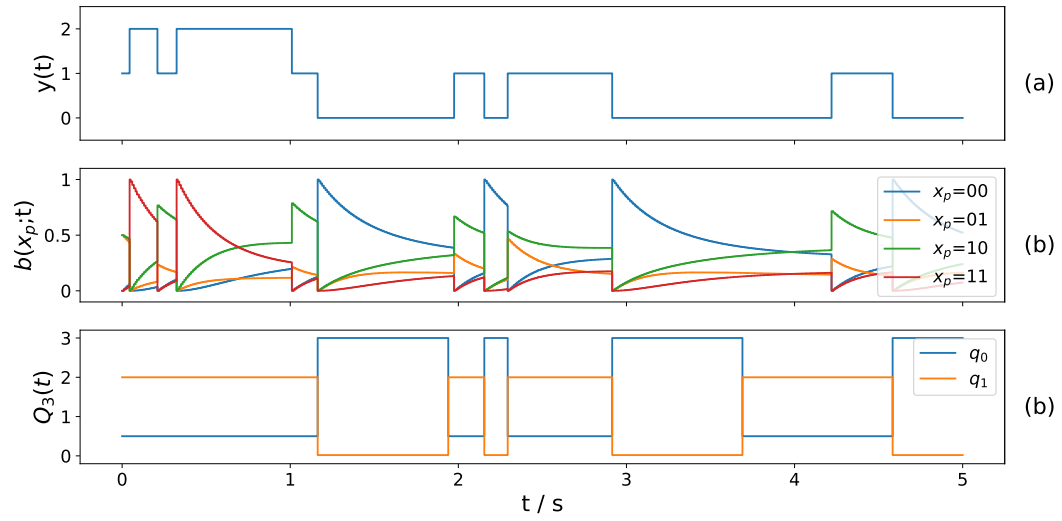


Figure C.2.: Observation, belief state and  $Q_3$  trajectories derived by  $\psi_1$  in Equation C.1 corresponding to parent trajectories in Figure C.1

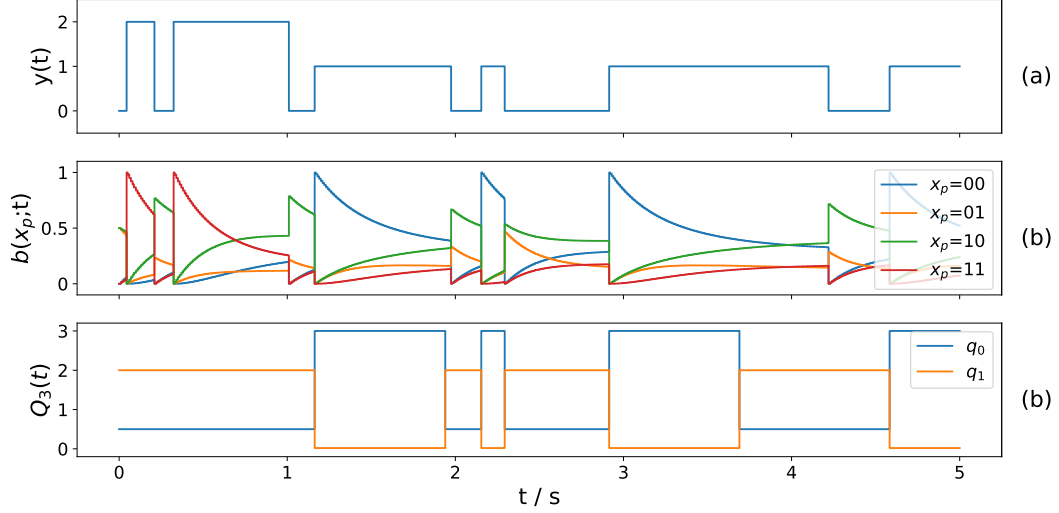


Figure C.3.: Observation, belief state and  $\mathbf{Q}_3$  trajectories derived by  $\psi_2$  in Equation C.2 corresponding to parent trajectories in Figure C.1

## Combination of Belief State and Policy

For some observation model, the reason of equivalence is that even though the belief state are different, the policy  $\pi(b)$  leads to same trajectory for  $Q_3$ . This case is exemplified below where the simulations are performed belief state update using exact update method as described in Section 3.1.2.2. Consider the problem of calculating the likelihood of one sample  $S^{[0,T]}$  given two observation models,  $\psi_1$  and  $\psi_2$  given below. Given a sample of parent trajectories shown in Figure C.4, these observation models lead to different observation trajectories as shown in Figure C.5(a) and Figure C.6(a), which results in different belief state trajectories given in Figure C.5(b) and Figure C.6(b). However, the policy leads to the exact same trajectories for  $Q_3$  and the likelihood of the sample given these two observation models,  $p(S^{[0,T]} | \psi_1)$  and  $p(S^{[0,T]} | \psi_2)$  end up exactly same.

$$\psi_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad p(S^{[0,T]} | \psi_1) = -80.648 \quad (\text{C.3})$$

$$\psi_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \quad p(S^{[0,T]} | \psi_2) = -80.648 \quad (\text{C.4})$$

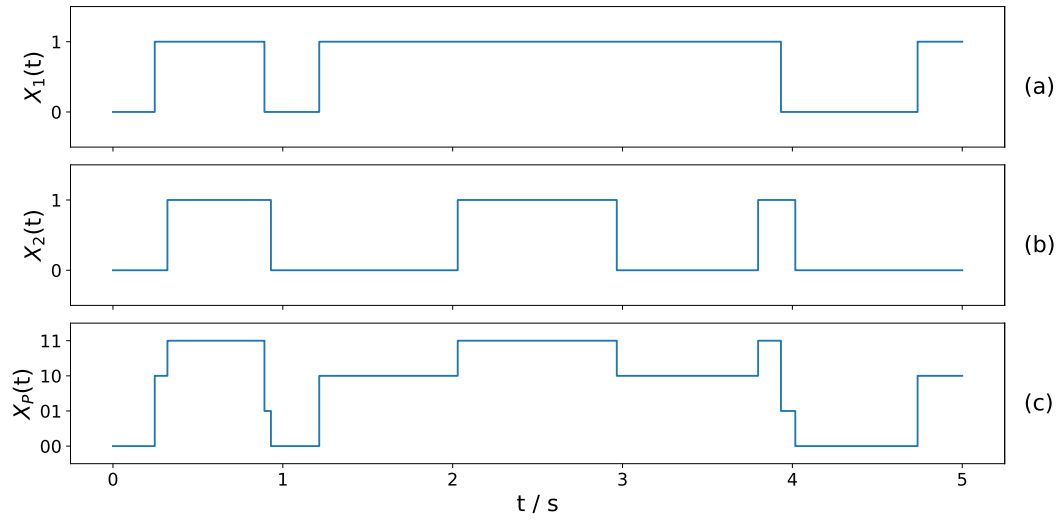


Figure C.4.: Parent trajectories for the models leading to the same behaviour

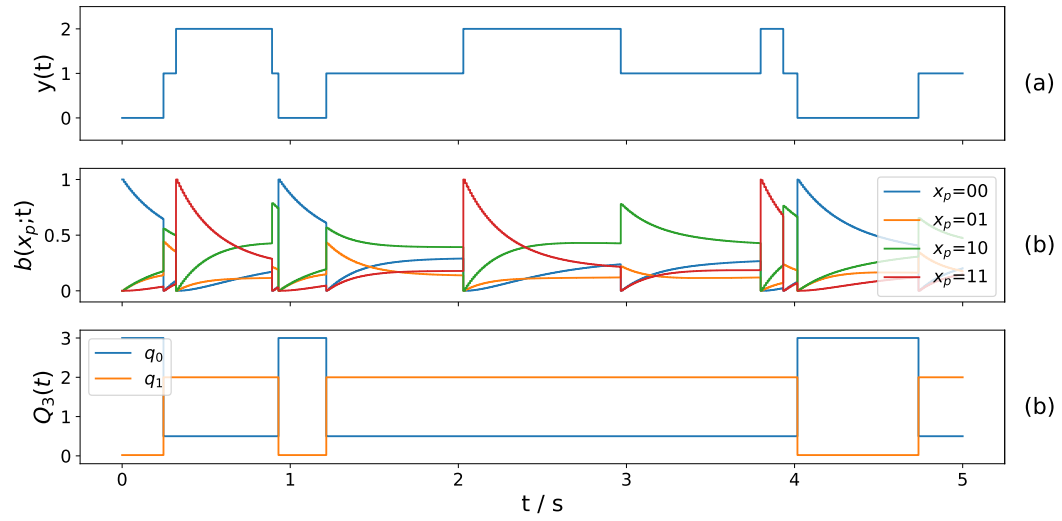


Figure C.5.: Observation, belief state and  $Q_3$  trajectories derived by  $\psi_1$  in Equation C.3 corresponding to parent trajectories in Figure C.4

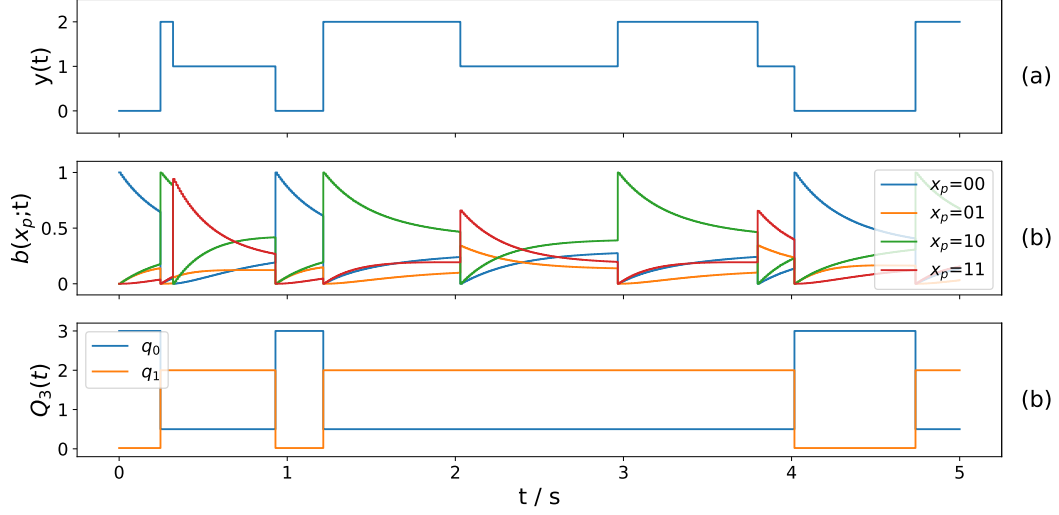


Figure C.6.: Observation, belief state and  $\mathbf{Q}_3$  trajectories derived by  $\psi_2$  in Equation C.4 corresponding to parent trajectories in Figure C.4

## C.1. Observation Models in Experiments

As mentioned in Section 4.3.1, the inference problem is reduced to maximum likelihood estimation between 10 classes. One observation model is picked as representative of each class and considered for the inference problem. These observation models are given below. This set of observation models are referred to as  $\psi$ .

$$\begin{aligned}
 \psi_{\text{true}} = \psi_0 &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} & \psi_1 &= \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} & \psi_2 &= \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} & \psi_3 &= \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix} & \psi_4 &= \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \\
 \psi_5 &= \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} & \psi_6 &= \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} & \psi_7 &= \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} & \psi_8 &= \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} & \psi_9 &= \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}
 \end{aligned} \tag{C.5}$$

## D. Additional Results

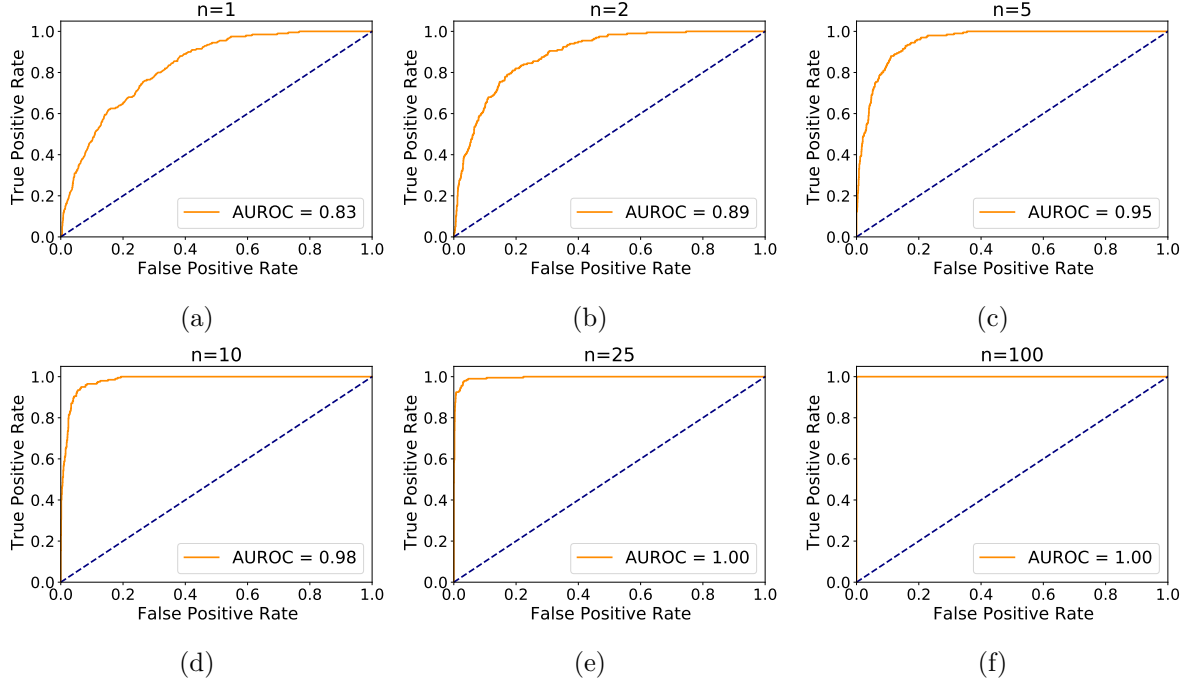


Figure D.1.: ROC curves with  $n$  number of trajectories for dataset generated using exact belief state update,  $\psi_0$ -vs-rest

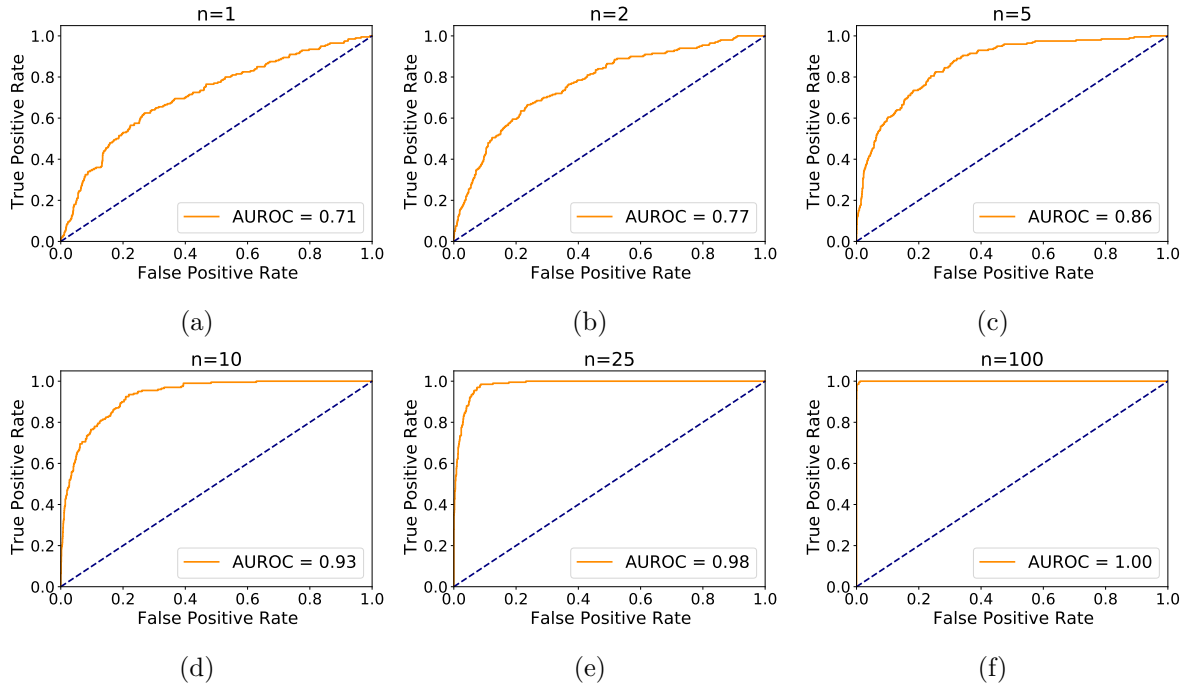


Figure D.2.: ROC curves with  $n$  number of trajectories for dataset generated using particle filtering,  $\psi_0$ -vs-rest