CZECH TECHNICAL UNIVERSITY IN PRAGUE
Faculty of Nuclear Sciences and Physical Engineering

# Mixture Ratios for Decision Making

# Podílové směsové modely pro rozhodování

Master's thesis

Author:            **Bc. Marko Ruman**

Supervisor:        **Ing. Miroslav Kárný, DrSc.**

Academic year:     2017/2018

*Název práce:*

**Podílové směsové modely pro rozhodování**

*Autor:* Bc. Marko Ruman

*Obor:* Matematické inženýrství

*Druh práce:* Diplomová práce

*Vedoucí práce:* Ing. Miroslav Kárný, DrSc., Oddělení Adaptivních systémů, Ústav Teorie Informace a Automatizace

*Abstrakt:* Směsi konečného počtu hustot pravděpodobností s komponentami z exponenciální rodiny slouží jako flexibilní modely vícerozměrných systémů. Avšak, s výjimkou pár úzce specializovaných případů, používané dynamické modely předpokládají datově nezávislé váhy jednotlivých směsových komponent. Použití takových vah je nelogické a omezuje praktickou použitelnost daných modelů. Tato práce řeší dané omezení použitím podílu konečných směsí jako parametrického modelu. V práci je použití tohoto modelu motivováno, je vybudováno jeho rekurzivní bayeovské odhadování a je ukázaná jeho vhodnost i pro modelování obecných diskrétně-spojitých systémů. Dále je provedeno porovnání podílových směsí se standardními směsmi s konstantními vahami, a to jak na simulovaných datech, tak na reálných datech vývoje cen futures. Pro případ podílu směsí Markovských řetězců je vybudován návrh rozhodovací strategie.

*Klíčová slova:* Dynamické systémy, bayesovské učení, směsové modely, podíl směsí

*Title:*

**Mixture Ratios for Decision Making**

*Author:* Bc. Marko Ruman

*Abstract:* Finite mixtures of probability densities with components from exponential family serve as flexible parametric models of high-dimensional systems. However, with a few specialized exceptions, these dynamic models assume data-independent weights of mixture components. Their use is illogical and restricts the modelling applicability. The work overcomes this restriction by considering ratios of finite mixtures as parametric models. It motivates them and elaborates their approximate Bayesian recursive estimation and shows its ability to model also general discrete-continuous systems. It also compares the mixture ratio model with the standard, fixed-weight mixtures, either on simulation data as well as on real futures price data. For the case of mixture ratio of Markov chains, the design of decision strategy is build.

*Key words:* Dynamic systems, Bayesian learning, mixture models, mixture ratio

# Contents

# Introduction

Decision making (DM) is a targeted choice among the available options[1]. This broadly understood DM covers many traditional fields including machine learning, [23], signal processing, [24], estimation and filtering, [17], hypothesis testing, [14], classification and pattern recognition, [11], knowledge sharing, [22], reinforcement learning, [27], control, [13], etc. The amount of relevant results is excessive as seen, e.g., in [16]. To cover all of the mentioned fields, this work uses its own vocabulary to formalize a general DM problem.

Any general solution of DM problems leads to a strategy, a collection of decision rules mapping the knowledge on actions, [26]. The chosen strategy should meet the DM aim in the best way under the encountered circumstances. Bayesian DM, [25], proved to be a powerful methodology of DM facing incomplete knowledge, uncertainty and randomness of the dynamic system to which actions relate. This work stays within it. Bayesian DM relates DM consequences to the acquired knowledge and models the used actions by conditional probabilities, here described by conditional probability densities (pd) or conditional probability functions (pf).

Generally, actions are chosen recursively while enriching the available knowledge. This gives a chance to learn, i.e. to improve gradually the system model. The learning redistributes belief (probability of) into the model adequacy within the set of considered models, [9]. It gradually extracts the knowledge and can be used in data-streams processing or employed as a feature extractor.

The need for learning arises whenever a good model is a priori unknown. A learning is possible iff the learnt relations practically do not change during the knowledge accumulation. This is assumed and a set of parametric models, distinguished by a constant multivariate parameter serving as a "pointer" to its respective members, is considered. Which one points to the best model is a priori unknown. Bayesian learning offers the unambiguous deductive way, Bayes' rule, [10], how to redistribute the belief concerning the model quality. It accumulates knowledge into the posterior pd and provides the predictive pd, which serves as the system model needed for DM. The achievable modeling and thus DM quality are predetermined by the considered set of parametric models.

This work offers *ratios* of finite mixtures, [4], with components from exponential family (EF), [7], as such black-box, [12], universally approximating, [15], models. The recursive learning of ratios of finite mixtures is inevitably approximate and consequently endangered by accumulation of approximation errors. The approximate, but feasible recursive Bayesian learning of models, equipped with a counteractor of approximation errors, from this *extremely flexible but yet unconsidered model set*, is developed, [30].

Often, the modelled observations as well as explanatory variables are discrete or discretized. Then, a high-order Markov chain provides their universal model. It relates the predicted observation to a finite-dimensional regression vector containing the past observations and explanatory variables. Its recursive Bayesian estimation, which provides a lossless compression of the available knowledge, is formally simple. Basically, it counts joint occurrences of the predicted variable and the corresponding regression vector. The applicability is, however, strongly limited by the curse of dimensionality, [8], as the size of

---

[1]A part of this work was submitted for publication [30].

the occurrence array blows up with the number of possible occurrence instances. Then, the observations insufficiently populate it.

The components of mixture of Markov chains belongs to EF, thus such high-order Markov chain can be well approximated by the developed mixture ratio model, which will counteract the curse of dimensionality of general Markov chain models, see [29] for comparison. This work develops the mentioned approximate Bayesian learning for the mixture ratios of Markov chains with conjugated Dirichlet distribution on its parameters.

Another important and commonly solved task is modelling of continuous variables. Almost any continuous probability density can be approximated by a finite mixture of Gaussian distribution to an arbitrary precision, [4]. Again, Gaussian components belong to EF, hence the approximate Bayesian learning of mixture ratios with such components is build in this work. These two types of EF components can be easily combined and can provide a general parametric models of mixed discrete-continuous systems.

Besides, this work provides comparison of learning of mixture ratio and standard mixture (developed in [18]) of Markov chain components in simulation set-up as well as on real futures trading data, where the suitability of mixture ratio model for decision making is demonstrated. The necessary design of decision strategies for the mixture ratio of Markov chain components is build.

Layout of the work is as follows: Chapter 1 introduces mathematical notation, functions and theorems used in the text. Chapter 2 formalizes a general model of a dynamic system and summarizes its approximate recursive Bayesian learning. Chapter 3 introduces finite mixture ratio as a parametric model. Chapter 4 describes the use of general learning algorithm from Chapter 2 for the finite mixture ratio parametric model. Chapter 5 introduces mixture ratio of Markov chains and builds its learning. Chapter 6 specifies mixture ratio of Gaussian components and derives its learning. Chapter 7 provides simulation and practical examples comparing the mixture ratio model with standard mixture model, where also the design of decision strategies for the mixture ratio of Markov chain components is created. Conclusion summarizes results of this work and outlines open questions for the further examination. Appendix provides derivation of some used formulas.

# Chapter 1

# Mathematical Preliminaries

This chapter introduces notation as well as theorems and functions used throughout the text. The list of notation symbols used in text:

- $\mathbb{N}$, $\mathbb{Z}$, $\mathbb{R}$, $\mathbb{R}^+$ stand for sets of *natural*, *integer*, *real* and *positive real* numbers, respectively,

- $\mathbb{R}^n$ denotes $n$-dimensional vector space, $\mathbb{R}^{nm}$ denotes vector space of matrices with dimensions $n \times m$; $n, m \in \mathbb{N}$,

- $a = (a_1, ..., a_n) \in \mathbb{R}^n$ stands for a $n$-dimensional column vector,

- for a matrix $A \in \mathbb{R}^{nn}$, $\text{tr}(A)$ denotes its trace, $\text{tr}(A) = \sum_{i=1}^{n} A_{ii}$,

- the bold symbols, for instance $\mathbf{A}$, will stand for the sets of all possible values of its members $A \in \mathbf{A}$,

- $|\mathbf{A}|$ denotes cardinality of $\mathbf{A}$,

- $\chi_{\mathbf{A}}(x)$ is a characteristic function of a set $\mathbf{A}$

$$\chi_{\mathbf{A}}(x) = \begin{cases} 1, & \text{if } x \in \mathbf{A}, \\ 0, & \text{otherwise,} \end{cases}$$

- $[a, b]^n$ stands for a $n$-ary Cartesian product, where $a, b \in \mathbb{R}$, $a < b$, $n \in \mathbb{N}$ and $[a, b]$ is a real interval,

- $\mathbb{N}^n$ stands for the $n$-ary Cartesian product of natural numbers, i.e. $\mathbb{N}^n = \bigtimes_{i=1}^{n} \mathbb{N}$,

- $\langle a \,|\, b \rangle$ denotes a scalar product

  - for vectors $a, b \in \mathbb{R}^n$, $\langle a \,|\, b \rangle$ is the dot product, $\langle a \,|\, b \rangle = \sum_{i=1}^{n} a_i b_i$
  - for matrices $a, b \in \mathbb{R}^{nn}$, $\langle a \,|\, b \rangle = \text{tr}(ab)$

- $0^0$ is defined as $0^0 = 1$.

Throughout the text, the functions with different arguments are assumed to be different, e.g. $\text{B}(\psi_t)$ and $\text{B}(\psi_{t;c})$.

If it is not necessary to distinguish realization and random variable, for the probability of the observation of a discrete random variable $\tilde{x}$ with a value $\tilde{x} = x$, the symbol $\text{P}(x)$ will be used[1]. The conditional

---

[1] For discrete valued random variables $\text{P}(x)$ is a probability function (pf), for a continuous random variables $\text{P}(x)$ is used for the probability density function (pd) of a random variable $\tilde{x}$. Alternatively, the symbols $\text{M}(x), \text{S}(x), \text{J}(x)$ will be used for pfs or pds.

probability of the observation of a random variable $\tilde{x}$ with a value $\tilde{x} = x$, which depends on a random variable $\tilde{y}$ with a known realization $\tilde{y} = y$, the symbol $\mathrm{P}(x|y)$ is used, more precisely:

$$\mathrm{P}(x|y) = \frac{\mathrm{P}(x,y)}{\mathrm{P}(y)},$$

where $\mathrm{P}(y) > 0$ denotes the probability, that random variable $\tilde{y}$ will be realized with a value $\tilde{y} = y$, similarly $\mathrm{P}(x,y)$ stands for the probability, that random variables $\tilde{x}$ and $\tilde{y}$ will be realized with a value $\tilde{y} = y$ and $\tilde{x} = x$.

The conditional probability $\mathrm{P}(a,b|c)$ where $\tilde{a}, \tilde{b}$ and $\tilde{c}$ are random variables, can be calculated via chain rule:

$$\mathrm{P}(a,b|c) = \mathrm{P}(a|b,c)\mathrm{P}(b|c). \tag{1.1}$$

The marginal probability $\mathrm{P}(a|c)$ can be calculated as:

$$\mathrm{P}(a|c) = \int_{\mathbf{b}} \mathrm{P}(a,b|c)db, \tag{1.2}$$

where $\mathbf{b}$ is a set of all values of the random variable $\tilde{b}$.

Let $t \in \mathbf{t}$ be a time index from the set of the time indices $\mathbf{t} \subset \mathbb{N}$. The notation $\mathrm{P}_t(\Theta) = \mathrm{P}(A_t|B_t,\Theta)$ will be used whenever realization of $A_t, B_t$ is inserted into the conditional pf/pd $\mathrm{P}(A_t|B_t,\Theta)$, i.e. when it is treated as likelihood function, $\Theta$ is usually a parameter vector.

**Definition 1** (Statistic). A *statistic* V is here (generally, the range of the statistic can be arbitrary) a finite-dimensional, measurable function of a data sequence
$D^t = \{D_t, D_{t-1}, ..., D_1, D_0\} \in \mathbf{D}^t$:

$$\mathrm{V} : \mathbf{D}^t \mapsto \mathbb{R}^n$$

The notation $\mathrm{V}_t = \mathrm{V}(D^t)$ will be used. □

**Definition 2** (Sufficient statistic, [21]). A statistic V is *sufficient* for a parameter $\Theta$ iff:

$$\mathrm{P}(\Theta|\mathrm{V}_t, D^t) = \mathrm{P}(\Theta|\mathrm{V}_t)$$

□

**Definition 3** (Exponential family). Let $\Theta$ be a finite-dimensional parameter vector (or matrix), $x = (x_1, ..., x_k)$ a $k$-dimensional random variable, then the conditional pd $\mathrm{P}(x|\Theta)$ belongs to *Exponential family* of distributions iff has the following form:

$$\mathrm{P}(x|\Theta) = \exp\langle \mathrm{B}(x) \,|\, \mathrm{C}(\Theta)\rangle, \tag{1.3}$$

where B and C are known finite-dimensional vector or matrix functions, $\langle \mathrm{B}(x)\,|\,\mathrm{C}(\Theta)\rangle$ denotes a scalar product. □

**Definition 4** (Gamma function). The gamma function $\Gamma$ is deifned for $x \in \mathbb{R}^+$ as follows

$$\Gamma(x) = \int_0^\infty t^{x-1}e^{-t}dt.$$

□

**Remark 1.** Gamma function satisfies $\Gamma(x+1) = x\Gamma(x)$, [19]. □

**Definition 5** (Digamma function). Digamma function $\psi(x)$ is defined for $x \in \mathbb{R}^+$ as follows:

$$\psi(x) = \frac{\partial}{\partial x} \ln(\Gamma(x)).$$

□

**Remark 2.** Digamma function satisfies $\psi(x+1) = \frac{1}{x} + \psi(x)$, [19]. □

**Definition 6** (Beta function,[19]). Beta function $\mathrm{Be}(x,y) : \mathbb{R}^+ \times \mathbb{R}^+ \mapsto \mathbb{R}^+$ is defined as follows

$$\mathrm{Be}(x,y) = \int_0^1 t^{x-1}(1-t)^{y-1}dt,$$

or equivalently

$$\mathrm{Be}(x,y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}.$$

The multivariate beta function $\mathrm{Be}(v) : \times_{i=1}^n \mathbb{R}^+ \mapsto \mathbb{R}^+$ is defined as follows

$$\mathrm{Be}(v) = \frac{\prod_{i=1}^n \Gamma(v_i)}{\Gamma(\sum_{i=1}^n v_i)}.$$

□

**Definition 7** (Kronecker delta function). The Kronecker delta function $\delta_{xy}$ is defined on $\mathbb{Z} \times \mathbb{Z}$ as

$$\delta_{xy} = \begin{cases} 1 & x = y, \\ 0 & x \neq y. \end{cases} \tag{1.4}$$

For a set $\mathbf{c} = \{1,..,|\mathbf{c}|\}$, the symbol $\delta(c)$ will stand for the following vector:

$$\delta(c) = (\delta_{1c},...,\delta_{|\mathbf{c}|c}).$$

□

**Definition 8** (Kullback–Leibler divergence, [1]). *Kullback–Leibler divergence* P to Q, KL(P‖Q), is a real-valued function mapping of two pfs[2] P and Q to a real number and it is defined as follows:

$$\mathrm{KL(P\|Q)} = \sum_{a \in \mathbf{a}} \mathrm{P}(a) \ln \frac{\mathrm{P}(a)}{\mathrm{Q}(a)}. \tag{1.5}$$

□

**Definition 9** (Kerridge inaccuracy, [20]). *Kerridge inaccuracy* K(P‖Q) is a real-valued function mapping of two pds P and Q to a real number and it is defined as:

$$\mathrm{K(P\|Q)} = -\int_{\boldsymbol{\Theta}} \mathrm{P}(\Theta) \ln(\mathrm{Q}(\Theta)) d\Theta. \tag{1.6}$$

□

---

[2]It can also be defined for pds via its integral form, but this work uses only the version for pfs.

**Theorem 1** (Bayes' rule, [10])**.** Let $a$ and $b$ be continuous random variables[3] and $\mathbf{a}$ a set of all possible values of random variable $a$. Then the conditional pd $P(a|b)$ can be expressed as follows (for $\int_{\mathbf{a}} P(b|a)P(a)da > 0$):

$$P(a|b) = \frac{P(b|a)P(a)}{\int_{\mathbf{a}} P(b|a)P(a)da} \propto P(b|a)P(a),\qquad(1.7)$$

where $\propto$ denotes equality up to normalizing constant. $\qquad\square$

---

[3]In the text, Bayes' rule will be used for continuous random variable $a$, thus the set $\mathbf{a}$ is a subset of the real numbers $\mathbb{R}$, the random variable $\tilde{b}$ can be either a discrete valued or continuous.

# Chapter 2

# Dynamic System Model and Its Approximate Learning

A model of a dynamic system in DM theory is formalized in Bayesian way as follows. Observations $O_t \in \mathbf{O}$, stimulated by actions $A_t \in \mathbf{A}$, are made at the discrete-time moments $t \in \mathbf{t} = \{1, 2, ..., |\mathbf{t}|\}$. A data sequence $D^t \in \mathbf{D}$ is formed by data records: $D^t = \{D_t, D_{t-1}, ..., D_1, D_0\}$, where a data record $D_t = (O_t, A_t)$, $t \in \mathbf{t}$, is formed by an observation $O_t$ and an action $A_t$, $D_0$ stands for the initial knowledge about the system. considered variables are random so that modelling must be done by pds.

To describe a behaviour of the dynamic system, the parametric pd of the observation $O_t$ is used.

**Definition 10** (Parametric system model). The parametric system model is a pd of the observation $O_t \in \mathbf{O}$ conditioned on the data sequence $D^{t-1} \in \mathbf{D}$, the action $A_t \in \mathbf{A}$ and on an unknown parameter $\Theta \in \mathbf{\Theta}$:

$$\mathrm{M}(O_t | A_t, D^{t-1}, \Theta), \tag{2.1}$$

where $\mathbf{\Theta}$ is a *parametric space* and it is a subset of a real $n$-dimensional vector space, i.e. $\mathbf{\Theta} \subset \mathbb{R}^n$, $n \in \mathbb{N}$. The notation $\mathrm{M}_t(\Theta) = \mathrm{M}(O_t | A_t, D^{t-1}, \Theta)$ will be used whenever realization of $A_t, D^{t-1}$ is inserted into the parametric system model, i.e. when it is treated as the likelihood function. $\qquad \square$

The goal of DM theory is to influence the dynamic system by taking appropriate actions to achieve desired observations (generally, states of the modelled system). To describe the next observation of the system, the pd $\mathrm{P}(O_t | D^{t-1})$ is used. Using (1.1) and (1.2) and incorporating the parametric system model (2.1), it can be expressed as follows:

$$\mathrm{P}(O_t | A_t, D^{t-1}) = \int_{\Theta} \mathrm{M}(O_t | A_t, D^{t-1}, \Theta) \pi(A_t | \Theta, D^{t-1}) \mathrm{P}(\Theta | D^{t-1}) d\Theta,$$

where $\pi(A_t | D^{t-1}, \Theta)$ is a pd that describes a randomized decision rule[1], $\mathrm{P}(\Theta | D^{t-1})$ is a posterior pd storing actual knowledge about the parameter $\Theta \in \mathbf{\Theta}$; the notation $\mathrm{P}_{t-1}(\Theta) = \mathrm{P}(\Theta | D^{t-1})$ will be used. Generally, the assumption of *natural conditions of control*, [32], is adopted, i.e. $\pi(A_t | \Theta, D^{t-1}) = \pi(A_t | D^{t-1})$. It means that parameters are unknown to considered decision rules.

---

[1] The action generator influences the pd $\P(O_t | D^{t-1})$ via decision rule $\pi(A_t | D^{t-1}, \Theta)$ according to a given goal.

After the interaction with the dynamic system, a new data record $D_t = (O_t, A_t)$ is created and $P_{t-1}(\Theta)$ is updated via Bayes' rule (1.7) as follows[2]:

$$\tilde{P}_t(\Theta) = \frac{M_t(\Theta)\pi(A_t|D^{t-1})P_{t-1}(\Theta)}{\int_\Theta M_t(\Theta)\pi(A_t|D^{t-1})P_{t-1}(\Theta)d\Theta} = \frac{M_t(\Theta)P_{t-1}(\Theta)}{\int_\Theta M_t(\Theta)P_{t-1}(\Theta)d\Theta} \propto M_t(\Theta)P_{t-1}(\Theta), \qquad \forall \Theta \in \Theta.$$

(2.2)

To ensure the computational feasibility of updating the pd $P_{t-1}(\Theta)$ during the whole interaction with a dynamic system, the set of feasible pds $\mathbf{P}$ needs to be considered. Generally, the posterior pd $\tilde{P}_t(\Theta)$ obtained by (2.2) does not belong to the set of feasible pds $\mathbf{P}$ and with growing $t \in \mathbf{t}$ it can become more and more complex function of $\Theta$. Therefore, a projection of $\tilde{P}_t(\Theta)$ on pd $\hat{P}_t(\Theta) \in \mathbf{P}$ from this set of feasible pds has to be made.

Works [1] and [2] suggest, that the adequate projection $\hat{P}_t(\Theta)$ is the minimizer of Kerridge inaccuracy (see Definition 9):

$$\hat{P}_t(\Theta) = \underset{P \in \mathbf{P}}{K}(\tilde{P}_t\|P) = \underset{P \in \mathbf{P}}{\mathrm{argmin}} \int_\Theta -\tilde{P}_t(\Theta)\ln(P(\Theta))d\Theta$$

(2.3)

The updating from $P_{t-1}(\Theta)$ to $\hat{P}_t(\Theta)$ can be interpreted as the application of Bayes' rule (1.7) on $P_{t-1}(\Theta)$ but using an unknown, different model than $M_t(\Theta)$. Therefore, using $\hat{P}_t(\Theta)$ as a prior pd for the next learning step may, in general, cause divergence $\hat{P}_t(\Theta)$ from $\tilde{P}_t(\Theta)$ obtained via (2.2) without projection (2.3). The solution preventing from the divergence is to include a data-depending forgetting factor $\lambda_t \in [0, 1]$. [3] implies the most plausible choice of the forgetting factor $\lambda_t$.

**Algorithm 1** (Learning algorithm). One learning step of the learning algorithm is summarized as follows:

1) $\tilde{P}_t(\Theta) \propto M_t(\Theta)P_{t-1}(\Theta)$,

2) $\hat{P}_t = \underset{P \in \mathbf{P}}{\mathrm{argmin}}\, K(\tilde{P}_t\|P)$,

3) $P_t(\Theta) \propto \hat{P}_t^{\lambda_t}(\Theta)P_{t-1}^{1-\lambda_t}(\Theta)$, where $\lambda_t = \dfrac{\left(\int_\Theta M_t(\Theta)P_{t-1}(\Theta)d\Theta\right)^2}{\int_\Theta M_t^2(\Theta)P_{t-1}(\Theta)d\Theta}$,

where $P_{t-1}(\Theta)$ is the prior pd, $K(\tilde{P}_t\|P)$ denotes Kerridge's inaccuracy (2.3) and $P_t(\Theta)$ is used as a prior pd for the next learning step. $\square$

---

[2]This particular form of Bayes' rule is valid only under the mentioned natural conditions of control.

# Chapter 3

# Finite Mixture Ratio Model

To achieve computational feasibility of the recursive learning summarized in Algorithm 1., the following simplifying assumption is made about the model (2.1):

**Assumption 1** (Markov property of the system model)**.** The system model (2.1) is assumed to be time-invariant $n$-order Markov model ($n \in \mathbb{N}$), i.e.

$$\text{M}(O_t|A_t, D^{t-1}, \Theta) = \text{M}(O_t|\psi_t, \Theta) \text{ , where} \tag{3.1}$$

$$\psi_t = \left\{ \begin{array}{ll} A_t, D_{t-1}, ..., D_{t-n} & n \geq 1, \\ A_t & n = 0. \end{array} \right. \text{ is a } regression \ vector.$$

The symbol $\phi_t \in \boldsymbol{\phi}$ will denote a data vector $\phi_t = (O_t, \psi_t)$. □

As stated in Introduction, this work focuses on modelling the dynamic system via mixtures of pds. There are many works on this topic, e.g. [4], however, most of them assume data-independence of the weights of the mixture components. The following approach overcome this limitation by introducing a ratio of finite mixtures.

The model (3.1) can be expressed as follows (using (1.1) and (1.2)):

$$\text{M}(O_t|\psi_t, \Theta) = \frac{\text{J}(O_t, \psi_t|\Theta)}{\int_{\mathbf{O}} \text{J}(O_t, \psi_t|\Theta) dO}. \tag{3.2}$$

Almost any practically met time-invariant joint pd $\text{J}(O_t, \psi_t|\Theta)$ in (3.2) can be approximated by a finite mixture of Gaussian pds, [19], to an arbitraty precision, [4]. A Gaussian pd belongs to the exponential family (EF) (see Definition 3), thus, following EF mixture form of $\text{J}(O_t, \psi_t|\Theta)$ can be considered:

**Definition 11** (Joint pd as a mixture)**.** The joint pd $\text{J}(O_t, \psi_t|\Theta)$ is modelled via a mixture of pds from EF (see Definition 3) as follows:

$$\text{J}(O_t, \psi_t|\Theta) = \text{J}(\phi_t|\Theta) = \sum_{c \in \mathbf{c}} \alpha_c \text{J}_c(\phi_t|\omega_c) = \sum_{c \in \mathbf{c}} \alpha_c \exp \langle \text{B}_c(\phi_t) \,|\, \text{C}_c(\omega_c) \rangle \tag{3.3}$$

$$= \sum_{c \in \mathbf{c}} \alpha_c \exp \langle \text{B}_c(\phi_{t;c}) \,|\, \text{C}_c(\omega_c) \rangle \text{P}_c(\phi_{t;c}^{\mathsf{c}}),$$

where

- $\text{J}_c(\phi_t|\omega_c) = \text{J}_c(O_t, \psi_t|\omega_c)$ is the $c$-th mixture component, it is a member of EF and it is a pd on $\boldsymbol{\phi}$. It is assumed to have the form: $\text{J}_c(\phi_t|\omega_c) = \exp \langle \text{B}_c(\phi_t) \,|\, \text{C}_c(\omega_c) \rangle = \exp \langle \text{B}_c(\phi_{t;c}) \,|\, \text{C}_c(\omega_c) \rangle \text{P}_c(\phi_{t;c}^{\mathsf{c}})$

with $\phi_{t;c}$ being a component-specific subvector of $\phi_{t;c}$ and $P_c(\phi_{t;c}^{\mathsf{c}})$ is a non-parametrized pd on its complement $\phi_{t;c}^{\mathsf{c}}$ with respect to $\phi_t$; $B_c(\phi_t)$, $C_c(\omega_c)$ and $B_c(\phi_{t;c})$ are known, real-valued vector functions, $B_c(\phi_{t;c}) = (B_{c1}(\phi_{t;c})\ldots B_{cm_c}(\phi_{t;c}))$, $C_c(\omega_c) = (C_{c1}(\omega_c),\ldots C_{cm_c}(\omega_c))$

- $\mathbf{c} = \{1,\ldots,|\mathbf{c}|\}$, $|\mathbf{c}|$ denotes the number of mixture components

- the parameter vector $\Theta$ has the following form: $\Theta = \left(\alpha_c, (\omega_{cj})_{j=1}^{m_c}\right)_{c\in\mathbf{c}} \in \left\{\alpha \underset{c\in\mathbf{c}}{\bigtimes} \omega_c\right\} = \boldsymbol{\Theta}$, where:

    - $\alpha = (\alpha_1,\ldots,\alpha_{|\mathbf{c}|}) \in \boldsymbol{\alpha}$ is the weight vector, $\boldsymbol{\alpha} = \left\{(\alpha_1,\ldots,\alpha_{|\mathbf{c}|}) \mid \alpha_c \geq 0, \forall c \in \mathbf{c}, \ \sum_{c\in\mathbf{c}} \alpha_c = 1\right\}$
    - $\omega_c = (\omega_{c1},\ldots,\omega_{cm_c}) \in \boldsymbol{\omega}_c$ is the component-specific parameter vector and $m_c$ stands for the number of parameters of the $c$-th mixture component.

$\square$

The general learning algorithm summarized in the Chapter 2, Algorithm 1, uses the system model $M(O_t|\psi_t,\Theta)$. By inserting right-hand side of (3.3) into (3.2), the desired parametric system model is obtained:

$$M(O_t|\psi_t,\Theta) = \sum_{c\in\mathbf{c}} \frac{\alpha_c J_c(O_t,\psi_t|\omega_c)}{\sum_{d\in\mathbf{c}} \alpha_d \underbrace{\int_{\mathbf{O}} J_d(O_t,\psi_t|\omega_d)dO_t}_{W_d(\psi_t,\omega_d)}} = \sum_{c\in\mathbf{c}} \frac{\alpha_c J_c(O_t,\psi_t|\omega_c)}{\sum_{d\in\mathbf{c}} \alpha_d W_d(\psi_t,\omega_d)}$$

$$= \sum_{c\in\mathbf{c}} \underbrace{\frac{\alpha_c W_c(\psi_t,\omega_c)}{\sum_{d\in\mathbf{c}} \alpha_d W_d(\psi_t,\omega_d)}}_{w_c(\psi_t,\Theta_c)} \underbrace{\frac{J_c(O_t,\psi_t|\omega_c)}{W_c(\psi_t,\omega_c)}}_{M_c(O_t|\psi_t,\Theta_c)} = \sum_{c\in\mathbf{c}} w_c(\psi_t,\Theta_c) M_c(O_t|\psi_t,\Theta_c), \qquad (3.4)$$

where

- $W_c(\psi_t,\Theta_c) = \int_{\mathbf{O}} J_c(O_t,\psi_t|\omega_c)dO_t$

- $w_c(\psi_t,\Theta_c) = \frac{\alpha_c W_c(\psi_t,\Theta_c)}{\sum_{d\in\mathbf{c}} \alpha_d W_d(\psi_t,\omega_d)}$

- $M_c(O_t|\psi_t,\Theta_c) = \frac{J_c(O_t,\psi_t|\omega_c)}{W_c(\psi_t,\Theta_c)}$

Equation (3.4) shows, that the parametric model can be interpreted as the mixture with data-dependent weights, while the data-dependence is not arbitrary and *does not introduce new parameters.*

As majority of joint pds is approximated by mixture and model (3.4) is derived from it by deductive rules of probability theory, it "universally" approximate any probabilistic dynamic model operating on finite-dimensional data vector.

# Chapter 4

# Learning with Mixture Ratio Model

This chapters uses the general approximate learning algorithm proposed in Chapter 2 and summarized in Algorithm 1 for the parametric system model in the form (3.3).

First of all, the set of computationally feasible pds $\mathbf{P}$ has to be specified.

Each component $\exp \langle B(\phi_{c,t}) \, | \, C_c(\omega_c) \rangle$ in (3.3) depends only on a part of the parameter vector $\Theta$. For each component, following *conjugated prior pds* to the respective components[1], determined by sufficient statistics $V_{t;c}$ (see Definition 2), are considered:

$$P(\omega_c|D^t) = P(\omega_c|V_t) = P_t(\omega_c) = \frac{\exp \langle V_{t;c} \, | \, C_c(\omega_c) \rangle}{N_c(V_{t;c})}, \qquad V_{t;c} = (V_{t;c1}, \dots V_{t;cm_c}) \qquad (4.1)$$

$$N_c(V_{t;c}) = \int_{\omega_c} \exp \langle V_{t;c} \, | \, C_c(\omega_c) \rangle \, d\omega_c.$$

Contiguously, the component weights in (3.3) can be seen as a pd $M(c_t = c|\Theta) = \alpha_c$ of an unobserved pointer $c_t \in \mathbf{c}$ to an active component. This pd belongs to EF:

$$M(c_t|\Theta) = \exp \langle \delta(c_t) \, | \, \ln(\alpha) \rangle$$

$$\delta(c_t) = (\delta_{1c_t}, \dots, \delta_{|\mathbf{c}|c_t}) \quad \ln(\alpha) = (\ln(\alpha_1), \dots, \ln(\alpha_{|\mathbf{c}|})) \quad \langle \delta(c_t) \, | \, \ln(\alpha) \rangle = \sum_{c \in \mathbf{c}} \delta_{cc_t} \ln(\alpha_c),$$

where $\delta_{cc_t}$ is Kronecker delta function (see Definition 7).

The conjugated prior for this pd is Dirichlet distribution, determined by the sufficient vector statistic $v = v_t$, in the following form:

$$P(\alpha|D^t) = P(\alpha|v_t) = P_t(\alpha) = \frac{\exp \langle v_t - 1 \, | \, \ln(\alpha) \rangle}{Be(v_t)}, \qquad v_t = \left( v_{t;1}, \dots v_{t;|\mathbf{c}|} \right), \quad v_c > 0, \quad \forall c \in \mathbf{c}, \qquad (4.2)$$

where $Be(v_t)$ is the multi-dimensional Beta function (see Definition 6).

The marginal pds $P_t(\alpha)$ (4.2), $P_t(\omega_c)$, $c \in \mathbf{c}$ (4.1) do not determine the joint pd $P_t(\Theta)$ unambiguously, [5]. This allows freedom in choosing the set of feasible joint pds $\mathbf{P}$. Since the information about relations between subparts of the parameter vector $\Theta$ is missing, the preferable choice of the set $\mathbf{P}$ is the set of joint pds $P_t(\Theta)$ that have the product form of their marginals:

$$\mathbf{P} = \left\{ P_t(\Theta) \, \middle| \, P_t(\Theta) = P_t(\alpha) \prod_{c \in \mathbf{c}} P_t(\omega_c) \right\}. \qquad (4.3)$$

---

[1]The term *conjugated pd* means that the pd preserves its functional form during the updating via Bayes' rule (1.7). The pds $P_t(\omega_c), c \in \mathbf{c}$, here are conjugated to the respective components $J_c(O_t, \psi_t|\omega_c)$ from Definition 11.

Thus, a pd $P_t(\Theta) \in \mathbf{P}$ (4.3) has the following form:

$$P_t(\Theta) = \frac{\exp\langle v_t - 1 \,|\, \ln(\alpha)\rangle \prod_{c\in\mathbf{c}} \exp\langle V_{t;c} \,|\, C_c(\omega_c)\rangle}{\mathrm{Be}(v_t) \prod_{c\in\mathbf{c}} N_c(V_{t;c})} \tag{4.4}$$

Since each member of the set of feasible pds $\mathbf{P}$ (4.3) is determined by the statistics $V_t = (v_t, (V_{t;c})_{c\in\mathbf{c}})$, it remains to convert the general algorithm summarized in Algorithm 1 to an algorithm updating $V_{t-1}$ to $V_t$.

**Updating step of Algorithm 1** The first step in general learning algorithm (Algorithm 1) is to update the feasible pd $P_{t-1}(\Theta) \in \mathbf{P}$, determined by statistics $V_{t-1} = (v_{t-1}, (V_{t-1;c})_{c\in\mathbf{c}})$, to $\tilde{P}_t(\Theta)$ via Bayes' rule (1.7) as follows:

$$\tilde{P}_t(\Theta) \overset{(1.7)}{=} \frac{M(O_t|\psi_t, \Theta)P_{t-1}(\Theta)}{\int_\Theta M(O_t|\psi_t, \Theta)P_{t-1}(\Theta)d\Theta} \tag{4.5}$$

$$\overset{(3.4)}{\propto} \sum_{c\in\mathbf{c}} \frac{\alpha_c J_c(O_t, \psi_t|\omega_c)P_{t-1}(\Theta)}{\sum_{d\in\mathbf{c}} \alpha_d W_d(\psi_t, \omega_d)}$$

$$\overset{(4.4)}{\propto} \sum_{c\in\mathbf{c}} \frac{\alpha_c P_c(\phi_{t;c}^\mathbf{c})}{\sum_{d\in\mathbf{c}} \alpha_d W_d(\psi_t, \omega_d)} \exp\langle B_c(\phi_{t;c}) \,|\, C_c(\omega_c)\rangle \exp\langle v_{t-1} - 1 \,|\, \ln(\omega)\rangle \prod_{d\in\mathbf{c}} \exp\langle V_{t-1;d} \,|\, C_d(\omega_d)\rangle$$

$$= \sum_{c\in\mathbf{c}} \frac{P_c(\phi_{t;c}^\mathbf{c})}{\sum_{d\in\mathbf{c}} \alpha_d W_d(\psi_t, \omega_d)} \exp\langle v_{t-1} - 1 + \delta(c) \,|\, \ln(\omega)\rangle \prod_{d\in\mathbf{c}} \exp\langle V_{t-1;d} + \delta_{cd} B_c(\phi_{t;c}) \,|\, C_d(\omega_d)\rangle$$

$$= \sum_{c\in\mathbf{c}} \frac{\gamma_{t;c} P_c(\phi_{t;c}^\mathbf{c}) Q_{t;c}(\Theta)}{\sum_{d\in\mathbf{c}} \alpha_d W_d(\psi_t, \omega_d)},$$

where $Q_{t;c}(\Theta)$, $c \in \mathbf{c}$, are pds of $\Theta \in \mathbf{\Theta}$ that equal

$$Q_{t;c}(\Theta) = \frac{1}{\gamma_{t;c}} \exp\langle v_{t-1} - 1 + \delta(c) \,|\, \ln(\alpha)\rangle \prod_{d\in\mathbf{c}} \exp\langle V_{t-1;d} + \delta_{cd} B_c(\phi_{t;c}) \,|\, C_d(\omega_d)\rangle \tag{4.6}$$

$$\gamma_{t;c} = \mathrm{Be}(v_{t-1} + \delta_c) \prod_{d\in\mathbf{c}} N_d(V_{t-1;d} + \delta_{cd} B_c(\phi_{t;c})) \qquad\qquad c \in \mathbf{c}.$$

$\tilde{P}_t(\Theta)$ then can be expressed as follows:

$$\tilde{P}_t(\Theta) = \frac{H_t(\Theta) \sum_{c\in\mathbf{c}} \gamma_{t;c} P_c(\phi_{t;c}^\mathbf{c}) Q_{t;c}(\Theta)}{L_{-1}}, \qquad \text{where} \tag{4.7}$$

$$H_t(\Theta) = H(\psi_t, \Theta) = \frac{1}{\sum_{d\in\mathbf{c}} \alpha_d W_d(\psi_t, \omega_d)} \tag{4.8}$$

$$L_{-1} = \sum_{c\in\mathbf{c}} \gamma_{t;c} P_c(\phi_{t;c}^\mathbf{c}) \int_\mathbf{\Theta} H_t(\Theta) Q_{t;c}(\Theta) d\Theta \tag{4.9}$$

**Projecting step of Algorithm 1** The second step of the learning algorithm (Algorithm 1) is to find a projection of the pd $\tilde{P}_t$ to a pd $\hat{P}_t$ from the feasible set $\mathbf{P}$ determined by the statistic $\hat{V}_t = (\hat{v}_t, (\hat{V}_{t;c})_{c\in\mathbf{c}})$, $\hat{v}_t = (\hat{v}_{t;1}, ..., \hat{v}_{t;|\mathbf{c}|})$, $\hat{V}_{t;c} = (\hat{V}_{t;c1}, ..., \hat{V}_{t;cm_c})$, $c \in \mathbf{c}$.

Kerridge's inaccuracy (Definition 9) $\tilde{P}_t$ to $P \in \mathbf{P}$ reads:

$$K(\tilde{P}_t \| P) = -\int_{\Theta} \tilde{P}_t(\Theta)\ln(P(\Theta))d\Theta = -\int_{\Theta} \tilde{P}_t(\Theta)\ln\left(\frac{\exp\langle v-1 \,|\, \ln(\alpha)\rangle \prod_{c\in\mathbf{c}} \exp\langle V_c \,|\, C(\Theta_c)\rangle}{\mathrm{Be}(v)\prod_{c\in\mathbf{c}} N_c(V_c)}\right)d\Theta$$

$$= \ln(\mathrm{Be}(v)) + \sum_{c\in\mathbf{c}} \ln\left(N_c(V_c)\right) - \left\langle v-1 \,\Big|\, \int_{\Theta} \ln(\alpha)\tilde{P}_t(\Theta)d\Theta\right\rangle - \sum_{c\in\mathbf{c}}\left\langle V_c \,\Big|\, \int_{\Theta} C_c(\omega_c)\tilde{P}_t(\Theta)d\Theta\right\rangle$$

$$= \ln(\mathrm{Be}(v)) + \sum_{c\in\mathbf{c}} \ln\left(N_c(V_c)\right) - \sum_{c\in\mathbf{c}}(v_c - 1)\int_{\Theta} \ln(\alpha_c)\tilde{P}_t(\Theta)d\Theta$$

$$- \sum_{c\in\mathbf{c}}\sum_{j=1}^{m_c} V_{cj}\int_{\Theta} C_{cj}(\omega_c)\tilde{P}_t(\Theta)d\Theta. \tag{4.10}$$

Evaluation of Kerridge's inaccuracy (4.10) requires evaluation of the following integrals[2] (see (4.7)):

$$L_{0c} = \sum_{d\in\mathbf{c}} \gamma_{t;d}P_d(\phi_{t;d}^{\mathsf{c}})\int_{\Theta} \ln(\alpha_c)H_t(\Theta)Q_{t;d}(\Theta)d\Theta \tag{4.11}$$

$$L_{cj} = \sum_{d\in\mathbf{c}} \gamma_{t;d}P_d(\phi_{t;d}^{\mathsf{c}})\int_{\Theta} C_{cj}(\omega_c)H_t(\Theta)Q_{t;d}(\Theta)d\Theta. \tag{4.12}$$

Inserting them into (4.10) the statistics $\hat{V}_t$ determining the projection $\hat{P}_t(\Theta)$ is obtained as follows:

$$\hat{V}_t = \left(\hat{v}_t, (\hat{V}_{t;c})_{c\in\mathbf{c}}\right) = \operatorname*{argmin}_{\left(v, (V_c)_{c\in\mathbf{c}}\right)}\left[\ln(\mathrm{Be}(v)) + \sum_{c\in\mathbf{c}} \ln\left(N_c(V_c)\right) - \left\langle v-1 \,\Big|\, \frac{L_0}{L_{-1}}\right\rangle - \sum_{c\in\mathbf{c}}\left\langle V_c \,\Big|\, \frac{L_c}{L_{-1}}\right\rangle\right]$$

$$\hat{v}_t = \operatorname*{argmin}_{v}\left(\ln(\mathrm{Be}(v)) - \left\langle v-1 \,\Big|\, \frac{L_0}{L_{-1}}\right\rangle\right) \tag{4.13}$$

$$\hat{V}_{t;c} = \operatorname*{argmin}_{V_c}\left(\ln\left(N_c(V_c)\right) - \left\langle V_c \,\Big|\, \frac{L_c}{L_{-1}}\right\rangle\right) \qquad \forall c \in \mathbf{c} \tag{4.14}$$

Thus, $|\mathbf{c}| + 1$ minimization tasks are solved. The minimisation (4.14) depends on specific form of $N_c$ and will be done in Chapter 5 for Markov chains and in Chapter 6 for Gaussian components.

## 4.1 Numerical Approximation of L's

The function $H_t(\Theta)$ (4.8) depends on the whole parameter vector $\Theta \in \mathbf{\Theta}$ and the integrals in (4.9), (4.11) and (4.12) cannot be expressed analytically. Light tails of pds $Q_{t;d}(\Theta)$ make integrands in (4.9), (4.11) and (4.12) well localized, thus the approximation based on linear Taylor expansion of $H_t(\Theta)$ around expected values of the parameter vector $\Theta$ (according to $Q_{t;d}(\Theta)$, $d \in \mathbf{c}$) is used. The expected

---

[2]Please note, that $L_{-1}$ (4.9), $L_{0c}$ (4.11) and $L_{cj}$ (4.12) depend on time $t$. To simplify the notation, this dependence is not stressed.

values of parameters equals (see (4.6)):

$$\tilde{\alpha}^d_{t;c} = \int_\Theta \alpha_c Q_{t;d}(\Theta) d\Theta = \int_\alpha \alpha_c \frac{\exp\langle v_{t-1} - 1 + \delta(c) \mid \ln(\alpha)\rangle}{\mathrm{Be}(v_{t-1} + \delta_d)} d\alpha = \frac{v_{t-1;c} + \delta_{cd}}{\sum_{c\in\mathbf{c}} v_{t-1;c} + 1} \qquad c,d \in \mathbf{c} \qquad (4.15)$$

$$\tilde{\omega}^d_{t;cj} = \int_\Theta \omega_{cj} Q_{t;d}(\Theta) d\Theta = \int_{\omega_c} \omega_{cj} \frac{\exp\langle V_{t-1;d} + \delta_{cd} B_c(\phi_{t;c}) \mid C_d(\omega_d)\rangle}{N_d(V_{t-1;d} + \delta_{cd} B_c(\phi_{t;c}))} d\omega_c \quad c,d \in \mathbf{c}, \; j \in \{1,2,...,m_c\}$$
$$(4.16)$$

$$\tilde{\alpha}^d_t = (\tilde{\alpha}^d_{t;1}, \tilde{\alpha}^d_{t;2}, ..., \tilde{\alpha}^d_{t;|\mathbf{c}|}) \qquad \tilde{\omega}^d_{t;c} = (\tilde{\omega}^d_{t;c1}, \tilde{\omega}^d_{t;c2}, ..., \tilde{\omega}^d_{t;cm_c}) \qquad \tilde{\Theta}^d_t = (\tilde{\alpha}^d_t, (\tilde{\omega}^d_{t;c})_{c\in\mathbf{c}}).$$

Hence, the function $H(\psi_t, \Theta)$ is approximated as follows (depending on particular pd $Q_{t;d}(\Theta)$, $d \in \mathbf{c}$, used in the integration):

$$H_t(\Theta) \approx H_t(\tilde{\Theta}^d_t) + \nabla H_t(\tilde{\Theta}^d_t)(\Theta - \tilde{\Theta}^d_t)$$

$$= H_t(\tilde{\Theta}^d_t) + \sum_{c\in\mathbf{c}} \frac{\partial H_t(\tilde{\Theta}^d_t)}{\partial \alpha_c}(\alpha_c - \tilde{\alpha}^d_{t;c}) + \sum_{c\in\mathbf{c}} \sum_{j=1}^{m_c} \frac{\partial H_t(\tilde{\Theta}^d_t)}{\partial \omega_{cj}}(\omega_{cj} - \tilde{\omega}^d_{t;cj}), \qquad (4.17)$$

The $L$'s are then computed as follows ($\forall c \in \mathbf{c}, j \in \{1,..,m_c\}$):

$$L_{-1} = \sum_{d\in\mathbf{c}} \gamma_{t;d} P_d(\phi^\mathbf{c}_{t;d}) \int_\Theta \left( H_t(\tilde{\Theta}^d_t) + \nabla H_t(\tilde{\Theta}^d_t)(\Theta - \tilde{\Theta}^d_t) \right) Q_{t;d}(\Theta) d\Theta$$

$$= \sum_{d\in\mathbf{c}} \gamma_{t;d} P_d(\phi^\mathbf{c}_{t;d}) H_t(\tilde{\Theta}^d_t), \qquad (4.18)$$

$$L_{0c} = \sum_{d\in\mathbf{c}} \gamma_{t;d} P_d(\phi^\mathbf{c}_{t;d}) \int_\Theta \ln(\alpha_c) \left( H_t(\tilde{\Theta}^d_t) + \nabla H_t(\tilde{\Theta}^d_t)(\Theta - \tilde{\Theta}^d_t) \right) Q_{t;d}(\Theta) d\Theta$$

$$= \sum_{d\in\mathbf{c}} \gamma_{t;d} P_d(\phi^\mathbf{c}_{t;d}) \int_\Theta \ln(\alpha_c) \left( H_t(\tilde{\Theta}^d_t) + \sum_{e\in\mathbf{c}} \frac{\partial H_t(\tilde{\Theta}^d_t)}{\partial \alpha_e}(\alpha_e - \tilde{\alpha}^d_{t;e}) \right) Q_{t;d}(\Theta) d\Theta, \qquad (4.19)$$

$$L_{cj} = \sum_{d\in\mathbf{c}} \gamma_{t;d} P_d(\phi^\mathbf{c}_{t;d}) \int_\Theta C_{cj}(\omega_c) \left( H_t(\tilde{\Theta}^d_t) + \nabla H_t(\tilde{\Theta}^d_t)(\Theta - \tilde{\Theta}^d_t) \right) Q_{t;d}(\Theta) d\Theta$$

$$= \sum_{d\in\mathbf{c}} \gamma_{t;d} P_d(\phi^\mathbf{c}_{t;d}) \int_\Theta C_{cj}(\omega_c) \left( H_t(\tilde{\Theta}^d_t) + \sum_{i=1}^{m_c} \frac{\partial H_t(\tilde{\Theta}^d_t)}{\partial \omega_{ci}}(\omega_{ci} - \tilde{\omega}^d_{t;ci}) \right) Q_{t;d}(\Theta) d\Theta. \qquad (4.20)$$

## 4.2 Forgetting

The final step of the learning algorithm (Algorithm 1) is to compute forgetting factors $\lambda_t$, $\lambda_{t;c}$, $c \in \mathbf{c}$, and to obtain the final statistic $V_t = (v_t, (V_{t;c})_{c\in\mathbf{c}})$ determining the posterior pd $P_t(\Theta)$. It remains to convert the general formula from the third step in Algorithm 1 for the parametric model (3.4).

Work [6] suggests that the forgetting should be applied component-wise. Therefore, the updated $P_t(\Theta) = P_t(\alpha) \prod_{c \in \mathbf{c}} P_t(\omega_c)$ is obtained as follows:

$$P_t(\alpha) \propto \hat{P}_t^{\lambda_t}(\alpha) P_{t-1}^{1-\lambda_t}(\alpha) = \left( \exp \langle \hat{v}_t - 1 \mid \ln(\alpha) \rangle \right)^{\lambda_t} \left( \exp \langle v_{t-1} - 1 \mid \ln(\alpha) \rangle \right)^{1-\lambda_t}$$
$$= \exp \langle \lambda_t \hat{v}_t + (1 - \lambda_t) v_{t-1} - 1 \mid \ln(\alpha) \rangle$$

$$P_t(\omega_c) \propto \hat{P}_t^{\lambda_{t;c}}(\omega_c) P_{t-1}^{1-\lambda_{t;c}}(\omega_c) = \left( \exp \left\langle \hat{V}_{t;c} \mid C_c(\omega_c) \right\rangle \right)^{\lambda_{t;c}} \left( \exp \langle V_{t-1;c} \mid C_c(\omega_c) \rangle \right)^{1-\lambda_{t;c}}$$
$$= \exp \left\langle \lambda_{t;c} \hat{V}_{t;c} + (1 - \lambda_{t;c}) V_{t-1;c} - 1 \mid C_c(\omega_c) \right\rangle \qquad c \in \mathbf{c}.$$

Thus, the final statistic $V_t = (v_t, (V_{t;c})_{c \in \mathbf{c}})$ is computed as follows:

$$v_t = \lambda_t \hat{v}_t + (1 - \lambda_t) v_{t-1}$$
$$V_{t;c} = \lambda_{t;c} \hat{V}_{t;c} + (1 - \lambda_{t;c}) V_{t-1;c} \qquad c \in \mathbf{c}.$$

It remains to evaluate the forgetting factors $\lambda_t$, $\lambda_{t;c}$, $c \in \mathbf{c}$. They are computed according to the formula in the third step of Algorithm 1 as follows:

$$\lambda_t = \frac{\left( \int_{\Theta} M_t(\Theta) P_{t-1}(\Theta) d\Theta \right)^2}{\int_{\Theta} M_t^2(\Theta) P_{t-1}(\Theta) d\Theta} \tag{4.21}$$

$$\lambda_{t;c} = \frac{\left( \int_{\omega_c} M_{t;c}(\omega_c) P_{t-1}(\omega_c) d\omega_c \right)^2}{\int_{\omega_c} M_{t;c}^2(\omega_c) P_{t-1}(\omega_c) d\omega_c}, \tag{4.22}$$

where $M_t(\Theta) = M(O_t | \psi_t, \Theta)$ is the parametric model (3.4) and $M_{t;c}(\omega_c) = M_c(O_t | \psi_{t;c}, \omega_c)$ is the $c$-th component of the parametric model (3.4), which equals:

$$M_c(O_t | \psi_{t;c}, \omega_c) = \frac{J_c(O_t, \psi_{t;c} | \omega_c)}{\int_{\mathbf{O}} J_c(O_t, \psi_{t;c} | \omega_c) dO_t} = \frac{\exp \langle B_c(O_t, \psi_{t;c}) \mid C_c(\omega_c) \rangle}{\int_{\mathbf{O}} \exp \langle B_c(O_t, \psi_{t;c}) \mid C_c(\omega_c) \rangle dO_t}. \tag{4.23}$$

The denominator in (4.21) can be expressed by using the updated posterior pd $\tilde{P}_t(\Theta) = \frac{M_t(\Theta) P_{t-1}(\Theta)}{\int_{\Theta} M_t(\Theta) P_{t-1}(\Theta) d\Theta}$ obtained via Bayes' rule (4.5). The forgetting factor $\lambda_t$ then equals:

$$\lambda_t = \frac{\left( \int_{\Theta} M_t(\Theta) P_{t-1}(\Theta) d\Theta \right)^2}{\int_{\Theta} M_t(\Theta) P_{t-1}(\Theta) d\Theta \int_{\Theta} M_t(\Theta) \tilde{P}_t(\Theta) d\Theta} = \frac{\int_{\Theta} M_t(\Theta) P_{t-1}(\Theta) d\Theta}{\int_{\Theta} M_t(\Theta) \tilde{P}_t(\Theta) d\Theta} \approx \frac{\int_{\Theta} M_t(\Theta) P_{t-1}(\Theta) d\Theta}{\int_{\Theta} M_t(\Theta) \hat{P}_t(\Theta) d\Theta}, \tag{4.24}$$

where the last approximation $\approx$ arises by using $\hat{P}_t(\Theta)$ instead of $\tilde{P}_t(\Theta)$ [3].

---

[3] The use of $\hat{P}_t(\Theta)$ instead of $\tilde{P}_t(\Theta)$ is motivated by the wish to make evaluations computationally cheap.

The integrals in (4.24) are computed as follows:

$$\int_{\Theta} M_t(\Theta) P_{t-1}(\Theta) d\Theta = \frac{1}{\mathrm{Be}(v_{t-1}) \prod_{d \in \mathbf{c}} N_d(V_{t-1;d})} \int_{\Theta} \Bigg( H_t(\Theta) \sum_{c \in \mathbf{c}} \omega_c P_c(\phi_{t;c}^{\mathbf{c}}) \exp \langle B_c(\phi_{t;c}) \,|\, C_c(\omega_c) \rangle$$

$$\times \exp \langle v_{t-1} - 1 \,|\, \ln(\omega) \rangle \prod_{d \in \mathbf{c}} \exp \langle V_{t-1;d} \,|\, C_d(\omega_d) \rangle d\Theta \Bigg)$$

$$\overset{(4.6)}{=} \frac{1}{\mathrm{Be}(v_{t-1}) \prod_{d \in \mathbf{c}} N_d(V_{t-1;d})} \sum_{c \in \mathbf{c}} P_c(\phi_{t;c}^{\mathbf{c}}) \gamma_{t;c} \int_{\Theta} H_t(\Theta) Q_{t;c}(\Theta) d\Theta$$

$$\overset{4.9}{=} \frac{1}{\mathrm{Be}(v_{t-1}) \prod_{d \in \mathbf{c}} N_d(V_{t-1;d})} L_{-1}$$

$$\int_{\Theta} M_t(\Theta) \hat{P}_t(\Theta) d\Theta = \frac{1}{\mathrm{Be}(\hat{v}_t) \prod_{d \in \mathbf{c}} N_d(\hat{V}_{t;d})} \hat{L}_{-1},$$

where $\hat{L}_{-1}$ is computed by the same formula as (5.5) with newly computed statistic $\hat{V}_t = \left( \hat{v}_t, (\hat{V}_{t;c})_{c \in \mathbf{c}} \right)$ (4.13), (4.14).

Practically, $L_{-1}$ and $\hat{L}_{-1}$ are evaluated in the same way as (4.18). The forgetting factor (4.24) is computed as follows:

$$\lambda_t = \frac{\mathrm{Be}(\hat{v}_t) \prod_{d \in \mathbf{c}} N_d(\hat{V}_{t;d})}{\mathrm{Be}(v_{t-1}) \prod_{d \in \mathbf{c}} N_d(V_{t-1;d})} \frac{L_{-1}}{\hat{L}_{-1}} \approx \frac{\mathrm{Be}(\hat{v}_t) \prod_{d \in \mathbf{c}} N_d(\hat{V}_{t;d})}{\mathrm{Be}(v_{t-1}) \prod_{d \in \mathbf{c}} N_d(V_{t-1;d})} \frac{\sum_{d \in \mathbf{c}} \gamma_{t;d} P_d(\phi_{t;d}^{\mathbf{c}}) H_t(\tilde{\Theta}_t^d)}{\sum_{d \in \mathbf{c}} \hat{\gamma}_{t;d} P_d(\phi_{t;d}^{\mathbf{c}}) H_t(\hat{\Theta}_t^d)}, \tag{4.25}$$

where

- $\hat{\gamma}_{t;c} = \mathrm{Be}(\hat{v}_t + \delta_c) \prod_{d \in \mathbf{c}} N_d(\hat{V}_{t;d} + \delta_{cd} B_c(\phi_{t;c}))$,

- $\hat{\Theta}_t^d = \left( \hat{\alpha}_t^d, (\hat{\omega}_{t;c}^d)_{c \in \mathbf{c}} \right)$ is expected value of parameter vector $\Theta$ computed in the same way as (4.15) and (4.16) with the statistic $\hat{V}_t = \left( \hat{v}_t, (\hat{V}_{t;c})_{c \in \mathbf{c}} \right)$.

Denoting $H_c(\psi_t, \omega_c) = \frac{1}{\int_{\mathbf{o}} J_c(O_t, \psi_{t;c} | \omega_c) dO_t}$ and using an analogical derivation as in (4.24) and (4.25), the remaining forgetting factors $\lambda_{t;c}, c \in \mathbf{c}$ are obtained as follows:

$$\lambda_{t;c} = \frac{\int_{\omega_c} M_{t;c}(\omega_c) P_{t-1}(\omega_c) d\omega_c}{\int_{\omega_c} M_{t;c}(\omega_c) \hat{P}_t(\omega_c) d\omega_c} = \frac{\frac{1}{N_c(V_{t-1;c})} \int_{\omega_c} H_c(\psi_t, \omega_c) \exp \langle B_c(\phi_t) \,|\, C_c(\omega_c) \rangle \exp \langle V_{t-1;c} \,|\, C_c(\omega_c) \rangle d\omega_c}{\frac{1}{N_c(\hat{V}_{t;c})} \int_{\omega_c} H_c(\psi_t, \omega_c) \exp \langle B_c(\phi_t) \,|\, C_c(\omega_c) \rangle \exp \langle \hat{V}_{t;c} \,|\, C_c(\omega_c) \rangle d\omega_c}$$

$$= \frac{\frac{1}{N_c(V_{t-1;c})} \int_{\omega_c} H_c(\psi_t, \omega_c) \exp \langle V_{t-1;c} + B_c(\phi_t) \,|\, C_c(\omega_c) \rangle d\omega_c}{\frac{1}{N_c(\hat{V}_{t;c})} \int_{\omega_c} H_c(\psi_t, \omega_c) \exp \langle \hat{V}_{t;c} + B_c(\phi_t) \,|\, C_c(\omega_c) \rangle d\omega_c}$$

$$\approx \frac{N_c(V_{t-1;c} + B_c(\phi_t))}{N_c(V_{t-1;c})} \frac{N_c(\hat{V}_{t;c})}{N_c(\hat{V}_{t;c} + B_c(\phi_t))} \frac{H_c(\psi_t, \tilde{\omega}_c^c)}{H_c(\psi_t, \hat{\omega}_c^c)}. \tag{4.26}$$

**Algorithm 2** (Learning algorithm for mixture ratio model). One complete step of the learning algorithm applied to the mixture ratio model is summarized as follows:

1. Compute the expected values of the parameter vector $\tilde{\Theta}_t^d$ via (4.15) and (4.16).

2. Evaluate $L_{-1}, L_{0c}, L_{ci}, c \in \mathbf{c}, i \in \{1, ..., m_c\}$ (4.9, 4.11, 4.12) via (4.18-4.20).

3. Perform minimization tasks (4.13, 4.14), see Chapter 5 and Chapter 5.

4. Apply forgetting with factors (4.25) and (4.26).

$\square$

**Remark 3.** Please note, that the dynamic forgetting factors (4.25) and (4.25) derived in this section will eventually not be used, instead the following algorithms will use fixed ones $\lambda_t = \lambda_{t;c} = 1$, $c \in \mathbf{c}$. The reason is that the proposed forgetting was tested in the research project in simulation learning of Markov chain and it did not provide sufficiently good results. The proper choice of the forgetting factor is subject to an ongoing research which is out of the scope of this work. $\square$

# Chapter 5

# Mixture of Markov Chain Components

One of the most important mixtures with the components from the exponential family are Markov chain mixtures, [18]. A Markov chain operates with discrete valued observations $O \in \mathbf{O} \subset \mathbb{N}$ as well as discrete valued regression vectors $\psi_t \in \psi \subset \mathbb{N}^n$, the joint pds $J_c(\psi_t|\omega_c)$, $J(\psi_t|\Theta)$ as well as the parametric system model $M(O_t|\psi_t, \Theta)$ are pfs (see Definition 11, (3.4)).

This chapter specialises the general mixture ratio model (Chapter 3) and its learning (Chapter 4) for the mixture ratio of Markov chain components with conjugated Dirichlet pd $P_t(\omega_c)$, $c \in \mathbf{c}$. The specialisation allows to work with specific function forms and to perform minimisation of Kerridge's inaccuracy used in projections of $\tilde{P}$ on $\hat{P}$.

**Definition 12** (Joint probability of Markov chain mixture model). The joint probability of a mixture of Markov chain models is defined as follows:

$$J(O_t, \psi_t|\Theta) = J(\phi_t|\Theta) = \sum_{c \in \mathbf{c}} \alpha_c J_c(O_t, \psi_{t;c}|\omega_c) P_c(\psi_c^{\mathsf{c}}) = \sum_{c \in \mathbf{c}} \alpha_c \chi_{\phi_c}(\phi_t) \frac{1}{|\psi_c^{\mathsf{c}}|} \prod_{j=1}^{m_c} \left(\frac{\omega_{cj}}{\tilde{K}_{cj}}\right)^{\Delta_{cj}(O_t, \psi_{t;c})} \tag{5.1}$$

where

- $\mathbf{c}$ is the set of component indices $c \in \mathbf{c}$, the symbol $|\mathbf{c}| < \infty$ stands for its cardinality, $\mathbf{c} = \{1, 2, .., |\mathbf{c}|\}$,

- $\chi_{\phi_c}(\phi_t)$ is the characteristic function of the set $\phi_c$; let $\phi$ be a set of all possible values of the data-vector $\phi \in \phi$, then $\phi_c$ is the subset of $\phi$ modelled by $c$-th component in the mixture model (5.1); the set $\phi_c$ can be rewritten as $\phi_c = \mathbf{O} \times \psi_c$, where $\mathbf{O}$ and $\psi_c$ are the sets of all possible values of the observation $O$ and the component-specific regression vector $\psi_c$ respectively, the $\psi_c^{\mathsf{c}}$ denotes the complement of the set $\psi_c$ with respect to the set $\psi$,

- $J_c(O_t, \psi_{t;c}|\omega_c) = \chi_{\phi_c}(\phi_t) \prod_{j=1}^{m_c} \left(\frac{\omega_{cj}}{\tilde{K}_{cj}}\right)^{\Delta_{cj}(O_t, \psi_{t;c})}$ is the $c$-th component in the mixture model (5.1),

  $m_c \in \mathbb{N}$ denotes the number of parameters of the $c$-th component,

- the non-parametric probability $P_c(\psi_c^{\mathsf{c}})$ (3.3) is chosen as uniform on $\psi_c^{\mathsf{c}}$, i.e. $P_c(\psi_c^{\mathsf{c}}) = \frac{1}{|\psi_c^{\mathsf{c}}|}$,

- $\alpha_c$ is the weight of $c$-th component; the set of all possible values of the weight-vectors $\alpha = \left(\alpha_1, \alpha_2, ..., \alpha_{|\mathbf{c}|}\right) \in \alpha$ is defined as follows - $\alpha = \left\{(\alpha_1, ..., \alpha_{|\mathbf{c}|}) \mid \alpha_c \geq 0, \forall c \in \mathbf{c}, \sum_{c \in \mathbf{c}} \alpha_c = 1\right\}$,

- $\Theta$ denotes the parameter vector, $\Theta$ is the set of all possible values of $\Theta$ and they have the form:

$$\Theta = (\alpha, (\omega_c)_{c \in \mathbf{c}}) \qquad \Theta = \alpha \underset{c \in \mathbf{c}}{\times} \omega_c, \tag{5.2}$$

24

where $\omega_c = (\omega_{c1}, \omega_{c2}, ..., \omega_{cm_c})$ is the parameter vector specific for the $c$-th component of the mixture and $\boldsymbol{\omega}_c$ the set of all its possible values, $\boldsymbol{\omega}_c = \left\{ \omega_c \in [0,1]^{m_c} \mid \sum_{j=1}^{m_c} \omega_{cj} = 1 \right\}$,

- $\Delta_{cj}(O_t, \psi_{t;c})$ is the indicator function[1] for the parameter $\omega_{cj}$, generally it has the form of a combination (sums, multiplications and compositions) of the Kronecker delta functions on subparts of the vector $(O_t, \psi_{t;c})$; the indicator outputs 0 or 1 for each vector $(O_t, \psi_{t;c})$; the indicator vector $\Delta_c(O_t, \psi_{t;c}) = (\Delta_{c1}, \Delta_{c2}, ..., \Delta_{cm_c})$ outputs 1 on at most one position, 0 on all remaining positions.

- $\tilde{K}_{cj}$ is the normalizing constant belonging to the parameter $\omega_{cj}$, generally, $\tilde{K}_{cj}$ equals to the number of values of the vector $(O_t, \psi_{t;c})$ for which $\Delta_{cj}(O_t, \psi_{t;c})$ outputs 1, i.e.

$$\tilde{K}_{cj} = \sum_{(O_t, \psi_{t;c}) \in \mathbf{O} \times \boldsymbol{\psi}_c} \Delta_{cj}(O_t, \psi_{t;c}).$$

$\square$

**Remark 4.** Throughout the subsequent text, following simplified notation will be used:

The values of functions $\chi_{\boldsymbol{\phi}_c}(\phi_t)$ and $\Delta_{cj}(O_t, \psi_{t;c})$ depend on $\phi_t$ and $(O_t, \psi_{t;c})$ respectively, but their arguments will not be highlighted, i.e. $\chi_{\boldsymbol{\phi}_c}(\phi_t) = \chi_{\boldsymbol{\phi}_c}$, $\Delta_{cj}(O_t, \psi_{t;c}) = \Delta_{cj}$, $\quad \forall c \in \mathbf{c}, j \in \{1, ..., m_c\}$. Consequently, let $\forall c \in \mathbf{c}, j \in \{1, .., m_c\}$, $K_{cj} = |\boldsymbol{\psi}_c^{\mathsf{C}}| \tilde{K}_{cj}$, then the joint probability of the parametric model (5.1) can be rewritten as follows:

$$J(O_t, \psi_t | \Theta) = \sum_{c \in \mathbf{c}} \chi_{\boldsymbol{\phi}_c} \alpha_c \prod_{j=1}^{m_c} \left( \frac{\omega_{cj}}{K_{cj}} \right)^{\Delta_{cj}}.$$

$\square$

**Remark 5.** The corresponding parametric model of the observation $O_t \in \mathbf{O}$ (cf. Definition 10) has the following mixture ratio form:

$$M(O_t | \psi_t, \Theta) = M_t(\Theta) = \frac{\sum\limits_{c \in \mathbf{c}} \chi_{\boldsymbol{\phi}_c} \alpha_c \prod_{j=1}^{m_c} \left( \frac{\omega_{cj}}{K_{cj}} \right)^{\Delta_{cj}}}{\sum\limits_{O_t \in \mathbf{O}} \sum\limits_{c \in \mathbf{c}} \chi_{\boldsymbol{\phi}_c} \alpha_c \prod_{j=1}^{m_c} \left( \frac{\omega_{cj}}{K_{cj}} \right)^{\Delta_{cj}}} = H_t(\Theta) \sum_{c \in \mathbf{c}} \chi_{\boldsymbol{\phi}_c} \alpha_c \prod_{j=1}^{m_c} \left( \frac{\omega_{cj}}{K_{cj}} \right)^{\Delta_{cj}}, \qquad (5.3)$$

$$H_t(\Theta) = H(\psi_t, \Theta) = \frac{1}{\sum\limits_{O_t \in \mathbf{O}} \sum\limits_{c \in \mathbf{c}} \chi_{\boldsymbol{\phi}_c} \alpha_c \prod_{j=1}^{m_c} \left( \frac{\omega_{cj}}{K_{cj}} \right)^{\Delta_{cj}}}$$

$\square$

**Remark 6.** One practical possibility how to use mixture ratios for modelling high-order Markov chains is to model its joint probability $J(O_t, \psi_t)$ (5.1) as a mixture, where each component models dependence between output $O_t$ and one variable in regression vector $\psi_t$, i.e.

$$J(O_t, \psi_t) = \sum_{c=1}^{|\psi_t|} \alpha_c J(O_t, \psi_{t;c}) \frac{1}{|\boldsymbol{\psi}_c^{\mathsf{C}}|} \qquad \psi_t = \left( \psi_{t;1}, ..., \psi_{t;|\psi_t|} \right) \in \boldsymbol{\psi}_t$$

$\square$

---

[1] The use of $\Delta_{cj}$ instead of $B_{cj}$ (as in Definition 11) stresses that observations $O$ and regression vectors $\psi$ are discrete valued.

# Learning with Ratio of Markov Mixtures

The set of feasible posterior pds (see Algorithm 1 and (4.3)) is chosen as follows:

$$\mathbf{P} = \left\{ P_t(\alpha) \prod_{c \in \mathbf{c}} P_t(\omega_c) \right\},$$

where the component-conjugated pds $P_t(\alpha)$ and $P_t(\omega_c)$ are determined by finite-dimensional sufficient statistics $v_t$ and $V_{t;c}$ respectively, $\forall c \in \mathbf{c}$, and they are defined as follows (see (5.2) and Definition 6):

$$P_t(\alpha) = P(\alpha | v_t) = \frac{\prod_{c \in \mathbf{c}} \alpha_c^{V_{t;c}-1}}{\mathrm{Be}(v_t)} \qquad P_t(\omega_c) = P(\omega_c | V_{t;c}) = \frac{\prod_{j=1}^{m_c} \omega_{cj}^{V_{t;cj}-1}}{\mathrm{Be}(V_{t;c})} \qquad \forall c \in \mathbf{c},$$

where $\alpha \in \boldsymbol{\alpha}, \quad \omega_c \in \omega_c, \quad c \in \mathbf{c}, \quad V_t = (v_t, V_{t;1}, ..., V_{t;|\mathbf{c}|}), \, v_t > 0, V_{t;c} > 0, c \in \mathbf{c}.$

**Updating step of Algorithm 1**   The posterior pd $\tilde{P}_t(\Theta)$ (prior to its projection to $\mathbf{P}$) is calculated via Bayes' rule (1.7) as follows (cf. (4.5)):

$$\tilde{P}_t(\Theta) = \frac{M_t(\Theta) P_{t-1}(\Theta)}{\int_\Theta M_t(\Theta) P_{t-1}(\Theta)}$$

$$\propto H(\psi_t, \Theta) \sum_{c \in \mathbf{c}} \chi_{\phi_c} \alpha_c \prod_{j=1}^{m_c} \left( \frac{\omega_{cj}}{K_{cj}} \right)^{\delta_{cj}} \prod_{d \in \mathbf{c}} \alpha_d^{V_{t-1;d}-1} \prod_{j=1}^{m_d} \omega_{dj}^{V_{t-1;dj}-1}$$

$$= H(\psi_t, \Theta) \sum_{c \in \mathbf{c}} \chi_{\phi_c} \prod_{j=1}^{m_c} \left( \frac{1}{K_{cj}} \right)^{\Delta_{cj}} \prod_{d \in \mathbf{c}} \alpha_d^{V_{t-1;d}-1+\delta_{cd}} \prod_{j=1}^{m_d} \omega_{dj}^{V_{t-1;dj}-1+\delta_{cd}\Delta_{dj}}$$

$$= H(\psi_t, \Theta) \sum_{c \in \mathbf{c}} \chi_{\phi_c} \prod_{j=1}^{m_c} \left( \frac{1}{K_{cj}} \right)^{\Delta_{cj}} \gamma_{t;c} Q_{t;c}(\Theta),$$

where, $\forall c \in \mathbf{c}$, $Q_{t;c}(\Theta)$ is a pd of the parameter $\Theta \in \boldsymbol{\Theta}$,

$$Q_{t;c}(\Theta) = \frac{\prod_{d \in \mathbf{c}} \alpha_d^{V_{t-1;d}-1+\delta_{cd}} \prod_{j=1}^{m_d} \omega_{dj}^{V_{t-1;dj}-1+\delta_{cd}\Delta_{dj}}}{\gamma_{t;c}}$$

$$\gamma_{t;c} = \mathrm{Be}(v_{t-1} + \delta(c)) \prod_{d \in \mathbf{c}} \mathrm{Be}(V_{t-1;d} + \delta_{cd}\Delta_c) \tag{5.4}$$

Thus, $\tilde{P}_t(\Theta)$ equals:

$$\tilde{P}_t(\Theta) = \frac{H(\psi_t, \Theta) \sum_{c \in \mathbf{c}} \chi_{\phi_c} \prod_{j=1}^{m_c} \left( \frac{1}{K_{cj}} \right)^{\Delta_{cj}} \gamma_{t;c} Q_{t;c}(\Theta)}{L_{-1}},$$

where

$$L_{-1} = \sum_{c \in \mathbf{c}} \chi_{\phi_c} \gamma_{t;c} \prod_{j=1}^{m_c} \left( \frac{1}{K_{cj}} \right)^{\Delta_{cj}} \int_{\boldsymbol{\Theta}} H(\psi_t, \Theta) Q_{t;c}(\Theta) d\Theta \tag{5.5}$$

**Projecting step of Algorithm 1**   Subsequently, the projection of $\tilde{P}_t(\Theta)$ to $\hat{P}_t(\Theta) \in \mathbf{P}$ is chosen as the minimizer of Kerridge's inaccuracy, which reads (cf. (4.10)):

$$K(\tilde{P}_t\|P) = -\int_\Theta \tilde{P}_t(\Theta)\ln(P(\Theta))d\Theta = -\int_\Theta \tilde{P}_t(\Theta)\ln\left(\prod_{c\in\mathbf{c}}\frac{\alpha_c^{v_c-1}\prod_{j=1}^{m_c}\omega_{cj}^{V_{cj}-1}}{Be(v)Be(V_c)}\right)d\Theta$$

$$= \ln(Be(v)) + \sum_{c\in\mathbf{c}}\ln(Be(V_c)) - \sum_{c\in\mathbf{c}}(v_c-1)\int_\Theta \tilde{P}_t(\Theta)\ln(\alpha_c)d\Theta$$

$$- \sum_{c\in\mathbf{c}}\sum_{j=1}^{m_c}(V_{cj}-1)\int_\Theta \tilde{P}_t(\Theta)\ln(\omega_{cj})d\Theta$$

$$= \ln(Be(v)) + \sum_{c\in\mathbf{c}}\ln(Be(V_c)) - \sum_{c\in\mathbf{c}}(v_c-1)\frac{L_{0c}}{L_{-1}} - \sum_{c\in\mathbf{c}}\sum_{j=1}^{m_c}(V_{cj}-1)\frac{L_{cj}}{L_{-1}},$$

where, $\forall c \in \mathbf{c}$,

$$L_{0c} = \sum_{c\in\mathbf{c}}\chi_{\phi_c}\gamma_{t;c}\prod_{j=1}^{m_c}\left(\frac{1}{K_{cj}}\right)^{\Delta_{cj}}\int_\Theta \ln(\alpha_c)H(\psi_t,\Theta)Q_{t;c}(\Theta)d\Theta \tag{5.6}$$

$$L_{cj} = \sum_{c\in\mathbf{c}}\chi_{\phi_c}\gamma_{t;c}\prod_{j=1}^{m_c}\left(\frac{1}{K_{cj}}\right)^{\Delta_{cj}}\int_\Theta \ln(\omega_{cj})H(\psi_t,\Theta)Q_{t;c}(\Theta)d\Theta, \quad \forall j\in\{1,2...,m_c\}. \tag{5.7}$$

The projection $\hat{P}_t(\Theta) \in \mathbf{P}$ is fully determined by the vector statistic $V_t = (v_t,(V_{t;c})_{c\in\mathbf{c}})$. Thus, searching for the desired projection $\hat{P}_t(\Theta)$ is converted into following convex minimization tasks:

$$\hat{v}_t = \underset{v\in(0,+\infty)^{|\mathbf{c}|}}{\operatorname{argmin}}\left(\ln(Be(v)) - \sum_{c\in\mathbf{c}}(v_c-1)\frac{L_{0c}}{L_{-1}}\right) \tag{5.8}$$

$$\hat{V}_{t;c} = \underset{V_c\in(0,+\infty)^{m_c}}{\operatorname{argmin}}\left(\ln(Be(V_c)) - \sum_{j=1}^{m_c}(V_{cj}-1)\frac{L_{cj}}{L_{-1}}\right), \quad \forall c\in\mathbf{c}. \tag{5.9}$$

## Numerical Approximation of *L*'s

The projection is possible when knowing *L*'s (5.5), (5.6) and (5.7). Their calculation for the considered Markov-chain case is done here. The numerical approximation of them is based on linear Taylor expansions of the function $H_t(\Theta) = H(\psi_t,\Theta)$ (5.3) around expected values $\tilde{\Theta}_t^d = (\tilde{\alpha}_t^d,(\tilde{\omega}_{t;c}^d)_{c\in\mathbf{c}})$ according to pds $Q_{t;d}(\Theta)$, $\forall d\in\mathbf{c}$, i.e. (cf. Section 4.1):

$$\tilde{\alpha}_{t;c}^d = \int_\Theta \alpha_c Q_{t;d}(\Theta)d\Theta = \int_\alpha \alpha_c\frac{\prod_{e\in\mathbf{c}}\alpha_e^{v_{t-1;e}-1+\delta_{de}}}{Be(v_{t-1}+\delta_d)}d\alpha = \frac{v_{t-1;c}+\delta_{cd}}{\sum_{c\in\mathbf{c}}v_{t-1;c}+1} \quad c\in\mathbf{c}$$

$$\tilde{\omega}_{t;cj}^d = \int_\Theta \omega_{cj}Q_{t;d}(\Theta)d\Theta = \int_{\omega_c}\omega_{cj}\frac{\prod_{i=1}^{m_c}\omega_{ci}^{V_{t-1;ci}-1+\delta_{dc}\Delta_{ci}}}{Be(V_{t-1;c}+\delta_{dc}\Delta_c)}d\omega_c = \frac{V_{t-1;cj}+\delta_{dc}\Delta_{cj}}{\sum_{i=1}^{m_c}(V_{t-1;i}+\delta_{dc}\Delta_{ci})}$$

$$c\in\mathbf{c}, j\in\{1,2,...,m_c\}$$

$$\tilde{\alpha}_t^d = (\tilde{\alpha}_{t;1}^d,\tilde{\alpha}_{t;2}^d,...,\tilde{\alpha}_{t;|\mathbf{c}|}^d) \quad \tilde{\omega}_{t;c}^d = (\tilde{\omega}_{t;c1}^d,\tilde{\omega}_{t;c2}^d,...,\tilde{\omega}_{t;cm_c}^d) \quad \tilde{\Theta}_t^d = (\tilde{\alpha}_t^d,(\tilde{\omega}_{t;c}^d)_{c\in\mathbf{c}}). \tag{5.10}$$

Hence, the function $H_t(\Theta)$ is approximated as follows (depending on particular pd $Q_{t;d}(\Theta)$, $d \in \mathbf{c}$, used in integration):

$$H_t(\Theta) \approx H_t(\tilde{\Theta}_t^d) + \sum_{c \in \mathbf{c}} \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \alpha_c}(\alpha_c - \tilde{\alpha}_{t;c}^d) + \sum_{c \in \mathbf{c}} \sum_{j=1}^{m_c} \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \omega_{cj}}(\omega_{cj} - \tilde{\omega}_{t;cj}^d),$$

where
$$\frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \alpha_c} = -H_t^2(\tilde{\Theta}_t^d) \sum_{O_t \in \mathbf{O}} \chi_{\phi_c} \prod_{j=1}^{m_c} \left(\frac{\tilde{\omega}_{cj}^d}{K_{cj}}\right)^{\Delta_{cj}}$$

$$\frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \omega_{cj}} = -H_t^2(\tilde{\Theta}_t^d)\tilde{\alpha}_{t;c}^d \sum_{O_t \in \mathbf{O}} \chi_{\phi_c} \left(\frac{1}{K_{cj}}\right)^{\Delta_{cj}} 0^{1-\Delta_{cj}}.$$

The $L$'s are then computed as follows (see Appendix for detailed derivation):

$$L_{-1} = \sum_{c \in \mathbf{c}} \chi_{\phi_c} \gamma_{t;c} \prod_{j=1}^{m_c} \left(\frac{1}{K_{cj}}\right)^{\Delta_{cj}} H_t(\tilde{\Theta}_t^c), \tag{5.11}$$

$$L_{0d} = \sum_{c \in \mathbf{c}} \chi_{\phi_c} \gamma_{t;c} \prod_{j=1}^{m_c} \left(\frac{1}{K_{cj}}\right)^{\Delta_{cj}} \left[H_t(\tilde{\Theta}_t^c)\left(\psi(v_{t-1;d} + \delta_{dc}) - \psi\left(\sum_{e \in \mathbf{c}} v_{t-1;e} + 1\right)\right) \right. \tag{5.12}$$
$$\left. + \frac{1}{\sum_{e \in \mathbf{c}} v_{t-1;e} + 1} \sum_{\substack{e \in \mathbf{c} \\ e \neq d}} \tilde{\alpha}_{t;e}^c H_t^2(\tilde{\Theta}_t^c) \sum_{O_t \in \mathbf{O}} \left(\chi_{\phi_e} \prod_{j=1}^{m_e} \left(\frac{\tilde{\omega}_{t;ej}^c}{K_{ej}}\right)^{\Delta_{ej}} - \chi_{\phi_d} \prod_{j=1}^{m_d} \left(\frac{\tilde{\omega}_{t;dj}^c}{K_{dj}}\right)^{\Delta_{dj}}\right)\right], \quad d \in \mathbf{c},$$

$$L_{di} = \sum_{c \in \mathbf{c}} \chi_{\phi_c} \gamma_{t;c} \prod_{j=1}^{m_c} \left(\frac{1}{K_{cj}}\right)^{\Delta_{cj}} \left\{H_t(\tilde{\Theta}_t^c)\left(\psi(V_{t-1;di} + \delta_{dc}\Delta_{di}) - \psi\left(\sum_{j=1}^{m_d}(V_{t-1;dj} + \Delta_{dj})\right)\right) \right. \tag{5.13}$$
$$\left. + \frac{1}{\sum_{j=1}^{m_d}(V_{t-1;dj} + \Delta_{dj})} \sum_{\substack{j=1 \\ j \neq i}}^{m_i} \left[\tilde{\omega}_{t;dj}^c \tilde{\alpha}_{t;j}^c H_t^2(\tilde{\Theta}_t^c) \sum_{O_t \in \mathbf{O}} \left(\left(\frac{1}{K_{dj}}\right)^{\Delta_{dj}} 0^{1-\Delta_{dj}} - \left(\frac{1}{K_{di}}\right)^{\Delta_{di}} 0^{1-\Delta_{di}}\right)\right]\right\},$$
$$d \in \mathbf{c}, i \in \{1, ..., m_d\}.$$

**Remark 7.** Since the values of $\gamma_{t;c}$ can be near zero, it is not wise to use them directly for numerical calculations. For practical calculations, the minimization tasks (5.8) and (5.9) requires values of $\frac{L_{0c}}{L_{-1}}$ and $\frac{L_{cj}}{L_{-1}}$, $c \in \mathbf{c}, j \in (1, 2, .., m_c)$.

The solution for this problem is to use attributes of *beta* function ($\gamma_{t;c}$ is the product of *beta* functions, see (5.4)) and expressing the ratios via $\tilde{L}_{-1}$, $\tilde{L}_{0c}$ and $\tilde{L}_{cj}$ while keeping $\frac{L_{0c}}{L_{-1}} = \frac{\tilde{L}_{0c}}{\tilde{L}_{-1}}$ and $\frac{L_{cj}}{L_{-1}} = \frac{\tilde{L}_{cj}}{\tilde{L}_{-1}}$. $\tilde{L}_{-1}$, $\tilde{L}_{0c}$ and $\tilde{L}_{cj}$ then equals (see Appendix for detailed derivation): □

$$\tilde{L}_{-1} = \sum_{c \in \mathbf{c}} \chi_{\boldsymbol{\phi}_c} \tilde{\gamma}_{t;c} \mathrm{H}_t(\tilde{\Theta}_t^c), \tag{5.14}$$

$$\tilde{L}_{0d} = \sum_{c \in \mathbf{c}} \chi_{\boldsymbol{\phi}_c} \tilde{\gamma}_{t;c} \left[ \mathrm{H}_t(\tilde{\Theta}_t^c) \left( \psi(\mathrm{v}_{t-1;d} + \delta_{dc}) - \psi\left(\sum_{e \in \mathbf{c}} \mathrm{v}_{t-1;e} + 1\right) \right) \right. \tag{5.15}$$
$$\left. + \frac{1}{\sum_{e \in \mathbf{c}} \mathrm{v}_{t-1;e} + 1} \sum_{\substack{e \in \mathbf{c} \\ e \neq d}} \tilde{\alpha}_{t;e}^c \mathrm{H}_t^2(\tilde{\Theta}_t^c) \sum_{O_t \in \mathbf{O}} \left( \chi_{\boldsymbol{\phi}_e} \prod_{j=1}^{m_e} \left(\frac{\tilde{\omega}_{t;ej}^c}{K_{ej}}\right)^{\Delta_{ej}} - \chi_{\boldsymbol{\phi}_d} \prod_{j=1}^{m_d} \left(\frac{\tilde{\omega}_{t;dj}^c}{K_{dj}}\right)^{\Delta_{dj}} \right) \right], \quad d \in \mathbf{c},$$

$$\tilde{L}_{di} = \sum_{c \in \mathbf{c}} \chi_{\boldsymbol{\phi}_c} \tilde{\gamma}_{t;c} \left\{ \mathrm{H}_t(\tilde{\Theta}_t^c) \left( \psi(\mathrm{V}_{t-1;di} + \delta_{dc}\Delta_{di}) - \psi\left(\sum_{j=1}^{m_d}(\mathrm{V}_{t-1;dj} + \Delta_{dj})\right) \right) \right. \tag{5.16}$$
$$\left. + \frac{1}{\sum_{j=1}^{m_d}(\mathrm{V}_{t-1;dj} + \Delta_{dj})} \sum_{\substack{j=1 \\ j \neq i}}^{m_i} \left[ \tilde{\omega}_{t;dj}^c \tilde{\alpha}_{t;j}^c \mathrm{H}_t^2(\tilde{\Theta}_t^c) \sum_{O_t \in \mathbf{O}} \left( \left(\frac{1}{K_{dj}}\right)^{\Delta_{dj}} 0^{1-\Delta_{dj}} - \left(\frac{1}{K_{di}}\right)^{\Delta_{di}} 0^{1-\Delta_{di}} \right) \right] \right\},$$

$$d \in \mathbf{c}, i \in \{1, ..., m_d\},$$

$$\text{where} \quad \tilde{\gamma}_{t;c} = \tilde{\tilde{\alpha}}_{t;c} \prod_{j=1}^{m_c} \left(\frac{\tilde{\tilde{\omega}}_{cj}}{K_{cj}}\right)^{\Delta_{cj}} \qquad \tilde{\tilde{\alpha}}_{t;c} = \frac{\mathrm{v}_{t-1;c}}{\sum_{d \in \mathbf{c}} \mathrm{v}_{t-1;d}} \qquad \tilde{\tilde{\omega}}_{t;cj} = \frac{\mathrm{V}_{t-1;cj}}{\sum_{i=1}^{m_c} \mathrm{V}_{t-1;ci}}.$$

**Finding the projection**

The minimization tasks (5.8) and (5.9) are convex. The proof can be found e.g. in [28]. The convexity of tasks (5.8) and (5.9) implies its equivalence to finding the root of the derivatives of the minimized functions:

$$\hat{\mathrm{v}}_t = \arg_{\mathrm{v} \in (0,+\infty)^{|\mathbf{c}|}} \left( \begin{pmatrix} \psi(\mathrm{v}_1) - \psi\left(\sum_{c \in \mathbf{c}} \mathrm{v}_c\right) - \frac{\tilde{L}_{01}}{\tilde{L}_{-1}} \\ \vdots \\ \psi(\mathrm{v}_{|\mathbf{c}|}) - \psi\left(\sum_{c \in \mathbf{c}} \mathrm{v}_c\right) - \frac{\tilde{L}_{0|\mathbf{c}|}}{\tilde{L}_{-1}} \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} \right) \tag{5.17}$$

$$\hat{\mathrm{V}}_{t;c} = \arg_{\mathrm{V}_c \in (0,\infty)^{m_c}} \left( \begin{pmatrix} \psi(\mathrm{V}_{c1}) - \psi\left(\sum_{j=1}^{m_c} \mathrm{V}_{cj}\right) - \frac{\tilde{L}_{c1}}{\tilde{L}_{-1}} \\ \vdots \\ \psi(\mathrm{V}_{cm_c}) - \psi\left(\sum_{j=1}^{m_c} \mathrm{V}_{cj}\right) - \frac{\tilde{L}_{cm_c}}{\tilde{L}_{-1}} \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} \right) \quad \forall c \in \mathbf{c}, \tag{5.18}$$

where $\psi(x)$ is the digamma function (see Definition 5), $\tilde{L}_{-1}, \tilde{L}_{0c}$ and $\tilde{L}_{cj}, c \in \mathbf{c}, j \in \{1, ..., m_c\}$, are computed via (5.14) - (5.16). The equations (5.17) and (5.18) are solved numerically, e.g. MATLAB function *fsolve* can be used to solve them.

**Forgetting**

The forgetting relies on the general formulas (4.25) and (4.26) applied on the Markov chain parametric model (5.3). The Markov chain parametric model (5.3) implies the following particular forms of needed functions for computing the forgetting factors (4.25) and (4.26):

- $N_c(V_c) = \text{Be}(V_c)$,

- $\gamma_{t;c} = \text{Be}(v_{t-1} + \delta(c)) \prod_{d \in \mathbf{c}} \text{Be}(V_{t-1;d} + \delta_{cd}\Delta_c)$,

- $\hat{\gamma}_{t;c} = \text{Be}(\hat{v}_t + \delta(c)) \prod_{d \in \mathbf{c}} \text{Be}(\hat{V}_{t;d} + \delta_{cd}\Delta_c)$,

- $P_c\left(\psi_{t;c}^{\mathbf{c}}\right) = \frac{1}{|\psi_c^{\mathbf{c}}|}$,

- $H_c(\psi_t, \omega_c) = \dfrac{1}{\sum_{o_t \in \mathbf{o}} \chi_{\phi_c} \prod_{j=1}^{m_c} \left(\frac{\omega_{cj}}{K_{cj}}\right)^{\Delta_{cj}}}$.

The forgetting factors $\lambda_t, \lambda_{t;c}, c \in \mathbf{c}$ then equal:

$$\lambda_t = \frac{\text{Be}(\hat{v}_t) \prod_{d \in \mathbf{c}} \text{Be}(\hat{V}_{t;d})}{\text{Be}(v_{t-1}) \prod_{d \in \mathbf{c}} \text{Be}(V_{t-1;d})} \frac{\sum_{c \in \mathbf{c}} \chi_{\phi_c} \gamma_{t;c} \prod_{j=1}^{m_c} \left(\frac{1}{K_{cj}}\right)^{\Delta_{cj}} H_t(\tilde{\Theta}_t^c)}{\sum_{c \in \mathbf{c}} \chi_{\phi_c} \hat{\gamma}_{t;c} \prod_{j=1}^{m_c} \left(\frac{1}{K_{cj}}\right)^{\Delta_{cj}} H_t(\hat{\Theta}_t^c)} \tag{5.19}$$

$$= \frac{\sum_{c \in \mathbf{c}} \chi_{\phi_c} \frac{\text{Be}(v_{t-1}+\delta(c)) \prod_{d \in \mathbf{c}} \text{Be}(V_{t-1;d}+\delta_{cd}\Delta_c)}{\text{Be}(v_{t-1}) \prod_{d \in \mathbf{c}} \text{Be}(V_{t-1;d})} \prod_{j=1}^{m_c} \left(\frac{1}{K_{cj}}\right)^{\Delta_{cj}} H_t(\tilde{\Theta}_t^c)}{\sum_{c \in \mathbf{c}} \chi_{\phi_c} \frac{\text{Be}(\hat{v}_t+\delta(c)) \prod_{d \in \mathbf{c}} \text{Be}(\hat{V}_{t;d}+\delta_{cd}\Delta_c)}{\text{Be}(\hat{v}_t) \prod_{d \in \mathbf{c}} \text{Be}(\hat{V}_{t;d})} \prod_{j=1}^{m_c} \left(\frac{1}{K_{cj}}\right)^{\Delta_{cj}} H_t(\hat{\Theta}_t^c)} = \frac{\sum_{c \in \mathbf{c}} \chi_{\phi_c} \tilde{\tilde{\alpha}}_c \prod_{j=1}^{m_c} \left(\frac{\tilde{\tilde{\omega}}_{cj}}{K_{cj}}\right)^{\Delta_{cj}} H_t(\tilde{\Theta}_t^c)}{\sum_{c \in \mathbf{c}} \chi_{\phi_c} \hat{\alpha}_c \prod_{j=1}^{m_c} \left(\frac{\hat{\omega}_{cj}}{K_{cj}}\right)^{\Delta_{cj}} H_t(\hat{\Theta}_t^c)}$$

$$\lambda_{t;c} = \frac{\text{Be}(V_{t-1;c} + \Delta_c)}{\text{Be}(V_{t-1;c})} \frac{\text{Be}(\hat{V}_{t;c})}{\text{Be}(\hat{V}_{t;c} + \Delta_c)} \frac{H_c(\psi_t, \tilde{\omega}_c^c)}{H_c(\psi_t, \hat{\omega}_c^c)} = \frac{H_c(\psi_t, \tilde{\omega}_c^c) \prod_{j=1}^{m_c} (\tilde{\tilde{\omega}}_{cj})^{\Delta_{cj}}}{H_c(\psi_t, \hat{\omega}_c^c) \prod_{j=1}^{m_c} (\hat{\omega}_{cj})^{\Delta_{cj}}}, \tag{5.20}$$

where $\tilde{\Theta}_t^c = \left(\tilde{\alpha}_t^c, (\tilde{\omega}_{t;d}^c)_{d \in \mathbf{c}}\right)$, $c \in \mathbf{c}$, are expected values of the parameter vector $\Theta$ computed in (5.10), $\hat{\Theta}_t^c = \left(\hat{\alpha}_t^c, (\hat{\omega}_{t;d}^c)_{d \in \mathbf{c}}\right)$, $c \in \mathbf{c}$, are expected values of the parameter vector $\Theta$ computed analogically to (5.10), but with using the statistic $\hat{V}_t = \left(\hat{v}_t, (\hat{V}_{t;c})_{c \in \mathbf{c}}\right)$ instead of the statistic $V_{t-1} = (v_{t-1}, (V_{t-1;c})_{c \in \mathbf{c}})$,

$$\tilde{\tilde{\alpha}}_{t;c} = \frac{v_{t-1;c}}{\sum_{d \in \mathbf{c}} v_{t-1;d}} \qquad\qquad \tilde{\omega}_{t;cj} = \frac{V_{t-1;cj}}{\sum_{i=1}^{m_c} V_{t-1;ci}}$$

$$\hat{\tilde{\alpha}}_{t;c} = \frac{\hat{v}_{t;c}}{\sum_{d \in \mathbf{c}} \hat{v}_{t;d}} \qquad\qquad \hat{\omega}_{t;cj} = \frac{\hat{V}_{t;cj}}{\sum_{i=1}^{m_c} \hat{V}_{t;ci}}$$

**Remark 8.** Please note, that the dynamic forgetting factors (5.19) and (5.20) derived in this section will eventually not be used, instead the following algorithms will use fixed ones $\lambda_t = \lambda_{t;c} = 1$, $c \in \mathbf{c}$. The reason is the same as in Remark 3. $\square$

**Algorithm 3** (Learning algorithm for Markov Chain mixture ratio model). One learning step of the learning algorithm converted for the Markov Chain mixture ratio model is summarized (similarly to Remark 2) as follows:

1. Compute the expected values of the parameter vector $\tilde{\Theta}_t^c, c \in \mathbf{c}$ via (5.10).

2. Evaluate $L_{-1}, L_{0c}, L_{ci}, c \in \mathbf{c}, i \in \{1, ..., m_c\}$ (5.5, 5.6, 5.7) via (5.14-5.16).

3. Perform minimization tasks, i.e. solve equations (5.17, 5.18).

4. Apply forgetting with factors $\lambda_t = \lambda_{t;c} = 1, c \in \mathbf{c}$.

$\square$

# Chapter 6

# Mixture of Gaussian Components

Besides Markov chain mixtures, introduced in the Chapter 5, another important case of exponential-family mixture is a mixture of Gaussian (normal) components. They are used to model continuous valued observations $O \in \mathbf{O}$ depending on continuous valued regression vectors $\psi \in \boldsymbol{\psi}$.

This chapter will convert the general mixture ratio model (Chapter 3) as well as its learning (Chapter 4) for this specific form of components.

As described in the Chapter 3, the mixture-ratio model requires joint probability $J(O_t, \psi_t | \Theta)$. To avoid using multi-dimensional Gaussian components and, therefore, avoid computational difficulties, the joint probability $J(O_t, \psi_t | \Theta)$ is factorized as described in the following definition. Moreover, it opens the way to handling mixed (countinuous and discrete) data vectors.

**Definition 13.** (Joint probability of Gaussian mixture model) The joint probability (3.3) for the case of Gaussian components has the following form:

$$J(O, \psi | \Theta) = J(\phi | \Theta) = \sum_{c \in \mathbf{c}} \alpha_c J_c(\phi_c | \Theta_c) P_c(\psi_c^{\mathsf{c}}) = \sum_{c \in \mathbf{c}} \alpha_c \prod_{i=1}^{n_c} J_c(\phi_{ci} | \psi_{ci}, r_{ci}, \Theta_{ci}) P_c(\psi_c^{\mathsf{c}}),$$

$$J_c(\phi_{ci} | \psi_{ci}, r_{ci}, \Theta_{ci}) = \frac{1}{\sqrt{2\pi r_{ci}}} \exp\left(-\frac{\left(\phi_{ci} - \Theta_{ci}^T \psi_{ci}\right)^2}{2 r_{ci}}\right) \qquad c \in \mathbf{c}, \ i = 1, 2, ..., n_c, \qquad (6.1)$$

where

- $\mathbf{c}$ is the set of component indices $c \in \mathbf{c}$, the symbol $|\mathbf{c}| < \infty$ stands for its cardinality, $\mathbf{c} = \{1, 2, .., |\mathbf{c}|\}$,

- $\phi = (O, \psi)$ denotes a data vector; to simplify the notation, for $c \in \mathbf{c}$, the observation $O \in \mathbf{O}$ will be sometimes denoted as $O = \phi_{c1}$,

- $\alpha_c$ is the weight of $c$-th component; the set of all possible values of the weight-vectors $\alpha = \left(\alpha_1, \alpha_2, ..., \alpha_{|\mathbf{c}|}\right) \in \boldsymbol{\alpha}$ is defined as follows - $\boldsymbol{\alpha} = \left\{\alpha \subset [0,1]^{|\mathbf{c}|} \mid \sum_{c \in \mathbf{c}} \alpha_c = 1\right\}$,

- $\phi_c = (O, \psi_c)$, where $\psi_c$ is a component-specific sub-vector of $\psi$,

- $n_c$ denotes the length of the data vector $\phi_c$, i.e. $n_c = |\boldsymbol{\psi}_c|$,

- $\phi_{ci}$ denotes $i$-th entry of the data vector $\phi_c$, $\psi_{ci}$ denotes the sub-vector of $\phi_c$ containing entries from $(i+1)$ to $n_c$, i.e. $\psi_{ci} = \left(\phi_{c(i+1)}, \phi_{c(i+2)}, ..., \phi_{cn_c}\right)$, for $i < n_c$, $\psi_{cn_c}$ is an empty vector,

- $P_c(\psi_c^{\complement})$ is non-parametrized pd on the complement $\psi_c^{\complement}$ of the sub-vector $\psi_c$ with respect to $\psi$; it can be chosen e.g. as Gaussian distribution with zero mean and big variance as an expression of the lack information about these entries,

- $\Theta = \left(\alpha, (r_{ci}, \Theta_{ci})_{c\in\mathbf{c}, i=1}^{n_c}\right)$ denotes the parameter vector, the set of all possible values of the parameter vector, $\boldsymbol{\Theta}$, has the following form:

$$\boldsymbol{\Theta} = \alpha \bigtimes_{c\in\mathbf{c}} \bigtimes_{i=1}^{n_c} \mathbf{r} \bigtimes \boldsymbol{\Theta}_{ci}, \quad \text{where} \quad \mathbf{r} = (0, +\infty), \quad \boldsymbol{\Theta}_{ci} = \mathbb{R}^{n_c - i + 1}$$

$\square$

**Remark 9.** The pd (6.1) can be rewritten into the following form:

$$J(\phi_{ci}|\psi_{ci}, \Theta_{ci}) = \frac{1}{\sqrt{2\pi r_{ci}}} \exp\left\{-\frac{1}{2r_{ci}} \mathrm{tr}\left([\phi_{ci}, \psi_{ci}][\phi_{ci}, \psi_{ci}]^T [-1, \Theta_{ci}^T]^T [-1, \Theta_{ci}^T]\right)\right\},$$

which corresponds with the form of an EF component (Definition 3), proof can be found in [32]. $\square$

**Remark 10.** The corresponding parametric model of the observation $O \in \mathbf{O}$ (cf. Definition 10) has the following mixture ratio form:

$$M(O|\psi, \Theta) = \frac{\sum_{c\in\mathbf{c}} \alpha_c \prod_{i=1}^{n_c} J_c(\phi_{ci}|\psi_{ci}, \Theta_{ci}) P_c(\psi_c^{\complement})}{\sum_{c\in\mathbf{c}} \alpha_c P_c(\psi_c^{\complement}) \prod_{i=2}^{n_c} J_c(\phi_{ci}|\psi_{ci}, \Theta_{ci}) \underbrace{\int_{\mathbf{O}} J_c(O|\phi_c, \Theta_{ci}) dO}_{=1}}$$

$$= H(\Theta, \psi) \sum_{c\in\mathbf{c}} \alpha_c P_c(\psi_c^{\complement}) \prod_{i=1}^{n_c} \frac{1}{\sqrt{2\pi r_{ci}}} \exp\left(-\frac{\left(\phi_{ci} - \Theta_{ci}^T \psi_{ci}\right)^2}{2r_{ci}}\right), \qquad (6.2)$$

$$\text{where} \quad H(\Theta, \psi) = \frac{1}{\sum_{c\in\mathbf{c}} \alpha_c P_c(\psi_c^{\complement}) \prod_{i=2}^{n_c} \frac{1}{\sqrt{2\pi r_{ci}}} \exp\left(-\frac{(\phi_{ci} - \Theta_{ci}^T \psi_{ci})^2}{2r_{ci}}\right)} \qquad (6.3)$$

$\square$

## Learning with Ratio of Gaussian Mixtures

First of all, the set of feasible posterior pds (see Algorithm 1 and (4.3) $\mathbf{P}$ needs to be specified. It has the following product form:

$$\mathbf{P} = \left\{P_t(\alpha) \prod_{c\in\mathbf{c}} P_t(\Theta_c)\right\} = \left\{P_t(\alpha) \prod_{c\in\mathbf{c}} \prod_{i=1}^{n_c} P_t(\Theta_{ci})\right\}. \qquad (6.4)$$

**Agreement.** In this section, the indices determining the specific component and its part will be often dropped. The absence of this indices means the expression has the same form for each component and its part. For instance,

$$N(\nu, V) = \frac{(2\pi)^{0.5|\psi|}}{|V_\psi|^{0.5}} \left(\frac{\Lambda}{2}\right)^{-0.5\nu} \Gamma(0.5\nu)$$

32

substitutes

$$N(\nu_{ci}, V_{ci}) = \frac{(2\pi)^{0.5|\psi_{ci}|}}{|V_{\psi ci}|^{0.5}} \left(\frac{\Lambda_{ci}}{2}\right)^{-0.5\nu_{ci}} \Gamma(0.5\nu_{ci}) \qquad c \in \mathbf{c}, i \in \{1, 2, ..., n_c\}.$$

$\square$

Pds $P_t(\Theta_{ci})$ are chosen as conjugated pds to the respective components $J_c(\phi_{ci}|\psi_{ci}, \Theta_{ci})$ (see Definition 13). The conjugated prior to a Gaussian component is called Gauss inverse-Wishart distribution, [32]. It has the following form (cf. 4.1):

$$P(r, \Theta|\nu, V) = \frac{r^{-0.5(\nu+|\psi|+2)}}{N(\nu, V)} \exp\left\{-\frac{1}{2r} \text{tr}\left(V[-1, \Theta^T]^T[-1, \Theta^T]\right)\right\}. \tag{6.5}$$

For the practical computing, mainly because of numerical stability, which will be explained later on, the following alternative expression of Gauss inverse-Wishart pd is more suitable and will be used (conversion to this alternative expression is proved in [32]):

$$P(r, \Theta|\nu, V) = \frac{r^{-0.5(\nu+|\psi|+2)}}{N(\nu, V)} \exp\left\{-\frac{1}{2r}\left((\Theta - \hat{\Theta})^T V_\psi(\Theta - \hat{\Theta}) + \Lambda\right)\right\}, \tag{6.6}$$

$$N(\nu, V) = \frac{(2\pi)^{0.5|\psi|}}{|V_\psi|^{0.5}} \left(\frac{\Lambda}{2}\right)^{-0.5\nu} \Gamma(0.5\nu), \tag{6.7}$$

where

$$V = \begin{bmatrix} V_{\phi_1} & V_{\phi_1\psi}^T \\ V_{\phi_1\psi} & V_\psi \end{bmatrix}, \qquad \hat{\Theta} = V_\psi^{-1} V_{\phi_1\psi},$$

$$\Lambda = V_{\phi_1} - V_{\phi_1\psi}^T V_\psi^{-1} V_{\phi_1\psi}.$$

The matrix $V$ needs to be positive definite. Otherwise in (6.5), there are unbounded combinations of $\Theta_{ci}$ that prevent the discussed function be integrable. In practical situations, numerical implementation of the standard evolution $V_t = V_{t-1} + \psi\psi^T$ can easily cause the indefiniteness of $V$. This is overcomed by considering the Choleski decomposition of the inversion of $V$, [32], i.e.

$$V^{-1} = GG^T.$$

By partitioning $G$ to

$$G = \begin{bmatrix} G_{\phi_1} & 0 \\ G_{\phi_1\psi} & G_\psi \end{bmatrix}, \tag{6.8}$$

the statistics $V_\psi$, $\hat{\Theta}$ and $\Lambda$ used in (6.6) can be obtained from $G$ as follows:

$$V_\psi = \left(G_\psi^T\right)^{-1} \left(G_\psi\right)^{-1} \tag{6.9}$$

$$\hat{\Theta} = -\frac{1}{G_{\phi_1}} G_{\phi_1\psi} \tag{6.10}$$

$$\Lambda = \frac{1}{(G_{\phi_1})^2}. \tag{6.11}$$

Hence, instead of updating the statistic $V$, it suffices to update the statistic $G$. This ensures that the statistic $V$ remains positive definite even when computed on computers with bounded precision. The inversion in (6.9) is never directly computed as the updating runs directly on $G$. All of the formulas with

$V_\psi$ contains its inversion $V_\psi^{-1}$ or the determinant $|V_\psi|^{0.5}$. These can be computed directly using $G_\psi$ by (6.9) and

$$|V_\psi|^{0.5} = \left|\left(G_\psi\right)^{-1}\right| = \frac{1}{|G_\psi|} = \frac{1}{\prod_i \left[G_\psi\right]_{ii}},\tag{6.12}$$

respectively.

Also computing the determinant of the triangular matrix $G_\psi$ is a trivial operation as seen in (6.12), which is another advantage of evolving $G_t$ instead of $V_t$.

**Updating step of Algorithm 1** The update of the statistic $G_t$ to $G_{t+1}$ is done by the following algorithm called REFIL, [32]. It takes $G_t$ and $\phi_t$ as arguments and produces updated $G_{t+1}$.

---

**Algorithm 4** REFIL

---

1: $G_{t+1} = \text{REFIL}(G_t, \phi_t)$
2: $n = \text{length}(\phi_t)$
3: $\sigma_{n+1} = 1$
4: **for** $j = 1, \ldots n$ **do**
5: $\quad f_j = \sum_{i=1}^{j} G_{t;ij}\phi_i$
6: $\quad \sigma_j = \sqrt{\sigma_{[j+1]}^2 + f_{[j]}^2}$
7: $\quad G_{t+1;jj} = \frac{\sigma_{j+1}}{\sigma_j} f_j$
8: $\quad g_j = G_{t+1;jj} f_j$
9: $\quad$ **for** $i = j+1, \ldots n$ **do**
10: $\qquad G_{t+1;jj} = \frac{\sigma_{j+1}}{\sigma_j}\left(G_{t;ij} - \frac{f_j}{\sigma_{j+1}^2} g_i\right)$
11: $\qquad g_i = G_{t;ij} f_j + g_i$
12: $\quad$ **end**
13: **end**

---

The update of the statistic $G_t$ is computed in the first step of the general learning algorithm (Algorithm 1) to obtain the updated pd $\tilde{P}_t(\Theta) = \tilde{P}(\Theta|\tilde{V}_t)$ from the pd $P_{t-1}(\Theta)$ available from the previous learning step, where

$$\tilde{P}(\Theta|\tilde{V}_t) = \frac{H_t(\Theta) \sum_{c \in \mathbf{c}} \gamma_{t;c} Q_{t;c}(\Theta)}{L_{-1}}$$

$$Q_{t;c}(\Theta) = P_t(\alpha|\tilde{v}_t^c) \prod_{d \in \mathbf{c}} \prod_{i=1}^{n_d} P(r_{di}, \Theta_{di}|\tilde{v}_{t;di}^c, \tilde{G}_{t;di}^c)\tag{6.13}$$

$$\gamma_{t;c} = \text{Be}(\tilde{v}_t^c) \prod_{d \in \mathbf{c}} \prod_{i=1}^{n_d} N(\tilde{v}_{t;di}^c, \tilde{G}_{t;di}^c)/P_c(\psi_c^\mathbf{c})$$

$$N(\nu, G) = (2\pi)^{0.5|\psi|}|G_\psi|\left(\frac{\Lambda}{2}\right)^{-0.5\nu}\Gamma(0.5\nu)$$

$$L_{-1} = \int_\Theta H_t(\Theta) \sum_{c\in\mathbf{c}} \gamma_{t;c} Q_{t;c}(\Theta) d\Theta$$

$$\tilde{V}_t = \left(\left(\tilde{v}^c_{t;d}\right)_{c,d\in\mathbf{c}}, \left(\tilde{v}^c_{t;di}, \tilde{G}^c_{t;di}\right)^{n_c}_{c,d\in\mathbf{c},i=1}\right).$$

The statistics $\tilde{v}^c_{t;d}$ and $\tilde{v}^c_{t;di}$, $c, d \in \mathbf{c}, i = 1, ..., n_c$ are obtained by the following recursion obtained from the Bayes' rule (1.7):

$$\tilde{v}^c_{t;d} = v_{t-1;d} + \delta_{cd} \tag{6.14}$$

$$\tilde{v}^c_{t;di} = v_{t-1;di} + \delta_{cd}, \tag{6.15}$$

which is the same for both $\tilde{v}^c_{t;d}$ and $\tilde{v}^c_{t;di}$, but the initial conditions are generally (quite) different.

The recursion for the statistics $\tilde{G}^c_{t;di}$ is done by the mentioned REFIL algorithm (Algorithm 4)[1].

**Projecting step of Algorithm 1** The second step in the learning algorithm (Algorithm 1) is the projection of $\tilde{P}_t(\Theta)$ to the set of feasible pds $\mathbf{P}$ (6.4). As mentioned earlier in Algorithm 1, the projection is found as minimizer of the Kerridge inaccuracy (Definition 9), which is given as follows:

$$K(\tilde{P}\|P) = -\int_\Theta \tilde{P}_t(\Theta)\ln(P(\Theta))\,d\Theta = \ln(Be(v)) + \sum_{c\in\mathbf{c}}\sum_{j=1}^{n_c}\ln(N(\nu_{cj}, G_{cj})) - \sum_{c\in\mathbf{c}}(v_c - 1)\frac{L_{0c}}{L_{-1}} - \sum_{c\in\mathbf{c}}\sum_{i=1}^{n_c}\frac{L_{ci}}{L_{-1}}, \tag{6.16}$$

where

$$L_{-1} = \int_\Theta H_t(\Theta)\sum_{d\in\mathbf{c}}\gamma_{t;d}Q_{t;d}(\Theta)d\Theta \tag{6.17}$$

$$L_{0c} = \int_\Theta \ln(\alpha_c)H_t(\Theta)\sum_{d}\gamma_{t;d}Q_{t;d}(\Theta)d\Theta \tag{6.18}$$

$$L_{ci} = \int_\Theta \left(-\frac{1}{2}(\nu_{ci} + |\psi_{ci}| + 2)\ln(r_{ci}) - \frac{1}{2r_{ci}}\left((\Theta_{ci} - \hat{\Theta}_{ci})^T V_{\psi ci}(\Theta_{ci} - \hat{\Theta}_{ci}) + \Lambda_{ci}\right)\right)H_t(\Theta)\sum_{d\in\mathbf{c}}\gamma_{t;d}Q_{t;d}(\Theta)d\Theta$$

$$= -\frac{1}{2}(\nu_{ci} + |\psi_{ci}| + 2)L^1_{ci} - 0.5\Lambda_{ci}L^2_{ci} - 0.5L^3_{ci} \tag{6.19}$$

$$L^1_{ci} = \int_\Theta \ln(r_{ci})H_t(\Theta)\sum_{d\in\mathbf{c}}\gamma_{t;d}Q_{t;d}(\Theta)d\Theta$$

$$L^2_{ci} = \int_\Theta \frac{1}{r_{ci}}H_t(\Theta)\sum_{d\in\mathbf{c}}\gamma_{t;d}Q_{t;d}(\Theta)d\Theta$$

$$L^3_{ci} = \int_\Theta \frac{1}{r_{ci}}\left((\Theta_{ci} - \hat{\Theta}_{ci})^T V_{\psi ci}(\Theta_{ci} - \hat{\Theta}_{ci})\right)H_t(\Theta)\sum_{d\in\mathbf{c}}\gamma_{t;d}Q_{t;d}(\Theta)d\Theta$$

The projection $\hat{P}_t(\Theta) \in \mathbf{P}$ is determined by the statistics $(v_t, V_t) = \left(v_t, (\nu_{t;ci}, G_{t;ci})^{n_c}_{c\in\mathbf{c},i=1}\right)$, thus, the Kerridge inaccuracy (6.16) can be viewed as a function of the statistics $(v, V)$, i.e. $K(\tilde{P}\|P) = K(v, V)$. The

---

[1]In fact, only statistics $\tilde{G}^c_{t;ci}$ are updated, remaining statistics $\tilde{G}^c_{t;di}$, where $c \neq d$, remains the same, i.e. $\tilde{G}^c_{t;di} = G_{t-1;di}$

statistics determining the feasible projection $\hat{P}_t(\Theta) \in \mathbf{P}$ is then found as follows:

$$(\mathrm{v}_t, V_t) = \underset{\mathrm{v},V}{\mathrm{argmin}}\; \mathrm{K}(\mathrm{v}, V). \tag{6.20}$$

## Numerical Approximation of $L$'s

The projection is possible when knowing $L$'s (6.17), (6.18), (6.19). Their calculation for the considered Markov-chain case is done here. The numerical approximation of them is based on linear Taylor expansions of the function $\mathrm{H}_t(\Theta) = \mathrm{H}(\psi_t, \Theta)$ (6.3) around expected values $\hat{\Theta}_t^d = \left(\hat{\alpha}_t^d, (\hat{r}_{t;ci}^d, \hat{\Theta}_{t;ci}^d)_{i=1,c\in\mathbf{c}}^{n_c}\right)$ according to pds $Q_{t;d}(\Theta)$ (6.13), $\forall d \in \mathbf{c}$, i.e. (cf. Section 4.1):

$$\hat{\alpha}_{t;c}^d = \int_{\Theta} \alpha_c Q_{t;d}(\Theta)d\Theta = \int_{\alpha} \alpha_c \frac{\prod_{e\in\mathbf{c}} \alpha_e^{\tilde{v}_{t;e}-1}}{\mathrm{Be}(\tilde{v}_{t-1})}d\alpha = \frac{\tilde{v}_{t;c}}{\sum_{d\in\mathbf{c}}\tilde{v}_{t;d}} \quad c \in \mathbf{c} \tag{6.21}$$

$$\hat{r}_{t;ci}^d = \int_{\Theta} r_{ci} Q_{t;d}(\Theta)d\Theta = \int_{r_{ci}} r_{ci} \frac{r_{ci}^{-0.5(\tilde{v}_{t;ci}^d+2)}}{\mathrm{N}(\tilde{\Lambda}_{t;ci}^d, \tilde{v}_{t;ci}^d)}\exp\left(-\frac{\tilde{\Lambda}_{t;ci}^d}{2r_{ci}}\right)dr_{ci} = \frac{\tilde{\Lambda}_{t;ci}^d}{\tilde{v}_{t;ci}^d - 2} \tag{6.22}$$

$$\hat{\Theta}_{t;ci}^d = \hat{\tilde{\Theta}}_{t;ci}^d, \tag{6.23}$$

where

$$\hat{\Theta}_{t;ci}^d = \left(\hat{\Theta}_{t;ci1}^d, ..., \hat{\Theta}_{t;ci|\psi_{ci}|}^d\right).$$

Hence, the function $\mathrm{H}_t(\Theta)$ is approximated as follows (depending on particular pd $Q_{t;d}(\Theta)$, $d \in \mathbf{c}$, used in integration):

$$\mathrm{H}_t(\Theta) \approx \mathrm{H}_t(\hat{\Theta}_t^d) + \sum_{c\in\mathbf{c}} \frac{\partial \mathrm{H}_t(\hat{\Theta}_t^d)}{\partial\alpha_c}(\alpha_c - \hat{\alpha}_{t;c}^d) + \sum_{c\in\mathbf{c}}\sum_{j=1}^{n_c}\left(\frac{\partial \mathrm{H}_t(\hat{\Theta}_t^d)}{\partial r_{cj}}\left(r_{cj} - \hat{r}_{t;cj}^d\right) + \frac{\partial \mathrm{H}_t(\hat{\Theta}_t^d)}{\partial\Theta_{cj}}\left(\Theta_{cj} - \hat{\Theta}_{t;cj}^d\right)\right),$$

$$\frac{\partial \mathrm{H}_t(\hat{\Theta}_t^d)}{\partial\alpha_c} = -\mathrm{H}_t^2(\hat{\Theta}_t^d)\mathrm{P}_c(\psi_c^{\mathbf{c}})\prod_{i=2}^{n_c}\frac{1}{\sqrt{2\pi\hat{r}_{t;ci}^d}}\exp\left(-\frac{\left(\phi_{ci} - \left(\hat{\Theta}_{t;ci}^d\right)^T\psi_{ci}\right)^2}{2\hat{r}_{t;ci}^d}\right)$$

$$\frac{\partial \mathrm{H}_t(\hat{\Theta}_t^d)}{\partial r_{cj}} = -\mathrm{H}_t^2(\hat{\Theta}_t^d)\frac{\hat{\alpha}_{t;c}}{\left(\hat{r}_{t;cj}^d\right)^2}\left(\left(\phi_{cj} - \left(\hat{\Theta}_{t;cj}^d\right)^T\psi_{cj}\right)^2 - \hat{r}_{t;cj}^d\right)\prod_{i=2}^{n_c}\frac{1}{\sqrt{2\pi\hat{r}_{t;ci}^d}}\exp\left(-\frac{\left(\phi_{ci} - \left(\hat{\Theta}_{t;ci}^d\right)^T\psi_{ci}\right)^2}{2\hat{r}_{t;ci}^d}\right)$$

$$\frac{\partial \mathrm{H}_t(\hat{\Theta}_t^d)}{\partial\Theta_{cjk}} = -\mathrm{H}_t^2(\hat{\Theta}_t^d)\frac{\hat{\alpha}_{t;c}}{\hat{r}_{t;cj}^d}\left(\phi_{cj} - \left(\hat{\Theta}_{t;cj}^d\right)^T\psi_{cj}\right)\psi_{cjk}\prod_{i=2}^{n_c}\frac{1}{\sqrt{2\pi\hat{r}_{t;ci}^d}}\exp\left(-\frac{\left(\phi_{ci} - \left(\hat{\Theta}_{t;ci}^d\right)^T\psi_{ci}\right)^2}{2\hat{r}_{t;ci}^d}\right)$$

$$\frac{\partial \mathrm{H}_t(\hat{\Theta}_t^d)}{\partial\Theta_{cj}} = \left(\frac{\partial \mathrm{H}_t(\hat{\Theta}_t^d)}{\partial\Theta_{cj1}}, ..., \frac{\partial \mathrm{H}_t(\hat{\Theta}_t^d)}{\partial\Theta_{cj|\psi_{cj}|}}\right).$$

The $L's$ are computed as follows (proofs of the following formulas are found in Appendix):

$$L_{-1} \approx \sum_{d \in \mathbf{c}} \mathrm{H}(\hat{\Theta}^d) \gamma_{t;d} \tag{6.24}$$

$$L_{0c} \approx \sum_{d \in \mathbf{c}} \gamma_{t;d} \left[ \mathrm{H}(\hat{\Theta}^d) \left( \psi(\mathbf{v}_{t-1;d} + \delta_{cd}) - \psi\left( \sum_{e \in \mathbf{c}} \mathbf{v}_{t-1;e} + 1 \right) \right) \right.$$
$$\left. + \frac{1}{\sum_{e \in \mathbf{c}} \mathbf{v}_{t-1;f} + 1} \sum_{\substack{e \in \mathbf{c} \\ e \neq c}} \hat{\alpha}_{t;e}^d \left( \frac{\partial \mathrm{H}(\hat{\Theta}^d)}{\partial \alpha_c} - \frac{\partial \mathrm{H}(\hat{\Theta}^d)}{\partial \alpha_e} \right) \right] \tag{6.25}$$

$$L_{cj} = -0.5(\nu_{cj} + |\psi_{cj}| + 2)L_{cj}^1 - 0.5\Lambda_{cj}L_{cj}^2 - 0.5L_{cj}^3 \tag{6.26}$$

$$L_{cj}^1 \approx \sum_{d \in \mathbf{c}} \gamma_{t;d} \left[ \mathrm{H}(\hat{\Theta}_t^d) \left( \ln\left( 0.5\tilde{\Lambda}_{t;cj}^d \right) - \psi\left( 0.5\tilde{\nu}_{t;cj}^d \right) \right) \right] \tag{6.27}$$

$$L_{cj}^2 \approx \sum_{d \in \mathbf{c}} \gamma_{t;d} \left( \mathrm{H}(\hat{\Theta}_t^d) \frac{\tilde{\nu}_{cj}^d}{\tilde{\Lambda}_{t;cj}^d} + \frac{\partial \mathrm{H}(\hat{\Theta}_t^d)}{\partial r_{cj}} \frac{2}{\tilde{\nu}_{cj}^d + 2} \right) \tag{6.28}$$

$$L_{cj}^3 \approx \sum_{d \in \mathbf{c}} \gamma_{t;d} \left[ \mathrm{H}(\hat{\Theta}_t^d) \mathrm{tr}\left( \left( \tilde{V}_{\psi cj}^d \right)^{-1} V_{\psi cj} \right) \right.$$
$$+ \left( \hat{\tilde{\Theta}}_{t;cj}^d - \hat{\Theta}_{t;cj} \right)^T V_{\psi cj} \left( \hat{\tilde{\Theta}}_{t;cj}^d - \hat{\Theta}_{t;cj} \right) \left( \mathrm{H}(\hat{\Theta}_t^d) \frac{\tilde{\nu}_{cj}^d}{\tilde{\Lambda}_{t;cj}^d} + \frac{\partial \mathrm{H}(\hat{\Theta}_t^d)}{\partial r_{cj}} \frac{2}{\tilde{\nu}_{cj}^d + 2} \right)$$
$$\left. + \frac{\partial \mathrm{H}(\hat{\Theta}_t^d)}{\partial \Theta_{cj}} \left( \tilde{V}_{\psi cj}^d \right)^{-1} \left( V_{\psi cj} + V_{\psi cj}^T \right) \left( \hat{\tilde{\Theta}}_{t;cj}^d - \hat{\Theta}_{t;cj} \right) \right] \tag{6.29}$$

## Finding the projection

The minimization task (6.20) is convex and, therefore, it is equivalent to finding the root of the derivative of $\mathrm{K}(\mathrm{v}, V)$ (6.16), i.e. the statistics $(\mathrm{v}_t, V_t)$, determining the optimal projection, satisfies the following system of equations

$$\nabla \mathrm{K}(\mathrm{v}_t, V_t) = 0. \tag{6.30}$$

Some parts of the statistics $(\mathrm{v}_t, V_t) = \left( \mathrm{v}_t, (\nu_{t;ci}, G_{t;ci})_{c \in \mathbf{c}, i=1}^{n_c} \right)$ minimizing the Kerridge inaccuracy (6.16) can be found analytically. Note that the statistics $G_{t;ci}$ can be defined equivalently via $\left( \hat{\Theta}_{t;ci}, \Lambda_{t;ci}, V_{t;\psi ci}^{-1} \right)$, see (6.8), (6.10), (6.11) and (6.9).

**Minimizing statistics $\hat{\Theta}_{t;ci}$** First of all, $\left(\hat{\Theta}_{t;ci}\right)_{c\in\mathbf{c},i=1}^{n_c}$ are found as follows:

$$
\frac{\partial K(v,V)}{\partial \hat{\Theta}_{ci}} = \sum_{d\in\mathbf{c}} a_d \left(V_{\psi cj}+V_{\psi cj}^T\right)\left(\hat{\bar{\Theta}}_{t;ci}^d - \hat{\Theta}_{ci}\right) - \sum_{d\in\mathbf{c}} b_d\left(V_{\psi cj}+V_{\psi cj}^T\right)[1,...,1]^T
$$

$$
= -\left(V_{\psi cj}+V_{\psi cj}^T\right)\hat{\Theta}_{ci}\sum_{d\in\mathbf{c}} a_d + \left(V_{\psi cj}+V_{\psi cj}^T\right)\sum_{d\in\mathbf{c}}\left(a_d\hat{\bar{\Theta}}_{t;ci}^d - b_d[1,...,1]^T\right),
$$

where
$$
a_d = H(\hat{\Theta}_t^d)\frac{\tilde{v}_{cj}^d}{\tilde{\Lambda}_{t;cj}^d} + \frac{\partial H(\hat{\Theta}_t^d)}{\partial r_{cj}}\frac{2}{\tilde{v}_{cj}^d+2}
$$

$$
b_d = \left[b_{d1},...,b_{d(n_c-j+1)}\right] = \frac{\partial H(\hat{\Theta}_t^d)}{\partial \Theta_{cj}}\left(\tilde{V}_{\psi cj}^d\right)^{-1}.
$$

The statistics $\left(\hat{\Theta}_{t;ci}\right)_{c\in\mathbf{c},i=1}^{n_c}$ satisfying (6.30) are given as follows:

$$
\hat{\Theta}_{t;ci} = \frac{1}{\sum_{d\in\mathbf{c}} a_d}\sum_{d\in\mathbf{c}} b_d[1,...,1]^T + \frac{1}{\sum_{d\in\mathbf{c}} a_d}\sum_{d\in\mathbf{c}} a_d\hat{\bar{\Theta}}_{t;ci}^d \tag{6.31}
$$

**Minimizing statistics $V_{\psi ci}^{-1}$** To find the analytical formula for the statistics $\left(V_{\psi ci}^{-1}\right)_{c\in\mathbf{c},i=1}^{n_c}$, the formula for $L_{cj}^3$ (6.29) is rewritten as follows[2]:

$$
L_{cj}^3 = \sum_{k,l}\left[V_{\psi cj}\right]_{kl}\sum_{d}\gamma_{t;d}\Bigg\{H(\hat{\Theta}_t^d)\left[\tilde{V}_{\psi cj}^d\right]_{kl}^{-1}
$$

$$
+\left(H(\hat{\Theta}_t^d)\frac{\tilde{v}_{cj}^d}{\tilde{\Lambda}_{t;cj}^d} + \frac{\partial H(\hat{\Theta}_t^d)}{\partial r_{cj}}\frac{2}{\tilde{v}_{cj}^d+2}\right)\left(\left[\hat{\bar{\Theta}}_{cj}^d\right]_k - \left[\hat{\Theta}_{cj}\right]_k\right)\left(\left[\hat{\bar{\Theta}}_{cj}^d\right]_l - \left[\hat{\Theta}_{cj}\right]_l\right)
$$

$$
+\sum_m \frac{\partial H(\hat{\Theta}_t^d)}{\partial\left[\Theta_{cj}\right]_m}\left(\left[\tilde{V}_{\psi cj}^d\right]_{mk}^{-1}\left(\left[\hat{\bar{\Theta}}_{cj}^d\right]_l - \left[\hat{\Theta}_{cj}\right]_l\right) + \left[\tilde{V}_{\psi cj}^d\right]_{ml}^{-1}\left(\left[\hat{\bar{\Theta}}_{cj}^d\right]_k - \left[\hat{\Theta}_{cj}\right]_k\right)\right)\Bigg\}
$$

$$
= \sum_{k,l}\left[V_{\psi cj}\right]_{kl}\left[L_{cj}^3\right]_{kl},
$$

where
$$
\left[L_{cj}^3\right]_{kl} = \sum_d \gamma_{t;d}\Bigg\{H(\hat{\Theta}_t^d)\left[\tilde{V}_{\psi cj}^d\right]_{kl}^{-1}
$$

$$
+\left(H(\hat{\Theta}_t^d)\frac{\tilde{v}_{cj}^d}{\tilde{\Lambda}_{t;cj}^d} + \frac{\partial H(\hat{\Theta}_t^d)}{\partial r_{cj}}\frac{2}{\tilde{v}_{cj}^d+2}\right)\left(\left[\hat{\bar{\Theta}}_{cj}^d\right]_k - \left[\hat{\Theta}_{cj}\right]_k\right)\left(\left[\hat{\bar{\Theta}}_{cj}^d\right]_l - \left[\hat{\Theta}_{cj}\right]_l\right)
$$

$$
+\sum_m \frac{\partial H(\hat{\Theta}_t^d)}{\partial\left[\Theta_{cj}\right]_m}\left(\left[\tilde{V}_{\psi cj}^d\right]_{mk}^{-1}\left(\left[\hat{\bar{\Theta}}_{cj}^d\right]_l - \left[\hat{\Theta}_{cj}\right]_l\right) + \left[\tilde{V}_{\psi cj}^d\right]_{ml}^{-1}\left(\left[\hat{\bar{\Theta}}_{cj}^d\right]_k - \left[\hat{\Theta}_{cj}\right]_k\right)\right)\Bigg\}.
$$

---

[2]The symbol $\left[\tilde{V}_{\psi cj}^d\right]_{kl}^{-1}$ denotes the $k,l$-th entry of the matrix $\left(\tilde{V}_{\psi cj}^d\right)^{-1}$, similarly $\left[\hat{\Theta}_{cj}\right]_k$ denotes $k$-th entry of the vector $\hat{\Theta}_{cj}$. The alternative expression of $L_{cj}^3$ is obtained from (6.29) by using basic matrix and vector multiplication and by appropriate repositioning all of the elements.

Then, $\frac{\partial K(v,V)}{\partial [V_{\psi ci}]_{kl}}$ equals as follows:

$$\frac{\partial K(v,V)}{\partial \left[V_{\psi ci}\right]_{kl}} = -0.5\frac{1}{|V_{\psi ci}|}\frac{\partial}{\partial \left[V_{\psi ci}\right]_{kl}}|V_{\psi ci}| + 0.5\frac{\left[L_{ci}^2\right]_{kl}}{L_{-1}} = -0.5\left[V_{\psi ci}^{-1}\right]_{kl} + 0.5\frac{\left[L_{ci}^2\right]_{kl}}{L_{-1}}.$$

The matrix statistics $\left(V_{\psi ci}^{-1}\right)_{c\in\mathbf{c},i=1}^{n_c}$ satisfying (6.30) are given by the following formula:

$$\left[V_{\psi ci}^{-1}\right]_{kl} = \frac{\left[L_{ci}^2\right]_{kl}}{L_{-1}} \tag{6.32}$$

**Minimizing statistics $\Lambda_{t;ci}$**    The statistics $(\Lambda_{t;ci})_{c\in\mathbf{c},i=1}^{n_c}$ are found by computing the following derivative:

$$\frac{\partial K(v,V)}{\partial \Lambda_{ci}} = -0.5\frac{v_{t;ci}}{\Lambda_{ci}} + 0.5\frac{L_{ci}^2}{L_{-1}}.$$

The statistics $(\Lambda_{t;ci})_{c\in\mathbf{c},i=1}^{n_c}$ satisfying (6.30) then equals:

$$\Lambda_{t;ci} = v_{t;ci}\frac{L_{-1}}{L_{ci}^2} \tag{6.33}$$

**Minimizing statistics $v_{t;c}, v_{t;ci}$**    The remaining statistics $\left((v_{t;c})_{c\in\mathbf{c}},(v_{t;ci})_{c\in\mathbf{c},i=1}^{n_c}\right)$ are computed numerically and the specific equation for numerical solving are obtained from the following derivatives:

$$\frac{\partial K(v,V)}{\partial v_{ci}} = -0.5 - 0.5\ln\left(v_{t;ci}\frac{L_{-1}}{L_{cj}^2}\right) + 0.5\ln(2) + 0.5\Gamma(0.5v_{ci})\psi(0.5v_{ci}) + 0.5\frac{L_{ci}^1}{L_{-1}}$$

$$\frac{\partial K(v,V)}{\partial v_c} = \psi(v_c) - \psi\left(\sum_{d\in\mathbf{c}}v_d\right) - \frac{L_{0c}}{L_{-1}}.$$

The statistics $\left((v_{t;c})_{c\in\mathbf{c}},(v_{t;ci})_{c\in\mathbf{c},i=1}^{n_c}\right)$ are then found as a numerical solution to the following equations:

$$\Gamma(0.5v_{ci})\psi(0.5v_{ci}) - \ln(v_{ci}) + \frac{L_{ci}^1}{L_{-1}} + \ln\left(\frac{2L_{ci}^2}{L_{-1}}\right) - 1 = 0 \tag{6.34}$$

$$\begin{pmatrix} \psi(v_1) - \psi\left(\sum_{c\in\mathbf{c}}v_c\right) - \frac{L_{01}}{L_{-1}} \\ \vdots \\ \psi(v_{|\mathbf{c}|}) - \psi\left(\sum_{c\in\mathbf{c}}v_c\right) - \frac{L_{0|\mathbf{c}|}}{L_{-1}} \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} \tag{6.35}$$

The tasks (6.34) and (6.35) can be easily solved, e.g. by MATLAB function *fsolve*.

**Proposition 1.**  The statistics $\left(v_t,(v_{t;ci},G_{t;ci})_{c\in\mathbf{c},i=1}^{n_c}\right) = \left(v_t,\left(v_{t;ci},\hat{\Theta}_{ci},\Lambda_{ci},V_{\psi ci}^{-1}\right)_{c\in\mathbf{c},i=1}^{n_c}\right)$ determining the projection to the feasible set of pds $P_t(\Theta) \in \mathbf{P}$ (6.4) are obtained via formulas (6.31), (6.32), (6.33) and by solving equations (6.34) and (6.35).

$\square$

**Algorithm 5** (Learning algorithm for Gaussian mixture ratio model). One learning step of the learning algorithm converted for the Gaussian mixture ratio model is summarized (similarly to Remark 2) as follows:

1. Compute the expected values of the parameter vector $\hat{\Theta}_t^d = \left( \hat{\alpha}_t^d, (\hat{r}_{t;ci}^d, \hat{\Theta}_{t;ci}^d)_{i=1,c\in\mathbf{c}}^{n_c} \right), d \in \mathbf{c}$, via (6.21 - 6.23).

2. Evaluate $L_{-1}, L_{0c}, L_{ci}, c \in \mathbf{c}, i \in \{1,...,n_c\}$ (5.5, 5.6, 5.7) via (6.24 -6.26).

3. Find the updated statistics $\left( v_t, (v_{t;ci}, G_{t;ci})_{c\in\mathbf{c},i=1}^{n_c} \right) = \left( \mathrm{v}_t, \left( v_{t;ci}, \hat{\Theta}_{ci}, \Lambda_{ci}, V_{\psi ci}^{-1} \right)_{c\in\mathbf{c},i=1}^{n_c} \right)$ via (6.31), (6.32), (6.33) and by solving equations (6.34) and (6.35).

4. Apply forgetting with factors $\lambda_t = \lambda_{t;c} = 1, c \in \mathbf{c}$ (see Remark 3).

$\square$

The way how the mixture ratio model is constructed, i.e. by modelling joint probability on output and regression vector, allows to seamlessly model also mixed discrete-continuous systems. The practical modelling of such systems is complicated mainly if the dependence of discrete variable on continuous regression vector is present. However, if the joint probability is modelled, it can always be factorized and modelled in the way that overcomes this. It is shown in the following remark.

**Remark 11** (General mixed component). Let $\mathrm{J}_c(\phi_t|\Theta), \phi_t = (O_t, \psi_t)$, be one component of joint pd from mixture ratio model (3.3). Without loss of generality, suppose that the data vector $\phi_t$ contain two entries $\phi_t = (\phi_{t;C_1}, \phi_{t;C_2}, \phi_{t;D_1}, \phi_{t;D_2})$, where $\phi_{t;C_1}, \phi_{t;C_2}$ are continuous variables and $\phi_{t;D_1}, \phi_{t;D_2}$ are discrete variables. Such component has then the following form

$$\mathrm{J}_c(O_t, \psi_t|\Theta) = (\mathrm{J}_c(\phi_{t;C_1}|\phi_{t;C_2}\Theta_1)\mathrm{J}_c(\phi_{t;C_2}|\Theta_2))^{\delta(\phi_{t;D_1}, i_{t;1})\delta(\phi_{t;D_2}, i_{2;t})} \mathrm{J}_c(\phi_{t;D_1}, \phi_{t;D_2}|\Theta_3),$$

where

- $\mathrm{J}_c(\phi_{t;C_1}|\phi_{t;C_2}\Theta_1), \mathrm{J}_c(\phi_{t;C_2}|\Theta_2)$ are Gaussian pd from Chapter 6,

- $\mathrm{J}_c(\phi_{t;D_1}, \phi_{t;D_2}|\Theta_3)$ is Markov chain pf from Chapter 5,

- $\delta(\phi_{t;D_1}, i_{t;1}), \delta(\phi_{t;D_2}, i_{t;2})$ are Kronecker delta functions,

- $i_{t;1} \in \{1, 2, ..., |\boldsymbol{\phi}_{t;D_1}|\}, i_{t;2} \in \{1, 2, ..., |\boldsymbol{\phi}_{t;D_2}|\}$ are observed values of the variables $\phi_{t;D_1}, \phi_{t;D_2}$, respectively.

$\square$

# Chapter 7

# Experimental Part

This chapter experimentally examines properties of mixture ratio model on several examples. It is compared with standard mixture model, [18], in three scenarios: i) Modelling strengths of both models are compared by simulating a dynamic systems by these two models. ii) Monte Carlo study on learning quality of both models on simulation data are made. iii) Both models is tested on real futures trading data and their ability to serve for creating DM strategies is compared.

All of the experiments are made for the case of Markov chain components. Similar experiments with Gaussian and mixed components are under preparation and will be subject to the future research.

## 7.1 Mixture Ratio Modelling Strength

A dynamic system with system model M, generating real scalar observations $O_t$ with no actions and regression vector $\psi_t = O_{t-1}$ was simulated. The joint probability was the mixture of joint Gaussian pds $\mathcal{N}_c(\mu_c, \omega_c^{0.5})$ with expectations $\mu_c$ and square roots $\omega_c^{0.5}$ of precision matrices, $c = \{1,2\}$, which had the following form

$$J(O_t, O_{t-1}) = \underbrace{0.5}_{\alpha_1} \mathcal{N}_1 \underbrace{O_t, O_{t-1}}_{\phi_t} \left( \underbrace{\begin{bmatrix} 1 \\ 1 \end{bmatrix}}_{\mu_1}, \underbrace{\begin{bmatrix} 1 & 3 \\ 0 & 0.5 \end{bmatrix}}_{\omega_1^{0.5}} \right) + \underbrace{0.5}_{\alpha_2} \mathcal{N}_2 \underbrace{O_t, O_{t-1}}_{\phi_t} \left( \underbrace{\begin{bmatrix} -1 \\ -1 \end{bmatrix}}_{\mu_2}, \underbrace{\begin{bmatrix} 1 & -2 \\ 0 & 0.5 \end{bmatrix}}_{\omega_2^{0.5}} \right). \quad (7.1)$$

The system model then equalled as follows:

$$M(O_t|O_{t-1}) = \frac{0.5\mathcal{N}_{1O_t,O_{t-1}} + 0.5\mathcal{N}_{2O_t,O_{t-1}}}{0.5\underbrace{\int_{\mathbf{O}} \mathcal{N}_{1O_t,O_{t-1}} dO_t}_{W_1(O_{t-1})} + 0.5\underbrace{\int_{\mathbf{O}} \mathcal{N}_{2O_t,O_{t-1}} dO_t}_{W_2(O_{t-1})}} = w_1(O_{t-1})\frac{\mathcal{N}_{1O_t,O_{t-1}}}{W_1(O_{t-1})} + w_2(O_{t-1})\frac{\mathcal{N}_{2O_t,O_{t-1}}}{W_2(O_{t-1})}$$

$$w_1(O_{t-1}) = \frac{W_1(O_{t-1})}{W_1(O_{t-1}) + W_2(O_{t-1})} \qquad w_2(O_{t-1}) = \frac{W_2(O_{t-1})}{W_1(O_{t-1}) + W_2(O_{t-1})} \quad (7.2)$$

The typical simulation results, differing only in initial values $O_0 \in \{-0.8, 0, 3\}$, are shown in Figure 7.1.

**Discussion** The mixture ratio model demonstrates the dynamic dependence of the components weight on the data realisation. This is the key feature of the truly dynamic model, which obviously allows to model non-linear dynamic effects and unbalanced activations of respective components. For comparison,

the same mixture was simulated as the conditional one, i.e. with the constant component weights. It lacks both mentioned features of the mixture ratio model.
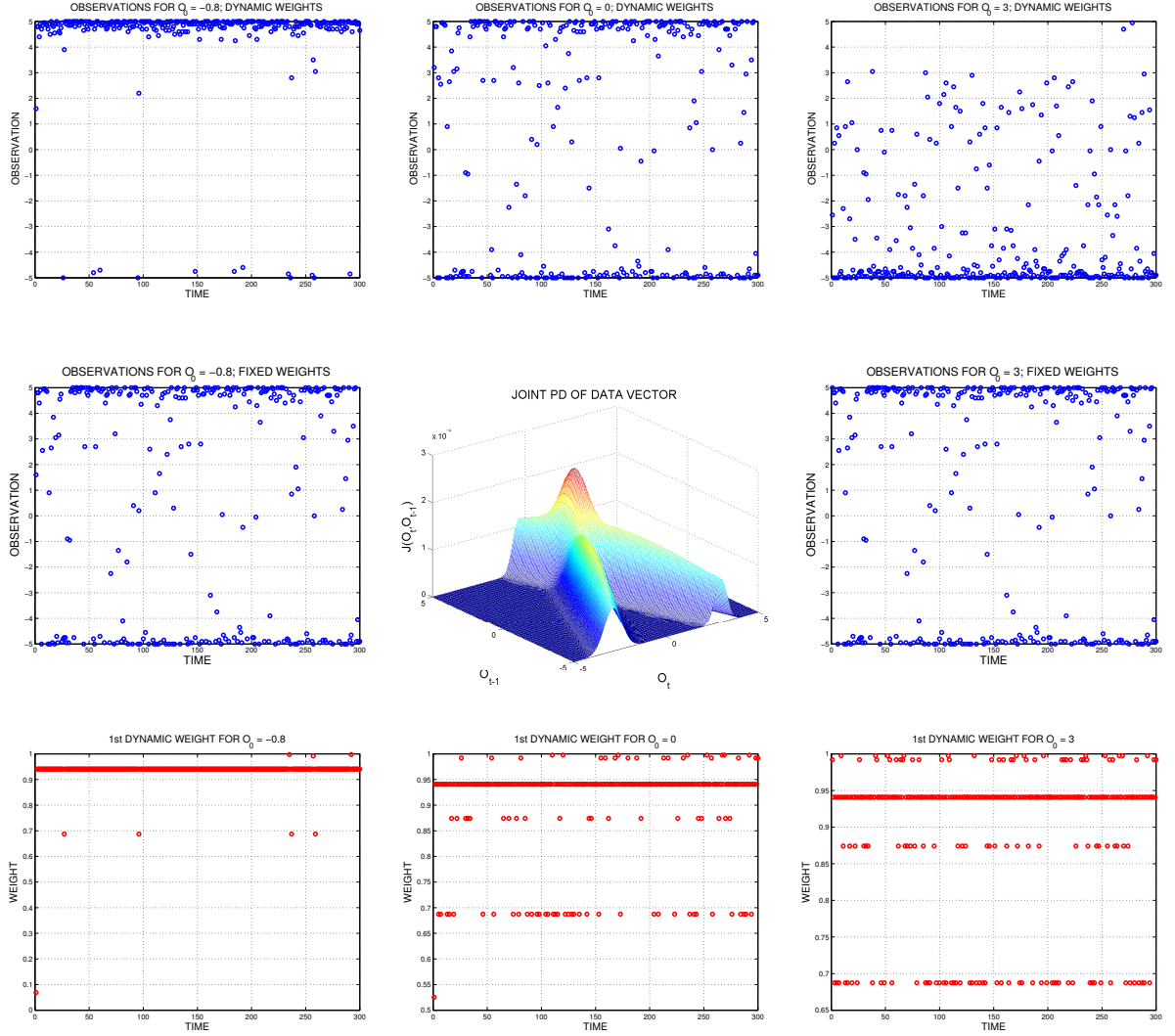


Figure 7.1: Simulation of system (7.1). Columns differ in the initial value of regression vector $\psi_1 = O_0 \in \{-0.8, 0, 3\}$. Pseudo-random realisations are the same in all cases. The 1st row shows the realized observations. The 2nd row shows realized observations when the same components as in (7.1) are used as conditional pds (constant component weights). The practically identical case for $O_0 = 0$ is replaced by the simulated joint pd. The 3rd row shows dynamic weights $w_1(O_{t-1})$ (7.2) corresponding to the 1st-row.

42

## 7.2 Learning Simulation with Mixture Ratio of Markov Chains

The simulation comparison of the Markov chain *mixture ratio model* discussed in Chapter 5 and the *standard* Markov chain *mixture model*, where the parametric conditional pf $M_C(O_t|\psi_t, \bar{\Theta})$ (compared to the joint pf $J(O_t, \psi_t|\Theta)$) is modelled, see [18], was made. The comparison was made as follows:

- A sequence of observations $(O_t)_{t=1}^{|\mathbf{t}|}$ was generated by the pf (5.3) with known parameters $\Theta_R$ (i.e. $P_R(O_t|\psi_t) = M(O_t|\psi_t, \Theta_R)$; $P_R(O_t|\psi_t)$ is referred as the *real model*), actions were not present in the system, $\psi_t = (O_{t-1}, ..., O_{t-n-1})$,

- both, the mixture ratio model and the standard mixture model were learnt recursively,

- the quality of both models during the simulation was compared using Kullback–Leibler divergence of pairs[1] of joint pfs $P(O_1, ..., O_t|\psi_0)$, $t \in \mathbf{t}$.

Let $P_J(O_t|\psi_t, \underline{D}^{t-1})$ and $P_C(O_t|\psi_t, \underline{D}^{t-1})$ be learnt predictors for the mixture ratio model and standard mixture model, respectively. The predictors depends on particular data realizations $\underline{D}^{t-1} = (\underline{O}_\tau)_{\tau=1}^{t-1}$ and are computed as follows:

$$P_N(O_t|\psi_t, \underline{D}^{t-1}) = \int_\Theta M_N(O_t|\psi_t, \Theta, \underline{D}^{t-1}) P(\Theta|\underline{D}^{t-1}) d\Theta, \qquad N \in \{J, C\}.$$

To simplify the notation, the dependence on particular data realization $\underline{D}^{t-1}$ will not be stressed, i.e. $P_J(O_t|\psi_t, \underline{D}^{t-1}) = P_J(O_t|\psi_t)$, $P_C(O_t|\psi_t, \underline{D}^{t-1}) = P_C(O_t|\psi_t)$.

Furthermore, let $P_R(O_1, ..., O_t|\psi_0)$, $P_J(O_1, ..., O_t|\psi_0)$, $P_C(O_1, ..., O_t|\psi_0)$, $t \in \mathbf{t}$, be joint pfs for the real model, the mixture ratio model and the standard mixture model respectively. They equal:

$$P_N(O_1, ..., O_t|\psi_0) = \prod_{\tau=1}^{t} P_N(O_\tau|\psi_\tau) P(\psi_0), \qquad N \in \{R, J, C\}, \tag{7.3}$$

with known $P(\psi_0)$. Kullback–Leibler divergences of the joint pfs then equal[2]:

$$KL^t(P_R \| P_N) = \sum_{O_1, ... O_t \in \mathbf{O}} \prod_{\tau=1}^{t} P_R(O_\tau|\psi_\tau) \ln\left(\prod_{\tau=1}^{t} \frac{P_R(O_\tau|\psi_\tau)}{P_N(O_\tau|\psi_\tau)}\right), \qquad N \in \{J, C\}, t \in \mathbf{t}. \tag{7.4}$$

(7.4) can be rewritten as follows:

$$KL^t(P_R \| P_L) = \sum_{\tau=1}^{t} \sum_{O_1, ... O_t \in \mathbf{O}} \prod_{\tau=1}^{t} P_R(O_\tau|\psi_\tau) \ln\left(\frac{P_R(O_\tau|\psi_\tau)}{P_N(O_\tau|\psi_\tau)}\right)$$

$$= \sum_{\tau=1}^{t} \sum_{\substack{O_\tau \in \mathbf{O} \\ \psi_\tau \in \psi}} P_R(O_\tau|\psi_\tau) P(\psi_\tau) \ln\left(\frac{P_R(O_\tau|\psi_\tau)}{P_N(O_\tau|\psi_\tau)}\right), \qquad N \in \{J, C\}. \tag{7.5}$$

where the joint pf $P(\psi_t)$ is computed recursively, starting with known $P(\psi_0)$, as follows:

$$P(\psi_t) = \sum_{O_{t-n-1} \in \mathbf{O}} P(\psi_t, O_{t-n-1}) = \sum_{O_{t-n-1} \in \mathbf{O}} P(O_{t-1}, \psi_{t-1}) = \sum_{O_{t-n-1} \in \mathbf{O}} P_R(O_{t-1}|\psi_{t-1}) P(\psi_{t-1}),$$

---

[1] In particular, Kullback–Leibler divergence was computed for two pairs of the "real" and "learnt" pfs (one learnt with mixture ratio model and the other with standard mixture model), the details follow.

[2] The time index $t \in \mathbf{t}$ in $KL^t(P_R \| P_N)$ denotes Kullback-Leibler divergence after the $t$-th observation.

with

$$\psi_t = O_{t-1}, O_{t-2}, ..., O_{t-n}$$
$$\psi_{t-1} = O_{t-2}, O_{t-3}, ..., O_{t-n-1}.$$

As stated above, the quality of both models is compared by comparing the values of $KL^t(P_R\|P_J)$ and $KL^t(P_R\|P_C)$ computed via (7.5), where the smaller value of the Kullback-Leibler divergence indicates better predicting and modelling quality of the respective model.

The learnt predictors $P_J(O_t|\psi_t)$ and $P_C(O_t|\psi_t)$ depend on the particular data realisations. Thus, to make a reliable comparison of the mentioned variables, it is necessary to make a Monte Carlo (MC) study of them.

### 7.2.1 Compared Models

In simulation, the system with 5 possible observations $O_t \in \mathbf{O} = \{1, 2, ..., |\mathbf{O}|\}$, where $|\mathbf{O}| = 5$ was modelled. Both mixture ratio model (denoted as $M_J(O_t|\psi_t, \Theta)$) and standard mixture model (denoted as $M_C(O_t|\psi_t, \bar{\Theta})$) were 2-component and used the same regression vector $\psi_t$ and respective components operated on the same subsets of data-vector $\boldsymbol{\phi}$ ($\boldsymbol{\phi}_1$ and $\boldsymbol{\phi}_2$, see Definition 5.1), which are specified as follows:

- $\Theta = \alpha \times \omega_1 \times \omega_2$,

- $\psi_t = (O_{t-1}, O_{t-2}) \in \boldsymbol{\psi} = \mathbf{O} \times \mathbf{O}$,

- $\phi_t = (O_t, \psi_t) \in \boldsymbol{\phi} = \mathbf{O} \times \mathbf{O} \times \mathbf{O}$,

- $\boldsymbol{\phi}_1 = \boldsymbol{\phi}_2 = \mathbf{O} \times \mathbf{O}$,

- $\phi_{t;1} = (O_t, O_{t-1}) \qquad \phi_{t;2} = (O_t, O_{t-2})$.

**Mixture ratio model**   The joint pf of the mixture ratio model equals (cf. (5.1)):

$$J(\phi_t|\Theta) = J(O_t, \psi_t|\Theta) = \alpha J(O_t, O_{t-1}|\Theta) + (1-\alpha)J(O_t, O_{t-2}|\Theta) = \alpha \frac{1}{25} \prod_{i=1}^{25} \omega_{1i}^{\Delta_{1i}} + (1-\alpha)\frac{1}{25} \prod_{i=1}^{25} \omega_{2i}^{\Delta_{2i}},$$

where

- $\Theta = \left(\alpha, (\omega_{1i})_{i=1}^{25}, (\omega_{2i})_{i=1}^{25}\right) \in \boldsymbol{\Theta}, \quad \sum_{i=1}^{25} \omega_{1i} = \sum_{i=1}^{25} \omega_{2i} = 1, \quad \alpha, \omega_{1i}, \omega_{2i} \in [0,1], \quad i \in \{1, 2, ..., 25\}$,

- $\Delta_{1i} = \Delta_{1i}(O_t, O_{t-1})$ and $\Delta_{2i} = \Delta_{2i}(O_t, O_{t-2})$ are indicators for all combinations ($O_t$, $O_{t-1}$) and ($O_t$, $O_{t-2}$), respectively.

The corresponding parametric system model then reads (cf. (3.4), (5.3)):

$$M_J(O_t|\psi_t, \Theta) = \frac{\alpha \prod_{i=1}^{25} \omega_{1i}^{\Delta_{1i}} + (1-\alpha)\prod_{i=1}^{25} \omega_{2i}^{\Delta_{2i}}}{\alpha \underbrace{\sum_{O_t \in \mathbf{O}} \prod_{i=1}^{25} \omega_{1i}^{\Delta_{1i}}}_{W_1(\psi_t, \omega_1)} + (1-\alpha)\underbrace{\sum_{O_t \in \mathbf{O}} \prod_{i=1}^{25} \omega_{2i}^{\Delta_{2i}}}_{W_2(\psi_t, \omega_2)}}$$

$$= w_1(\psi_t, \Theta)\frac{\prod_{i=1}^{25} \omega_{1i}^{\Delta_{1i}}}{W_1(\psi_t, \omega_1)} + w_2(\psi_t, \Theta)\frac{\prod_{i=1}^{25} \omega_{2i}^{\Delta_{2i}}}{W_2(\psi_t, \omega_2)}, \qquad (7.6)$$

where

$$w_1(\psi_t, \Theta) = \frac{\alpha\, W_1(\psi_t, \omega_1)}{\alpha W_1(\psi_t, \omega_1) + (1-\alpha)W_2(\psi_t, \omega_2)} \qquad w_2(\psi_t, \Theta) = \frac{(1-\alpha)\, W_2(\psi_t, \omega_2)}{\alpha W_1(\psi_t, \omega_1) + (1-\alpha)W_2(\psi_t, \omega_2)}.$$
(7.7)

**Standard mixture model**   The standard mixture model had the following form:

$$M_C(O_t|\psi_t, \bar{\Theta}) = \bar{\alpha} \prod_{i=1}^{25} \bar{\omega}_{1i}^{\Delta_{1i}} + (1-\bar{\alpha}) \prod_{i=1}^{25} \bar{\omega}_{2i}^{\Delta_{2i}},$$
(7.8)

where

- $\bar{\Theta} = \left(\bar{\alpha}, (\bar{\omega}_{1i})_{i=1}^{25}, (\bar{\omega}_{2i})_{i=1}^{25}\right) \in \bar{\Theta}, \qquad \bar{\alpha}, \bar{\omega}_{1i}, \bar{\omega}_{2i} \in [0,1], \quad i \in \{1, 2, ..., 25\},$

- $\sum_{j=1}^{5} \bar{\omega}_{1(5i+j)} = \sum_{j=1}^{5} \bar{\omega}_{2(5i+j)} = 1 \quad i \in \{0, 1, 2, 3, 4\},$

- $\Delta_{1i}$ and $\Delta_{2i}$ are indicators for all combinations $(O_t, O_{t-1})$ and $(O_t, O_{t-2})$, respectively.

### 7.2.2   Results of Comparison

**Compared quality of models**   The quality of models was compared by a MC study. Three sets of 200 stochastic simulations (with identical initial conditions) of the dynamic system were made. In each simulation, both, the mixture ratio and the standard mixture models were learnt recursively on 500 observations. Those were generated by the real model $P_{R_i}(O_t|\psi_t)$, $i \in \{1, 2, 3\}$, which was different for each set of simulations:

1. The real model was the mixture ratio model (7.6) with parameters $\Theta_R \in \Theta$ generated randomly for each simulation, i.e.
$$P_{R_1}(O_t|\psi_t) = M_J(O_t|\psi_t, \Theta_R).$$
(7.9)

2. The real model was the standard mixture model (7.8) with parameters $\Theta_S \in \bar{\Theta}$ generated randomly for each simulation, i.e.
$$P_{R_2}(O_t|\psi_t) = M_C(O_t|\psi_t, \Theta_S).$$
(7.10)

3. The real model was the mixture ratio model (7.6) with parameters $\Theta_F \in \Theta$ fixed for all simulations giving true dynamic weights $w_c(\psi_t)$, i.e.
$$P_{R_3}(O_t|\psi_t) = M_J(O_t|\psi_t, \Theta_F).$$
(7.11)

The initial statistics $V_0 = (v_0, (V_{0;c})_{c \in \mathbf{c}})$ describing prior information about the parameters were set to
$$v_0 = 1 \qquad V_{0;c} = 1 \ \ (\text{entrywise}) \qquad c \in \mathbf{c}$$

for both models, which express no prior information about the parameters. The initial regression vector $\psi_0$ was set to
$$\psi_0 = (O_0, O_{-1}, O_{-2}) = (|\mathbf{O}|, |\mathbf{O}|, |\mathbf{O}|).$$

The values of Kullback-Leibler divergence for both models (at the end of the simulation), $KL^{500}(P_R||P_J)$ and $KL^{500}(P_{R_i}||P_C)$, were studied (see (7.5)). They were compared by their differences

$$\Delta KL_i = KL^{500}(P_{R_i}\|P_J) - KL^{500}(P_{R_i}\|P_C), \qquad i \in \{1,2,3\}, \tag{7.12}$$

where negative values of $\Delta KL_i$ indicates better predicting performance of mixture ratio model. $P_{R_i}(O_1,...,O_t|\psi_0)$, $P_J(O_1,...,O_t|\psi_0)$ and $P_C(O_1,...,O_t|\psi_0)$, $i \in \{1,2,3\}$, denote the joint pfs for the real models $P_{R_1}, P_{R_2}$ and $P_{R_3}$, the mixture ratio model and the standard mixture model, respectively (cf. (7.3)).

The results are summarized in Table 7.1. For the illustration, Figure 7.2 displays time-evolution of Kullback-Leibler divergence for both models, $KL^t(P_{R_i}\|P_J)$ and $KL^t(P_{R_i}\|P_C)$, $i \in \{1,2,3\}$ (7.5), in one of the simulations.

| Simulated Case | $P_{R_1}$ | $P_{R_2}$ | $P_{R_3}$ |
|---|---|---|---|
| Mean | -0.7001 | -0.6029 | -9.1191 |
| Median | -0.7818 | -0.5957 | -9.0598 |
| Minimum | -3.2138 | -4.0519 | -16.1047 |
| Maximum | 4.9534 | 1.9341 | -4.5774 |
| Standard Deviation | 1.0583 | 0.9377 | 2.0317 |

Table 7.1: Sample statistics of $\Delta KL_i$ (7.12), $i \in \{1,2,3\}$. The columns belong to the simulations carried with different real models $P_{R_1}, P_{R_2}, P_{R_3}$ (7.9-7.11).
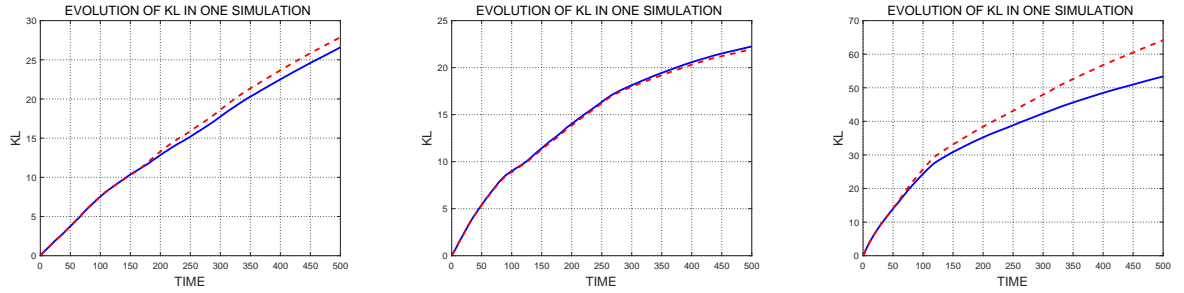


Figure 7.2: Time evolutions of $KL^t(P_R\|P_J)$ (blue line —) and $KL^t(P_R\|P_C)$ (red line - -) (7.5) in one of the simulations. Left to right: $P_{R_1}$ with randomly chosen $\Theta_R$ (7.9); $P_{R_2}$ with randomly chosen $\Theta_S$ (7.10); $P_{R_3}$ with fixed $\Theta_F$ (7.11) causing truly dynamic weights $w_c(\psi_t)$ (7.7).
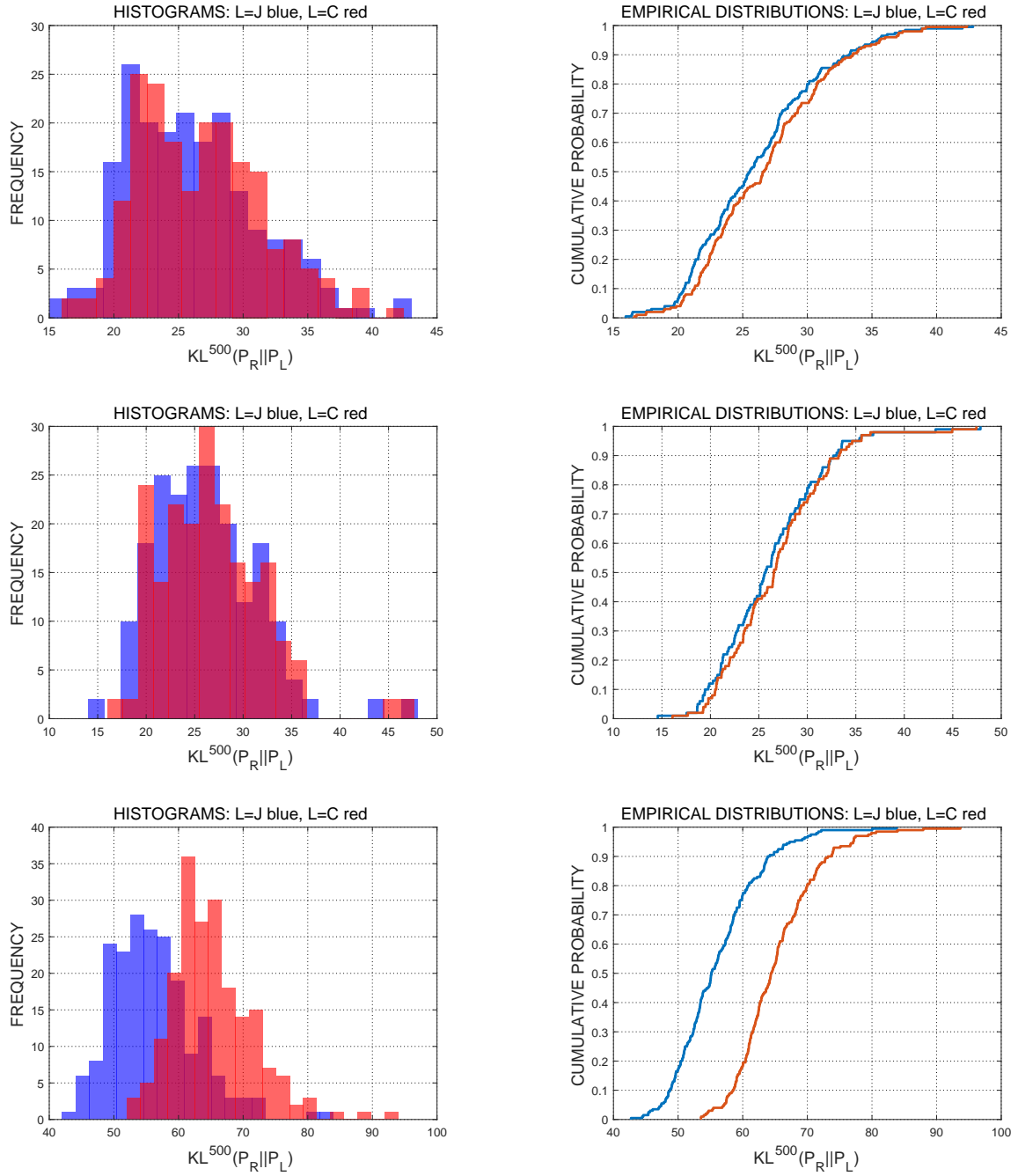
Figure 7.3: Histograms of $KL^{500}(P_{R_i} \| P_J)$ (blue) and $KL^{500}(P_{R_i} \| P_C)$ (7.5) (left) and their empirical distribution function (right). From the top row to the bottom one: $P_{R_1}$ with randomly chosen $\Theta_R$ (7.9); $P_{R_2}$ with randomly chosen $\Theta_S$ (7.10); $P_{R_3}$ with fixed $\Theta_F$ (7.11) causing truly dynamic weights $w_c(\psi_t)$ (7.7)..

**Discussion** The carried MC study shows two important results, which are implied by Table 7.1 and Figure 7.1. The first is that the mixture ratio model serves well as a parametric system model even if the system is simulated by a standard mixture model. It is even slightly better, in terms of the studied Kullback-Leibler divergence 7.5, then the standard mixture model as it (probably) copes better with errors caused by approximate learning (Algorithm 2).

The second result is that the expectation that the mixture ratio outperforms the standard mixture if the mixture ratio is simulated was confirmed. The numerical approximations used in learning of this model

47

allowed feasible and approximate sufficiently-precise learning of it. The simulation results suggest that it can be used for learning and predicting real-life stochastic systems.

## 7.3 Trading with Futures

The mixture ratio and the standard mixture models were also compared on a real data containing daily closing prices of various futures contracts traded in public market for 20 years. The prices was discretized into 15 equally spaced segments between their minimum and maximum. Both models were used to create an automatic adaptive agent who optimizes her trading strategy.

The key ingredient of a successful trading strategy is the price prediction. Traders with futures enters a so-called long position when they expect the futures price to rise and to short position when they expect it to fall. Their gain is then the differences in prices between the price when they entered a position and the price when they closed their position. Importing aspect for modelling is also that the agent do not influence the future price by her actions. More details about trading with futures contracts can be found e.g. in [34].

### 7.3.1 General strategy design

Let $S_t \in \mathbf{S}$ denote a state of the system after the $t$-th decision epoch. To quantify the agent's aim, a real-valued reward function $R = R(S_t, S_{t-1}, A_t)$ is introduced. In general, it depends on the actual state $S_{t-1}$, the agent's action $A_t$ and on the next state $S_t$.

The aim of the agent is to optimize her expected reward $E[R(S_t, S_{t-1}, A_t)]$ accumulated through $h \in \mathbb{N}$ days, it equals as follows:

$$E[R(S_\tau, S_{\tau-1}, A_\tau)] = \sum_{\substack{S_\tau \in \mathbf{S} \\ S_{\tau-1} \in \mathbf{S}}} \sum_{A_\tau \in \mathbf{A}} R(S_\tau, S_{\tau-1}, A_\tau) P(S_\tau | S_{\tau-1}, A_\tau) \pi(A_\tau | S_{\tau-1}) P(S_{\tau-1}), \tag{7.13}$$

where $t$ is the initial decision epoch, $P(S_\tau | S_{\tau-1}, A_\tau)$ is the transition pf, $\pi(A_\tau | S_{\tau-1}) \in \boldsymbol{\pi}$ is the decision rule and $P(S_{\tau-1})$ is the pf of the state $S_{\tau-1}$.

The agent does so by choosing the set of optimal decision rules $(\pi^*(A_\tau | S_{\tau-1}))_{\tau=t}^{t+h}$ (optimal policy) maximizing the accumulated reward, i.e.

$$(\pi^*(A_\tau | S_{\tau-1}))_{\tau=t}^{t+h} = \operatorname*{argmax}_{(\pi(A_\tau | S_{\tau-1}))_{\tau=t}^{t+h}} \sum_{\tau=t}^{t+h} E[R(S_\tau, S_{\tau-1}, A_\tau)]. \tag{7.14}$$

The optimal policy $(\pi^*(A_\tau | S_{\tau-1}))_{\tau=t}^{t+h}$ is found by dynamic programming, [8].

**Algorithm 6** (Dynamic programming). Let $S_t \in \mathbf{S}$ denotes the state of the system. The solution $(A_\tau^*(S_{t-1}))_{\tau=t}^{t+h}$ to the following optimization problem

$$(\pi^*(A_\tau | S_{\tau-1}))_{\tau=t}^{t+h} = \operatorname*{argmax}_{(\pi(A_\tau | S_{\tau-1}))_{\tau=t}^{t+h}} \sum_{\tau=t}^{t+h} E[R(S_\tau, S_{\tau-1}, A_\tau)] \tag{7.15}$$

is given by the following algorithm. Starting with $\varphi_{t+h}(S_{t+h}) = 0$. For $\tau = t + h, t + h - 1, ..., t$

- $A_\tau^*(S_{\tau-1}) = \underset{A_\tau(S_{t-1})}{\operatorname{argmax}} \operatorname{E}[\operatorname{R}(S_\tau, S_{\tau-1}, A_\tau) + \varphi_\tau(S_\tau)|S_{\tau-1}, A_\tau]$

$$= \underset{A_\tau(S_{t-1})}{\operatorname{argmax}} \sum_{S_\tau \in \mathbf{S}} (\operatorname{R}(S_\tau, S_{\tau-1}, A_\tau) + \varphi_\tau(S_\tau)) \operatorname{P}(S_\tau|S_{\tau-1}, A_\tau)$$

- $\pi^*(A_\tau|S_{\tau-1}) = \delta_{A_\tau A_\tau^*(S_{\tau-1})}$
- $\varphi_{\tau-1}(S_{\tau-1}) = \operatorname{E}[\operatorname{R}(S_\tau, S_{\tau-1}, A_\tau) + \varphi_\tau(S_\tau)|S_{\tau-1}, A_\tau^*(S_\tau)]$

The predictive pf $\operatorname{P}(S_\tau|S_{\tau-1}, A_\tau)$ are assumed to be constant during optimization. Generally the pf $\operatorname{P}(S_\tau|S_{\tau-1}, A_\tau)$ are learned from data on-line, thus it is not constant. This approximation is included to assure feasibility of the optimisation. More details and proof of optimality can be found in [8]. □

### 7.3.2 Trading strategy

For the trading problem, the reward function, respecting the trading described in the beginning of this section, is specified as follows.

**Definition 14** (Reward function). The agent's reward function is defined as follows

$$\operatorname{R}(O_{t-1}, C_{t-1}, A_t) = -\delta_{A_t 1} 2^{\delta_{C_{t-1} 2}} O_{t-1} + \delta_{A_t 2} 2^{\delta_{C_{t-1} 1}} O_{t-1} + \delta_{A_t 3} (-1)^{\delta_{C_{t-1} 2}} O_{t-1}, \tag{7.16}$$

where

- times indices $t$ denote days,

- $\delta$ denotes Kronecker delta function (Definition 7),

- $O_t \in \mathbf{O}$ denotes closing price after $t$-th day,

- $A_t \in \mathbf{A} = \{1, 2, 3, 4\}$ denotes action at the beginning of the $t$-th day,

  - $A_t = 1$ means entering long position,
  - $A_t = 2$ denotes entering short position,
  - $A_t = 3$ stands for closing actual position,
  - $A_t = 4$ means doing nothing,

- $C_t \in \mathbf{C} = \{1, 2, 3\}$ denotes opened position in the $t$-th day,

  - $C_t = 1$ means long position,
  - $C_t = 2$ denotes short position,
  - $C_t = 3$ stands for no open position.

□

As mentioned in the beginning of this section, the key element in the successful trading strategy is the prediction of the future price. Here it will be expressed by a transition pf $\operatorname{P}(O_t|\psi_t)$, where $O_t \in \mathbf{O}$

is the future price and $\psi_t$ is the regression vector[3], cf. (3.1). It will contain the actual price $O_{t-1}$, the previous price $O_{t-2}$ and the average price for the past 7 days denotes as $\overline{O}_{t-1}$, i.e.

$$\psi_t = \left( O_{t-1}, O_{t-2}, \overline{O}_{t-1} \right). \tag{7.17}$$

The transition pf $P(O_t|\psi_t)$ will be learned from the data by using the mixture ratio model (Chapter 5) as well as the standard mixture model, [18]. It is discussed below.

For the specific reward function (7.16) quantifying trading goals, Algorithm 6 is used with the following elements:

- the state $S_t$ from the algorithm corresponds to $(\psi_{t+1}, C_t)$,

- the transition pf $P(S_t|S_{t-1}, A_t) = P(\psi_{t+1}, C_t|\psi_t, C_{t-1}, A_t)$, respecting the structure and dependencies of the trading problem, has the following form[4]
  $P(\psi_{t+1}, C_t|\psi_t, C_{t-1}, A_t) = P(O_t|\psi_t)d_1(C_t|A_t, C_{t-1})d_2(\overline{O}_t|O_t, \overline{O}_{t-1})\delta_{O_{t-1}\psi_{t;1}}$ with

  o $d_1(S_t|A_t, S_{t-1}) = \delta_{S_t S_{new}}$

$$S_{new} = \begin{cases} 1, & A_t = 1 \\ 2, & A_t = 2 \\ 3, & A_t = 3 \\ S_{t-1} & A_t = 4 \end{cases}$$

  o $d_2(\overline{O}_t|\overline{O}_{t-1}) = \delta_{\overline{O}_t \overline{O}_{new}}$ $\qquad \overline{O}_{new} = \overline{O}_{t-1} + \frac{O_t - O_{t-7}}{7}$

Hence, it remains to specify the transition pf $P(O_t|\psi_t)$. It is learned on-line by both agents using the mixture ratio and the standard mixture models.

**Mixture ratio model** The joint probability in mixture ratio model $J(O_t, \psi_t|\Theta)$ (5.1) has the following form:

$$J(O_t, \psi_t|\Theta) = \alpha_1 \frac{1}{225} \prod_{i=1}^{225} \omega_{1i}^{\Delta_{1i}} + \alpha_2 \frac{1}{225} \prod_{i=1}^{225} \omega_{2i}^{\Delta_{2i}} + \alpha_3 \frac{1}{225} \prod_{i=1}^{225} \omega_{3i}^{\Delta_{3i}} \tag{7.18}$$

where

- $\Theta = \left( (\alpha_i)_{i=1}^3, (\omega_{1i})_{i=1}^{225}, (\omega_{2i})_{i=1}^{225}, (\omega_{3i})_{i=1}^{225} \right) \in \Theta, \quad \sum_{i=1}^3 \alpha_i = \sum_{i=1}^{225} \omega_{1i} = \sum_{i=1}^{225} \omega_{2i} = 1$,

- $\Delta_{1i} = \Delta_{1i}(O_t, O_{t-1})$, $\Delta_{2i} = \Delta_{2i}(O_t, O_{t-2})$ and $\Delta_{3i} = \Delta_{3i}(O_t, \overline{O}_{t-1})$ are indicators for all combinations $(O_t, O_{t-1}),(O_t, O_{t-2})$ and $(O_t, \overline{O}_{t-1})$, respectively.

The parametric model $M_J(O_t|\psi_t, \Theta)$ is then given as follows:

$$M_J(O_t|\psi_t, \Theta) = \frac{\sum_{c=1}^3 \alpha_c \prod_{i=1}^{225} \omega_{ci}^{\Delta_{ci}}}{\sum_{O_t \in \mathbf{O}} \sum_{c=1}^3 \alpha_c \prod_{i=1}^{225} \omega_{ci}^{\Delta_{ci}}} = H_t(\Theta) \sum_{c=1}^3 \alpha_c \prod_{i=1}^{225} \omega_{ci}^{\Delta_{ci}} \tag{7.19}$$

$$H_t(\Theta) = \frac{1}{\sum_{O_t \in \mathbf{O}} \sum_{c=1}^3 \alpha_c \prod_{i=1}^{225} \omega_{ci}^{\Delta_{ci}}}$$

---

[3]The regression vector will not contain action $A_t$ as in (cf. 3.1). The action does not influence the predicting price, thus omitting $A_t$ from regression vector simplifies the notation.

[4]The functions $d_1(C_t|A_t, C_{t-1},)d_2(\overline{O}_t|O_t, \overline{O}_{t-1}),\delta_{O_{t-1}\psi_{t;1}}$ describes deterministic evolution of $C_t$, $\overline{O}_t$ and $O_{t-1}$, only the evolution of the future price $O_t$ is assumed to be stochastic and is described by the transition pf $P(O_t|\psi_t)$.

The key transition probability $P_J(O_t|\psi_t)$ is then obtained from the parametric model $M_J(O_t|\psi_t, \Theta)$ (7.19) as follows:

$$P_J(O_t|\psi_t) = \int_{\Theta} M_J(O_t|\psi_t, \Theta) P_{Jt}(\Theta) d\Theta, \tag{7.20}$$

where pd $P_{Jt}(\Theta) = P_J(\Theta|v_t, V_t)$ is obtained recursively as described in Section 6. The integration in (7.20) is approximated similarly as the integral $L_{-1}$ (5.5), i.e. the function $H_t(\Theta)$ is approximated by its linear Taylor expansion around the expected values of parameter $\hat{\Theta}_t^c$, $c \in \{1, 2, 3\}$ computed by (5.10). The transition probability $P_J(O_t|\psi_t)$ is then computed as follows:

$$\begin{aligned}
P_J(O_t|\psi_t) &\approx \sum_{c=1}^{3} H_t(\tilde{\Theta}_t) \frac{\int_{\Theta} \alpha_c \prod_{i=1}^{225} \omega_{ci}^{\Delta_{ci}} \prod_{d=1}^{3} \alpha_d^{v_d - 1} \prod_{i=1}^{225} \omega_{di}^{V_{di} - 1} d\Theta}{\mathrm{Be}(v) \prod_{d=1}^{3} \mathrm{Be}(V_c)} \\
&= \sum_{c=1}^{3} H_t(\tilde{\Theta}_t) \frac{\mathrm{Be}(v + \delta_c) \prod_{d=1}^{3} \mathrm{Be}(V_d + \delta_c \Delta_c)}{\mathrm{Be}(v) \prod_{d=1}^{3} \mathrm{Be}(V_d)} \\
&= \sum_{c=1}^{3} H_t(\tilde{\Theta}_t) \frac{\mathrm{Be}(v + \delta_c) \prod_{d=1}^{3} \mathrm{Be}(V_d + \delta_c \Delta_c)}{\mathrm{Be}(v) \prod_{d=1}^{3} \mathrm{Be}(V_d)} = \sum_{c=1}^{3} H_t(\tilde{\Theta}_t) \tilde{\alpha}_{t;c} \prod_{i=1}^{225} \tilde{\omega}_{t;ci}^{\Delta_{ci}} \tag{7.21}
\end{aligned}$$

$$\tilde{\alpha}_{t;c} = \frac{v_{t-1;c}}{\sum_{d \in \mathbf{c}} v_{t-1;d}} \qquad \tilde{\omega}_{t;cj} = \frac{V_{t-1;cj}}{\sum_{i=1}^{m_c} V_{t-1;ci}} \qquad \tilde{\Theta}_t = (\tilde{\alpha}_t, (\tilde{\omega}_{t;c})_{c \in \mathbf{c}}) \tag{7.22}$$

**Standard mixture model** The conditional probability $M_C(O_t|\psi_t, \Theta)$ modelled in standard mixture model has the following form:

$$M_C(O_t|\psi_t, \Theta) = \bar{\alpha}_1 \prod_{i=1}^{225} \bar{\omega}_{1i}^{\Delta_{1i}} + \bar{\alpha}_2 \prod_{i=1}^{225} \bar{\omega}_{2i}^{\Delta_{2i}} + \bar{\alpha}_3 \prod_{i=1}^{225} \bar{\omega}_{3i}^{\Delta_{3i}} \tag{7.23}$$

where

- $\bar{\Theta} = \left( \bar{\alpha}, (\bar{\omega}_{1i})_{i=1}^{225}, (\bar{\omega}_{2i})_{i=1}^{225}, (\bar{\omega}_{3i})_{i=1}^{225} \right) \in \Theta$, $\quad \sum_{j=1}^{15} \bar{\omega}_{1(15i+j)} = \sum_{j=1}^{15} \bar{\omega}_{2(15i+j)} = \sum_{j=1}^{15} \bar{\omega}_{3(15i+j)} = 1$ $i \in \{0, 1, ..., 14\}$,

- $\Delta_{1i} = \Delta_{1i}(O_t, O_{t-1})$, $\Delta_{2i} = \Delta_{2i}(O_t, O_{t-2})$ and $\Delta_{3i} = \Delta_{3i}(O_t, \overline{O}_{t-1})$ are indicators for all combinations $(O_t, O_{t-1})$,$(O_t, O_{t-2})$ and $(O_t, \overline{O}_{t-1})$, respectively.

The objective transition probability $P_C(O_t|\psi_t)$ is obtained as follows:

$$P_C(O_t|\psi_t) = \int_{\Theta} M_C(O_t|\psi_t, \Theta) P_{Ct}(\Theta) d\Theta = \sum_{c=1}^{3} \tilde{\bar{\alpha}}_{t;c} \prod_{i=1}^{225} \tilde{\bar{\omega}}_{t;ci}^{\Delta_{ci}}, \tag{7.24}$$

where pd $P_{Ct}(\Theta) = P_C(\Theta|\bar{v}_t, \bar{V}_t)$ is learned recursively, $\bar{v}_t, \bar{V}_t$ are sufficient statistics determining this pd and

$$\tilde{\bar{\alpha}}_{t;c} = \frac{\bar{v}_{t-1;c}}{\sum_{d \in \mathbf{c}} \bar{v}_{t-1;d}} \qquad \tilde{\bar{\omega}}_{t;cj} = \frac{\bar{V}_{t-1;cj}}{\sum_{i=1}^{m_c} \bar{V}_{t-1;ci}}.$$

Learning of standard mixture model is described in detail in [18].

### 7.3.3 Results of Trading

Both agents traded with 16 different futures contracts. Each of the time series was discretized (see Figure 7.4) into 15 values and divided into two halves. The first half of the data was used to learning only, while in the second half, both agents also traded.

The accumulated reward (7.16) of both agents was compared. The results are summarized in Table 7.2 and Figure 7.5. Trading with three picked futures can be seen on Figure 7.6.
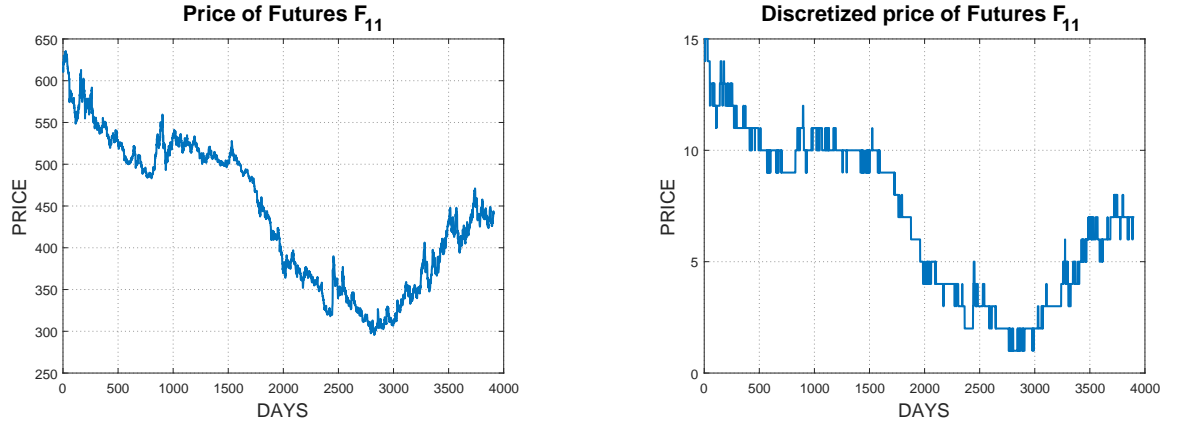


Figure 7.4: Price evolution of futures $F_{11}$ for 3910 trading days; the real prices on the left, the discretized ones on the right.
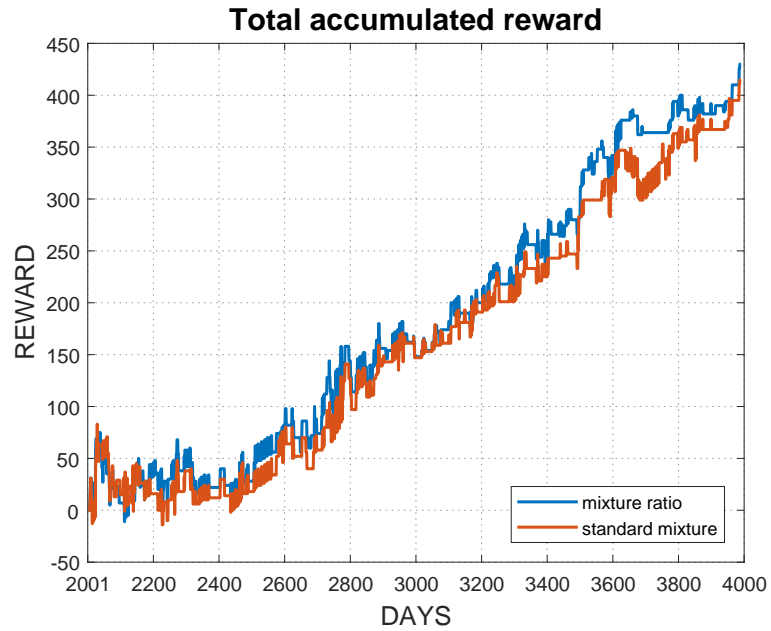


Figure 7.5: Total accumulated reward (7.16) while trading, i.e. for the past 1989 days.

| futures | mixture ratio reward | standard mixture reward |
|---|---|---|
| $F_1$ | 32 | 18 |
| $F_2$ | -1 | -3 |
| $F_3$ | 0 | 0 |
| $F_4$ | 20 | 10 |
| $F_5$ | -6 | -2 |
| $F_6$ | 17 | 25 |
| $F_7$ | 37 | 55 |
| $F_8$ | 70 | 74 |
| $F_9$ | 42 | 38 |
| $F_{10}$ | 14 | 10 |
| $F_{11}$ | 22 | 22 |
| $F_{12}$ | 13 | 15 |
| $F_{13}$ | 42 | 37 |
| $F_{14}$ | 60 | 56 |
| $F_{15}$ | -6 | -6 |
| $F_{16}$ | 74 | 66 |
| **TOTAL** | 430 | 415 |
| **MEAN** | 26.875 | 25.938 |
| **MEDIAN** | 21 | 20 |
| **MIN** | -6 | -6 |
| **MAX** | 74 | 74 |

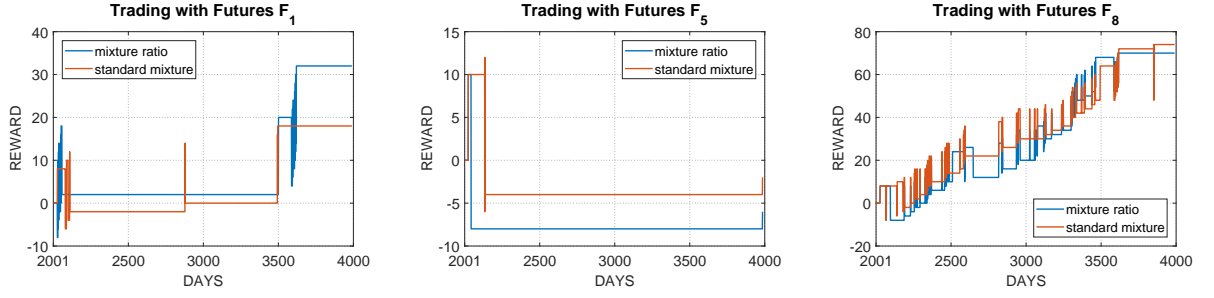Table 7.2: The results from trading with all of the traded futures in terms of accumulated reward (7.16).

Figure 7.6: Accumulated reward while trading with three futures - from the left $F_1$, $F_5$ and $F_8$.

**Discussion** As Figure 7.5 and Table 7.2 show, the performance of agent with mixture ratio model is similar to the agent with standard mixture model. Predicting time series of futures prices is in general complicated task and it is very hard to beat the useless prediction that the price will stay the same the next day. Both agents ended in positive numbers in terms of their reward, however this does not mean that the agents would earn proportional amount of real money if they traded on a futures market. Our reward function (7.16) does not include transaction costs, which complicate real algorithmic trading even more.

Despite this, the results on the futures data leave an important message on mixture ratio model. As in simulation in Section 7.2, the mixture ratio model is at least as good as the standard mixture model even when real data is predicted and it is suitable for real decision making.

# Conclusions

In this work, a very flexible, but yet unconsidered model set of mixture ratios with components from exponential family was build. Commonly used dynamic mixture models mostly uses fixed weights of components. This work overcame this limitation and provided universally approximating parametric model for a general dynamic system. Also, the approximate Bayesian learning, presented in [3] and summarized in Chapter 2, was specified for the mixture ratio models.

Two important types of components from EF were closely studied - Markov chains components (Chapter 5) and Gaussian components (Chapter 6). The first type serves as a model for discrete variables dependent on discrete regression vector while the latter one models continuous variables dependent on continuous regression vector. For mixtures with these components, the general learning algorithm described in Chapter 4 was converted into particular formulas and summarized in Algorithm 3 and Algorithm 5. The formulas include necessary numerical approximations of integrals occurring in Bayes' theorem (Theorem 1) and in Kullback-Leibler divergence (Definition 8) used for projection of pds.

The way how the mixture ratio model was constructed, i.e. by modelling joint pd on output and regression vector, allowed to model also mixed discrete-continuous systems. The practical modelling of such systems is complicated mainly if the dependence of discrete variable on continuous regression vector is present. The mixture ratio model can cope easily with such cases, because the joint pd can always be factorized and modelled in the way that overcomes this as it shown in Remark 11.

Chapter 7 provided simulation as well as real-life examples demonstrating simulation and learning strength of mixture ratio model. It was compared with standard mixtures, introduced in [18], in three setups - 1) basic simulation of mixture of two Gaussian components was made, 2) Monte Carlo simulation study of learning and predicting quality of both mixtures with Markov chain components, 3) both models was compared in real-life scenario involving trading with futures on real data. The main outcomes are following:

- When both models are simulated, mixture ratio model shows true dynamic behaviour compared to the standard mixture, see Figure 7.3.

- Mixture ratio model showed similar predictive quality (even slightly better) to standard mixture model when the system was simulated from standard mixture model.

- When the system was simulated from mixture ratio model with parameters giving true dynamic weights, the standard mixture model showed significantly worse predictive quality compared with the mixture ratio model.

- Mixture ratio model showed its applicability for decision making in real-life scenario - trading with futures. It produced slightly better results in terms of accumulated reward then the standard mixture model.

The conversion of the general approximate learning algorithm from Chapter 2 for these specific components required a lot of detailed technical work present throughout the thesis and in Appendix.

Although, the theoretical and experimental parts of this thesis show that it was worthwhile as the mixture ratio models provide a universal approach for modelling general dynamic systems and designing DM strategies.

Many challenging tasks remains to be studied, in particular: i) simulation and real-life testing mixture ratio model with Gaussian and mixed components, ii) examining the suitability of the proposed numerical approximations for high-dimensional models, iii) refining the data-dependent choice of forgetting factors, iv) including structure estimation of the mixture ratio model as well as its use for feature extraction similarly as in [33].

# Appendix

## Derivation of L's (5.11 - 5.13)

In this section, the formulas (5.11 - 5.13) are proved. The derivation starts with recalling definitions of treated objects:

$$L_{-1} = \sum_{c \in \mathbf{c}} \chi_{\phi_c} \gamma_{t;c} \prod_{j=1}^{m_c} \left( \frac{1}{K_{cj}} \right)^{\Delta_{cj}} \int_{\Theta} H(\psi_t, \Theta) Q_{t;c}(\Theta) d\Theta \tag{7.25}$$

$$L_{0c} = \sum_{c \in \mathbf{c}} \chi_{\phi_c} \gamma_{t;c} \prod_{j=1}^{m_c} \left( \frac{1}{K_{cj}} \right)^{\Delta_{cj}} \int_{\Theta} \ln(\alpha_c) H(\psi_t, \Theta) Q_{t;c}(\Theta) d\Theta \tag{7.26}$$

$$L_{cj} = \sum_{c \in \mathbf{c}} \chi_{\phi_c} \gamma_{t;c} \prod_{j=1}^{m_c} \left( \frac{1}{K_{cj}} \right)^{\Delta_{cj}} \int_{\Theta} \ln(\omega_{cj}) H(\psi_t, \Theta) Q_{t;c}(\Theta) d\Theta \quad \forall j \in \{1, 2..., m_c\} \tag{7.27}$$

$$H_t(\Theta) \approx H_t(\tilde{\Theta}_t^d) + \sum_{c \in \mathbf{c}} \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \alpha_c} (\alpha_c - \tilde{\alpha}_{t;c}^d) + \sum_{c \in \mathbf{c}} \sum_{j=1}^{m_c} \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \omega_{cj}} (\omega_{cj} - \tilde{\omega}_{t;cj}^d), \tag{7.28}$$

$$\text{where} \quad \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \alpha_c} = -H_t^2(\tilde{\Theta}_t^d) \sum_{O_t \in \mathbf{O}} \chi_{\phi_c} \prod_{j=1}^{m_c} \left( \frac{\tilde{\omega}_{cj}^d}{K_{cj}} \right)^{\Delta_{cj}}$$

$$\frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \omega_{cj}} = -H_t^2(\tilde{\Theta}_t^d) \tilde{\alpha}_{t;c}^d \sum_{O_t \in \mathbf{O}} \chi_{\phi_c} \left( \frac{1}{K_{cj}} \right)^{\Delta_{cj}} 0^{1-\Delta_{cj}}$$

To prove the formulas (5.11 - 5.13), the following integrals will be used (see Definition 6 and Definition 5):

$$\int_{\mathbf{x}} \prod_{i=1}^{n} x_i^{p_i - 1} dx = \text{Be}(p), \tag{7.29}$$

$$\int_{\mathbf{x}} \ln(x_j) \prod_{i=1}^{n} x_i^{p_i - 1} dx = \frac{\partial}{\partial p_j} \int_{\mathbf{x}} \prod_{i=1}^{n} x_i^{p_i - 1} dx = \frac{\partial}{\partial p_j} \text{Be}(p) = \frac{\partial}{\partial p_j} \frac{\prod_{i=1}^{n} \Gamma(p_i)}{\Gamma(\sum_{i=1}^{m} p_i)}$$

$$= \frac{\Gamma'(p_i) \prod_{\substack{i=1 \\ i \neq j}}^{n} \Gamma(p_i)}{\Gamma(\sum_{i=1}^{m} p_i)} - \frac{\prod_{i=1}^{n} \Gamma(p_i) \Gamma'(\sum_{i=1}^{m} p_i)}{\Gamma^2(\sum_{i=1}^{m} p_i)} = \text{Be}(p) \left( \psi(p_j) - \psi\left( \sum_{i=1}^{n} p_i \right) \right), \tag{7.30}$$

where

$$p \in (\mathbb{R}^+)^n, \quad x \in \boldsymbol{x} = \left\{ x \in [0,1]^n, \sum_{i=1}^{n} x_i = 1 \right\}$$

By inserting (7.28) into (7.26) it is obtained:

$$
\begin{aligned}
L_{-1}(\phi_t, \mathbf{V}_{t-1}) &= \sum_{d \in \mathbf{c}} \chi_{\phi_d} \gamma_{t;d} \prod_{j=1}^{m_d} \left( \frac{1}{K_{dj}} \right)^{\Delta_{dj}} \int_{\Theta} \Big[ \Big( \mathrm{H}_t(\tilde{\Theta}_t^d) + \sum_{e \in \mathbf{c}} \frac{\partial \mathrm{H}_t(\tilde{\Theta}_t^d)}{\partial \alpha_e} (\alpha_e - \tilde{\alpha}_{t;e}^d) \\
&\qquad + \sum_{e \in \mathbf{c}} \sum_{i=1}^{m_e} \frac{\partial \mathrm{H}_t(\tilde{\Theta}_t^d)}{\partial \omega_{ei}} (\omega_{ei} - \tilde{\omega}_{t;ei}^d) \Big) \mathrm{Q}_{t;d}(\Theta) d\Theta \Big] \\
&= \sum_{d \in \mathbf{c}} \chi_{\phi_d} \gamma_{t;d} \prod_{j=1}^{m_d} \left( \frac{1}{K_{dj}} \right)^{\Delta_{dj}} \mathrm{H}_t(\tilde{\Theta}_t^d),
\end{aligned}
$$

which proves (5.11).

Subsequently, insertion of (7.28) into (7.26) gives

$$
\begin{aligned}
L_{0c} &= \sum_{d \in \mathbf{c}} \chi_{\phi_d} \gamma_{t;d} \prod_{j=1}^{m_d} \left( \frac{1}{K_{dj}} \right)^{\Delta_{dj}} \int_{\Theta} \Big[ \ln(\alpha_c) \Big( \mathrm{H}_t(\tilde{\Theta}_t^d) + \sum_{e \in \mathbf{c}} \frac{\partial \mathrm{H}_t(\tilde{\Theta}_t^d)}{\partial \alpha_e} (\alpha_e - \tilde{\alpha}_{t;e}^d) \\
&\qquad + \sum_{e \in \mathbf{c}} \sum_{i=1}^{m_e} \frac{\partial \mathrm{H}_t(\tilde{\Theta}_t^d)}{\partial \omega_{ei}} (\omega_{ei} - \tilde{\omega}_{t;ei}^d) \Big) \mathrm{Q}_{t;d}(\Theta) d\Theta \Big] \\
&= \sum_{d \in \mathbf{c}} \chi_{\phi_d} \prod_{j=1}^{m_d} \left( \frac{1}{K_{dj}} \right)^{\Delta_{dj}} \Big( \mathrm{H}_t(\tilde{\Theta}_t^d) \prod_{f \in \mathbf{c}} \mathrm{Be}(\mathrm{V}_{t-1;f} + \delta_{df} \Delta_d) \int_{\alpha} \ln(\alpha_c) \prod_{f \in \mathbf{c}} \alpha_f^{\mathrm{v}_{t-1;f} - 1 + \delta_{fd}} d\alpha \\
&\qquad + \underbrace{\sum_{e \in \mathbf{c}} \frac{\partial \mathrm{H}_t(\tilde{\Theta}_t^d)}{\partial \alpha_e} \prod_{f \in \mathbf{c}} \mathrm{Be}(\mathrm{V}_{t-1;f} + \delta_{df} \Delta_d) \int_{\boldsymbol{\alpha}} \ln(\alpha_c)(\alpha_e - \tilde{\alpha}_{t;e}^d) \prod_{f \in \mathbf{c}} \alpha_f^{\mathrm{v}_{t-1;f} - 1 + \delta_{fd}} d\alpha}_{\mathrm{I}} \Big) \\
&= \sum_{d \in \mathbf{c}} \chi_{\phi_d} \prod_{j=1}^{m_d} \left( \frac{1}{K_{dj}} \right)^{\Delta_{dj}} \Big( \gamma_{t;d} \mathrm{H}_t(\tilde{\Theta}_t^d) (\psi(\mathrm{v}_c + \delta_{cd}) - \psi(\Sigma_{e \in \mathbf{c}} \mathrm{v}_e + 1)) + \mathrm{I} \Big). \quad (7.31)
\end{aligned}
$$

The remaining part of $L_{0c}$, I, is computed as follows:

$$I = \sum_{e \in \mathbf{c}} \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \alpha_e} \prod_{f \in \mathbf{c}} \mathrm{Be}(V_{t-1;f} + \delta_{df}\Delta_d) \Big[ \mathrm{Be}(v + \delta_d + \delta_e)\Big(\psi(v_c + \delta_{cd} + \delta_{ce}) - \psi(\Sigma_{f \in \mathbf{c}} v_f + 2)\Big)$$

$$- \Big(\tilde{\alpha}_{t;e}^d \mathrm{Be}(v + \delta_d)\Big(\psi(v_c + \delta_{cd}) - \psi(\Sigma_{f \in \mathbf{c}} v_f + 1)\Big)\Big)\Big]$$

$$= \sum_{e \in \mathbf{c}} \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \alpha_e} \tilde{\alpha}_{t;e}^d \gamma_{t;d}[\frac{\delta_{ce}}{v_c + \delta_{cd}} - \frac{1}{\sum_{f \in \mathbf{c}} v_f + 1}]$$

$$= \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \alpha_c} \tilde{\alpha}_{t;c}^d \gamma_{t;d}[\frac{1}{v_c + \delta_{cd}} - \frac{1}{\sum_{f \in \mathbf{c}} v_f + 1}] - \sum_{\substack{e \in \mathbf{c} \\ e \neq c}} \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \alpha_e} \tilde{\alpha}_{t;e}^d \gamma_{t;d}[\frac{1}{\sum_{f \in \mathbf{c}} v_f + 1}]$$

$$= \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \alpha_c} \gamma_{t;d} \frac{v_c + \delta_{cd}}{\sum_{f \in \mathbf{c}} v_f + 1} \frac{\sum_{f \in \mathbf{c}} v_f + 1 - (v_c + \delta_{cd})}{(v_c + \delta_{cd})(\sum_{f \in \mathbf{c}} v_f + 1)} - \sum_{\substack{e \in \mathbf{c} \\ e \neq c}} \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \alpha_e} \tilde{\alpha}_{t;e}^d \gamma_{t;d}[\frac{1}{\sum_{f \in \mathbf{c}} v_f + 1}]$$

$$= \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \alpha_c} \gamma_{t;d} \Big(\frac{1}{\sum_{f \in \mathbf{c}} v_f + 1}\Big) \sum_{\substack{e \in \mathbf{c} \\ e \neq c}} \tilde{\alpha}_{t;e}^d - \sum_{\substack{e \in \mathbf{c} \\ e \neq c}} \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \alpha_e} \tilde{\alpha}_{t;e}^d \gamma_{t;d} \Big(\frac{1}{\sum_{f \in \mathbf{c}} v_f + 1}\Big)$$

$$= \frac{\gamma_{t;d}}{\sum_{f \in \mathbf{c}} v_f + 1} \sum_{\substack{e \in \mathbf{c} \\ e \neq c}} \tilde{\alpha}_{t;e}^d \Big(\frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \alpha_c} - \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \alpha_e}\Big)$$

$$= \frac{\gamma_{t;d}}{\sum_{f \in \mathbf{c}} v_f + 1} \sum_{\substack{e \in \mathbf{c} \\ e \neq c}} \tilde{\alpha}_{t;e}^d H_t^2(\tilde{\Theta}_t^d) \Big( \sum_{O_t \in \mathbf{O}} \chi_{\phi_e} \prod_{j=1}^{m_e} \Big(\frac{\tilde{\omega}_{ej}^d}{K_{ej}}\Big)^{\Delta_{ej}} - \sum_{O_t \in \mathbf{O}} \chi_{\phi_c} \prod_{j=1}^{m_c} \Big(\frac{\tilde{\omega}_{cj}^d}{K_{cj}}\Big)^{\Delta_{cj}} \Big). \tag{7.32}$$

Hence, (7.31) and (7.32) prove the formula (5.12).

Finally, by inserting (7.28) into (7.27), it is obtained:

$$L_{cj} = \sum_{d \in \mathbf{c}} \chi_{\phi_d} \gamma_{t;d} \prod_{j=1}^{m_d} \Big(\frac{1}{K_{dj}}\Big)^{\Delta_{dj}} \int_{\Theta} \Big[\ln(\omega_{cj})\Big(H_t(\tilde{\Theta}_t^d) + \sum_{e \in \mathbf{c}} \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \alpha_e}(\alpha_e - \tilde{\alpha}_{t;e}^d)$$

$$+ \sum_{e \in \mathbf{c}} \sum_{i=1}^{m_e} \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \omega_{ei}}(\omega_{ei} - \tilde{\omega}_{t;ei}^d)\Big) Q_{t;d}(\Theta) d\Theta\Big]$$

$$= \sum_{d \in \mathbf{c}} \chi_{\phi_d} \prod_{j=1}^{m_d} \Big(\frac{1}{K_{dj}}\Big)^{\Delta_{dj}} \Big(H_t(\tilde{\Theta}_t^d)\mathrm{Be}(v_{t-1} + \delta_d) \prod_{\substack{f \in \mathbf{c} \\ f \neq c}} \mathrm{Be}(V_{t-1;f} + \delta_{df}\Delta_d) \int_{\omega_c} \ln(\omega_{cj}) \prod_{i=1}^{m_c} \omega_{ci}^{V_{t-1;ci}-1+\delta_{cd}\Delta_{ci}} d\omega$$

$$+ \sum_{i=1}^{m_c} \frac{\partial H_t(\tilde{\Theta}_t^d)}{\partial \omega_{ci}} \mathrm{Be}(v_{t-1} + \delta_d) \prod_{\substack{f \in \mathbf{c} \\ f \neq c}} \mathrm{Be}(V_{t-1;f} + \delta_{df}\Delta_d) \underbrace{\int_{\omega_c} \ln(\omega_{cj})(\omega_{ci} - \tilde{\omega}_{t;ci}^d) \prod_{i=1}^{m_c} \omega_{ci}^{V_{t-1;ci}-1+\delta_{cd}\Delta_{ci}} d\omega}_{\text{II}}\Big)$$

$$= \sum_{d \in \mathbf{c}} \chi_{\phi_d} \prod_{j=1}^{m_d} \Big(\frac{1}{K_{dj}}\Big)^{\Delta_{dj}} \Big(\gamma_{t;d} H_t(\tilde{\Theta}_t^d)\Big(\psi(V_{t-1;cj} + \delta_{cd}\Delta_{cj}) - \psi\Big(\Sigma_{i=1}^{m_c}(V_{t-1;ci} + \delta_{cd}\Delta_{ci})\Big)\Big) + \mathrm{II}\Big), \tag{7.33}$$

where the remaining part II is computed as follows:

$$\text{II} = \sum_{i=1}^{m_c} \frac{\partial \text{H}_t(\tilde{\Theta}_t^d)}{\partial \omega_{ci}} \text{Be}(v_{t-1} + \delta_d) \prod_{\substack{f \in \mathbf{c} \\ f \neq c}} \text{Be}(V_{t-1;f} + \delta_{df}\Delta_d) \Big[ \text{Be}(V_{t-1;c} + \delta_{cd}\Delta_c + \delta_i)$$

$$\times \Big( \psi(V_{cj} + \delta_{cd}\Delta_{cj} + \delta_{ij}) - \psi\big(\Sigma_{i=1}^{m_c}(V_{ci} + \delta_{cd}\Delta_{ci}) + 1\big) \Big)$$

$$- \tilde{\omega}_{t;ci}^d \text{Be}(V_{t-1;c} + \delta_{cd}\Delta_c) \Big( \psi(V_{cj} + \delta_{cd}\Delta_{cj}) - \psi\big(\Sigma_{i=1}^{m_c}(V_{ci} + \delta_{cd}\Delta_{ci})\big) \Big) \Big]$$

$$= \sum_{i=1}^{m_c} \frac{\partial \text{H}_t(\tilde{\Theta}_t^d)}{\partial \omega_{ci}} \gamma_{t;d} \tilde{\omega}_{t;ci}^d \Big[ \psi(V_{cj} + \delta_{cd}\Delta_{cj} + \delta_{ij}) - \psi\big(\Sigma_{i=1}^{m_c}(V_{ci} + \delta_{cd}\Delta_{ci}) + 1\big)$$

$$- \Big( \psi(V_{cj} + \delta_{cd}\Delta_{cj}) - \psi\big(\Sigma_{i=1}^{m_c}(V_{ci} + \delta_{cd}\Delta_{ci})\big) \Big) \Big]$$

$$= \sum_{i=1}^{m_c} \frac{\partial \text{H}_t(\tilde{\Theta}_t^d)}{\partial \alpha_e} \gamma_{t;d} \tilde{\omega}_{t;ci}^d \Big[ \frac{\delta_{ij}}{V_{t-1;cj} + \delta_{cd}\Delta_{cj}} - \frac{1}{\Sigma_{i=1}^{m_c}(V_{ci} + \delta_{cd}\Delta_{ci})} \Big]$$

$$= \frac{\partial \text{H}_t(\tilde{\Theta}_t^d)}{\partial \omega_{cj}} \gamma_{t;d} \frac{V_{t-1;cj} + \delta_{cd}\Delta_{cj}}{\Sigma_{i=1}^{m_c}(V_{ci} + \delta_{cd}\Delta_{ci})} \Big[ \frac{\Sigma_{i=1}^{m_c}(V_{ci} + \delta_{cd}\Delta_{ci}) - (V_{t-1;cj} + \delta_{cd}\Delta_{cj})}{(V_{t-1;cj} + \delta_{cd}\Delta_{cj})\Sigma_{i=1}^{m_c}(V_{ci} + \delta_{cd}\Delta_{ci})} \Big]$$

$$- \sum_{\substack{i=1 \\ i \neq j}}^{m_c} \frac{\partial \text{H}_t(\tilde{\Theta}_t^d)}{\partial \omega_{ci}} \gamma_{t;d} \tilde{\omega}_{t;ci}^d \Big[ \frac{1}{\Sigma_{i=1}^{m_c}(V_{ci} + \delta_{cd}\Delta_{ci})} \Big]$$

$$= \frac{\partial \text{H}_t(\tilde{\Theta}_t^d)}{\partial \omega_{cj}} \gamma_{t;d} \sum_{\substack{i=1 \\ i \neq j}}^{m_c} \tilde{\omega}_{t;ci}^d \Big[ \frac{1}{\Sigma_{i=1}^{m_c}(V_{ci} + \delta_{cd}\Delta_{ci})} \Big] - \sum_{\substack{i=1 \\ i \neq j}}^{m_c} \frac{\partial \text{H}_t(\tilde{\Theta}_t^d)}{\partial \omega_{ci}} \gamma_{t;d} \tilde{\omega}_{t;ci}^d \Big[ \frac{1}{\Sigma_{i=1}^{m_c}(V_{ci} + \delta_{cd}\Delta_{ci})} \Big]$$

$$= \gamma_{t;d} \Big[ \frac{1}{\Sigma_{i=1}^{m_c}(V_{ci} + \delta_{cd}\Delta_{ci})} \Big] \sum_{\substack{i=1 \\ i \neq j}}^{m_c} \Big( \frac{\partial \text{H}_t(\tilde{\Theta}_t^d)}{\partial \omega_{cj}} - \frac{\partial \text{H}_t(\tilde{\Theta}_t^d)}{\partial \omega_{ci}} \Big) \tilde{\omega}_{t;ci}^d$$

$$= \frac{\chi_{\phi_c} \gamma_{t;d} \text{H}_t^2(\tilde{\Theta}_t^d)}{\Sigma_{i=1}^{m_c}(V_{ci} + \delta_{cd}\Delta_{ci})} \sum_{\substack{i=1 \\ i \neq j}}^{m_c} \tilde{\omega}_{t;ci}^d \tilde{\alpha}_{t;c}^d \left( \sum_{O_t \in \mathbf{O}} \Big( \frac{1}{K_{ci}} \Big)^{\Delta_{ci}} 0^{1-\Delta_{ci}} - \sum_{O_t \in \mathbf{O}} \Big( \frac{1}{K_{cj}} \Big)^{\Delta_{cj}} 0^{1-\Delta_{cj}} \right), \tag{7.34}$$

which, along with (7.33), proves (5.13).

# Derivation of (5.14- 5.16)

The formulas (5.14 - 5.16) (with $\tilde{\gamma}_{t;c}$) is derived by applying the properties of beta function (see Definition 6). Schematically, the derivation progresses as follows:

$$\frac{\sum_{c\in\mathbf{c}}C_c\prod_{j=1}^{m_c}\left(\frac{1}{K_{cj}}\right)^{\Delta_{cj}}\gamma_{t;c}}{\sum_{c\in\mathbf{c}}D_c\prod_{j=1}^{m_c}\left(\frac{1}{K_{cj}}\right)^{\Delta_{cj}}\gamma_{t;c}} = \frac{\sum_{c\in\mathbf{c}}C_c\prod_{j=1}^{m_c}\left(\frac{1}{K_{cj}}\right)^{\Delta_{cj}}\mathrm{Be}(v_{t-1}+\delta(c))\prod_{d\in\mathbf{c}}\mathrm{Be}(V_{t-1;d}+\delta_{cd}\Delta_c)}{\sum_{c\in\mathbf{c}}D_c\prod_{j=1}^{m_c}\left(\frac{1}{K_{cj}}\right)^{\Delta_{cj}}\mathrm{Be}(v_{t-1}+\delta(c))\prod_{d\in\mathbf{c}}\mathrm{Be}(V_{t-1;d}+\delta_{cd}\Delta_c)}$$

$$= \frac{\sum_{c\in\mathbf{c}}C_c\tilde{\tilde{\alpha}}_{t;c}\prod_{j=1}^{m_c}\left(\frac{\tilde{\tilde{\omega}}_{t;cj}}{K_{cj}}\right)^{\Delta_{cj}}\mathrm{Be}(v_{t-1})\prod_{d\in\mathbf{c}}\mathrm{Be}(V_{t-1;d})}{\sum_{c\in\mathbf{c}}D_c\tilde{\tilde{\alpha}}_{t;c}\prod_{j=1}^{m_c}\left(\frac{\tilde{\tilde{\omega}}_{t;cj}}{K_{cj}}\right)^{\Delta_{cj}}\mathrm{Be}(v_{t-1})\prod_{d\in\mathbf{c}}\mathrm{Be}(V_{t-1;d})}$$

$$= \frac{\sum_{c\in\mathbf{c}}C_c\tilde{\tilde{\alpha}}_{t;c}\prod_{j=1}^{m_c}\left(\frac{\tilde{\tilde{\omega}}_{t;cj}}{K_{cj}}\right)^{\Delta_{cj}}}{\sum_{c\in\mathbf{c}}D_c\tilde{\tilde{\alpha}}_{t;c}\underbrace{\prod_{j=1}^{m_c}\left(\frac{\tilde{\tilde{\omega}}_{t;cj}}{K_{cj}}\right)^{\Delta_{cj}}}_{\tilde{\gamma}_{t;c}}},$$

$$\text{where}\quad \tilde{\tilde{\alpha}}_{t;c} = \frac{v_{t-1;c}}{\sum_{d\in\mathbf{c}}v_{t-1;d}}\qquad \tilde{\tilde{\omega}}_{t;cj} = \frac{V_{t-1;cj}}{\sum_{i=1}^{m_c}V_{t-1;ci}}.$$

The form of $C_c, D_c, c\in\mathbf{c}$ depends on the particular formula (5.14 - 5.16), but the derivation remains unchanged for all of formulas (5.14 - 5.16).

# Derivation of $L_{ci}$ (6.26)

In this section the derivation of $L_{ci}$ (6.26) will be provided. See Chapter 6 for the meaning and explanation of all of the used symbols here. It is computed from the following general formulas (6.19):

$$L_{ci} = -\frac{1}{2}(v_{ci}+|\psi_{ci}|+2)L_{ci}^1 - 0.5\Lambda_{ci}L_{ci}^2 - 0.5L_{ci}^3 \tag{7.35}$$

$$L_{ci}^1 = \int_{\Theta}\ln(r_{ci})\mathrm{H}_t(\Theta)\sum_{d\in\mathbf{c}}\gamma_{t;d}\mathrm{Q}_{t;d}(\Theta)d\Theta \tag{7.36}$$

$$L_{ci}^2 = \int_{\Theta}\frac{1}{r_{ci}}\mathrm{H}_t(\Theta)\sum_{d\in\mathbf{c}}\gamma_{t;d}\mathrm{Q}_{t;d}(\Theta)d\Theta \tag{7.37}$$

$$L_{ci}^3 = \int_{\Theta}\frac{1}{r_{ci}}\left((\Theta_{ci}-\hat{\Theta}_{ci})^T V_{\psi ci}(\Theta_{ci}-\hat{\Theta}_{ci})\right)\mathrm{H}_t(\Theta)\sum_{d\in\mathbf{c}}\gamma_{t;d}\mathrm{Q}_{t;d}(\Theta)d\Theta, \tag{7.38}$$

where the function $\mathrm{H}_t(\Theta)$ is approximated via linear Taylor expansion around the expected values of parameter $\Theta$ (6.21) - (6.23) as follows

$$\mathrm{H}_t(\Theta) \approx \mathrm{H}_t(\hat{\Theta}_t^d) + \sum_{c\in\mathbf{c}}\frac{\partial\mathrm{H}_t(\hat{\Theta}_t^d)}{\partial\alpha_c}(\alpha_c-\hat{\alpha}_{t;c}^d) + \sum_{c\in\mathbf{c}}\sum_{j=1}^{n_c}\left(\frac{\partial\mathrm{H}_t(\hat{\Theta}_t^d)}{\partial r_{cj}}\left(r_{cj}-\hat{r}_{t;cj}^d\right) + \frac{\partial\mathrm{H}_t(\hat{\Theta}_t^d)}{\partial\Theta_{cj}}\left(\Theta_{cj}-\hat{\Theta}_{t;cj}^d\right)\right).$$

$$\tag{7.39}$$

Following integrals will be useful and the derivation will refer to them:

$$\int_{r_{ci}} r_{ci}^{-0.5(\tilde{v}_{t;ci}^d+2)} \exp\left\{-\frac{\tilde{\Lambda}_{t;ci}^d}{2r_{ci}}\right\} dr_{ci} = \left|x = \frac{\tilde{\Lambda}_{t;ci}^d}{2r_{ci}}, \; dr_{ci} = -\frac{\tilde{\Lambda}_{t;ci}^d}{2x^2}dx\right| = \left(\frac{\tilde{\Lambda}_{t;ci}^d}{2}\right)^{-0.5\tilde{v}_{t;ci}^d} \int_0^\infty x^{0.5(\tilde{v}_{t;ci}^d-2)} e^{-x} dx$$

$$= \left(\frac{\tilde{\Lambda}_{t;ci}^d}{2}\right)^{-0.5\tilde{v}_{t;ci}^d} \Gamma\left(\frac{\tilde{v}_{t;ci}^d}{2}\right) \qquad (7.40)$$

$$\int_{r_{ci}} r_{ci}^{-0.5(\tilde{v}_{t;ci}^d)} \exp\left\{-\frac{\tilde{\Lambda}_{t;ci}^d}{2r_{ci}}\right\} dr_{ci} = \left|x = \frac{\tilde{\Lambda}_{t;ci}^d}{2r_{ci}}, \; dr_{ci} = -\frac{\tilde{\Lambda}_{t;ci}^d}{2x^2}dx\right| = \left(\frac{\tilde{\Lambda}_{t;ci}^d}{2}\right)^{-0.5\left(\tilde{v}_{t;ci}^d-2\right)} \int_0^\infty x^{0.5(\tilde{v}_{t;ci}^d-4)} e^{-x} dx$$

$$= \left(\frac{\tilde{\Lambda}_{t;ci}^d}{2}\right)^{-0.5\left(\tilde{v}_{t;ci}^d-2\right)} \Gamma\left(\frac{\tilde{v}_{t;ci}^d}{2}-1\right) = \left(\frac{\tilde{\Lambda}_{t;ci}^d}{2}\right)^{-0.5\left(\tilde{v}_{t;ci}^d-2\right)} \Gamma\left(\frac{\tilde{v}_{t;ci}^d}{2}\right) \frac{2}{\tilde{v}_{t;ci}^d-2}$$

$$(7.41)$$

$$\int_{r_{ci}} r_{ci}^{-0.5(\tilde{v}_{t;ci}^d+4)} \exp\left\{-\frac{\tilde{\Lambda}_{t;ci}^d}{2r_{ci}}\right\} dr_{ci} = \left|x = \frac{\tilde{\Lambda}_{t;ci}^d}{2r_{ci}}, \; dr_{ci} = -\frac{\tilde{\Lambda}_{t;ci}^d}{2x^2}dx\right| = \left(\frac{\tilde{\Lambda}_{t;ci}^d}{2}\right)^{-0.5\left(\tilde{v}_{t;ci}^d-2\right)} \int_0^\infty x^{0.5(\tilde{v}_{t;ci}^d-4)} e^{-x} dx$$

$$= \left(\frac{\tilde{\Lambda}_{t;ci}^d}{2}\right)^{-0.5\left(\tilde{v}_{t;ci}^d+2\right)} \Gamma\left(\frac{\tilde{v}_{t;ci}^d}{2}+1\right) = \left(\frac{\tilde{\Lambda}_{t;ci}^d}{2}\right)^{-0.5\left(\tilde{v}_{t;ci}^d+2\right)} \Gamma\left(\frac{\tilde{v}_{t;ci}^d}{2}\right) \frac{\tilde{v}_{t;ci}^d}{2}$$

$$(7.42)$$

$L_{ci}^1$ **6.27**

By inserting (7.39) into (7.36) it is obtained:

$$L_{ci}^1 \approx \sum_{d\in\mathbf{c}} \mathrm{Be}(\tilde{v}_t^d) \prod_{\substack{e\in\mathbf{c}\\e\neq c}} \prod_{j=1}^{n_e} \mathrm{N}(\tilde{v}_{t;ej}^d, \tilde{G}_{t;ej}^d)/\mathrm{P}_e(\psi_e^{\mathsf{c}})$$

$$\times \int_{r_{ci}} \int_{\Theta_{ci}} \ln(r_{ci})\left(\mathrm{H}_t(\hat{\Theta}_t^d) + \frac{\partial \mathrm{H}_t(\hat{\Theta}_t^d)}{\partial r_{ci}}\left(r_{ci} - \hat{r}_{t;ci}^d\right)\right) r_{ci}^{-0.5(v_{ci}+|\psi_{ci}|+2)} \mathrm{P}(r_{ci}, \Theta_{ci}|\tilde{v}_{t;ci}^d, \tilde{V}_{t;ci}^d) d\Theta_{ci} dr_{ci} \quad (7.43)$$

The integral in (7.43) is computed in the following steps:

$$\int_{r_{ci}} \int_{\Theta_{ci}} \left[ \ln(r_{ci}) \left( H_t(\hat{\Theta}_t^d) + \frac{\partial H_t(\hat{\Theta}_t^d)}{\partial r_{ci}} (r_{ci} - \hat{r}_{t;ci}^d) \right) r_{ci}^{-0.5(\nu_{ci} + |\psi_{ci}| + 2)} \right.$$

$$\left. \times \exp \left\{ -\frac{1}{2r_{ci}} \left( (\Theta_{ci} - \hat{\hat{\Theta}}_{t;ci}^d)^T \tilde{V}_{t;\psi ci}^d (\Theta_{ci} - \hat{\hat{\Theta}}_{t;ci}^d) + \tilde{\Lambda}_{t;ci}^d \right) \right\} d\Theta_{ci} dr_{ci} \right]$$

$$= \left| x = r_{ci}^{-0.5} (\tilde{V}_{t;\psi ci}^d)^{0.5} (\Theta_{ci} - \hat{\hat{\Theta}}_{t;ci}^d), \ dx = r_{ci}^{-0.5|\psi_{ci}|} \left| \tilde{V}_{t;\psi ci}^d \right|^{0.5} d\Theta_{ci} \right|$$

$$= \frac{1}{\left| \tilde{V}_{t;\psi ci}^d \right|^{0.5}} \int_{r_{ci}} \int_{\Theta_{ci}} \ln(r_{ci}) \left( H_t(\hat{\Theta}_t^d) + \frac{\partial H_t(\hat{\Theta}_t^d)}{\partial r_{ci}} (r_{ci} - \hat{r}_{t;ci}^d) \right) r_{ci}^{-0.5(\tilde{\nu}_{t;ci}^d + 2)} \exp \left\{ -\frac{\tilde{\Lambda}_{t;ci}^d}{2r_{ci}} \right\} \exp \left\{ -\frac{x^T x}{2} \right\} dx dr_{ci}$$

$$= \frac{(2\pi)^{0.5|\psi_{ci}|}}{\left| \tilde{V}_{t;\psi ci}^d \right|^{0.5}} \int_{r_{ci}} \ln(r_{ci}) \left( H_t(\hat{\Theta}_t^d) + \frac{\partial H_t(\hat{\Theta}_t^d)}{\partial r_{ci}} (r_{ci} - \hat{r}_{t;ci}^d) \right) r_{ci}^{-0.5(\tilde{\nu}_{t;ci}^d + 2)} \exp \left\{ -\frac{\tilde{\Lambda}_{t;ci}^d}{2r_{ci}} \right\} dr_{ci}$$

$$= -2 \frac{(2\pi)^{0.5|\psi_{ci}|}}{\left| \tilde{V}_{t;\psi ci}^d \right|^{0.5}} \frac{\partial}{\partial \tilde{\nu}_{t;ci}^d} \int_{r_{ci}} \left( H_t(\hat{\Theta}_t^d) + \frac{\partial H_t(\hat{\Theta}_t^d)}{\partial r_{ci}} (r_{ci} - \hat{r}_{t;ci}^d) \right) r_{ci}^{-0.5(\tilde{\nu}_{t;ci}^d + 2)} \exp \left\{ -\frac{\tilde{\Lambda}_{t;ci}^d}{2r_{ci}} \right\} dr_{ci}$$

$$= -2 \frac{(2\pi)^{0.5|\psi_{ci}|}}{\left| \tilde{V}_{t;\psi ci}^d \right|^{0.5}} H_t(\hat{\Theta}_t^d) \frac{\partial}{\partial \tilde{\nu}_{t;ci}^d} \int_0^\infty r_{ci}^{-0.5(\tilde{\nu}_{t;ci}^d + 2)} \exp \left\{ -\frac{\tilde{\Lambda}_{t;ci}^d}{2r_{ci}} \right\} dr_{ci}$$

$$\overset{(7.40)}{=} -2 \frac{(2\pi)^{0.5|\psi_{ci}|}}{\left| \tilde{V}_{t;\psi ci}^d \right|^{0.5}} H_t(\hat{\Theta}_t^d) \frac{\partial}{\partial \tilde{\nu}_{t;ci}^d} \left( \frac{\tilde{\Lambda}_{t;ci}^d}{2} \right)^{-0.5\tilde{\nu}_{t;ci}^d} \Gamma \left( \frac{\tilde{\nu}_{t;ci}^d}{2} \right)$$

$$= \frac{(2\pi)^{0.5|\psi_{ci}|}}{\left| \tilde{V}_{t;\psi ci}^d \right|^{0.5}} H_t(\hat{\Theta}_t^d) \left( \frac{\tilde{\Lambda}_{t;ci}^d}{2} \right)^{-0.5\tilde{\nu}_{t;ci}^d} \Gamma \left( \frac{\tilde{\nu}_{t;ci}^d}{2} \right) \left( \ln \left( 0.5 \tilde{\Lambda}_{t;ci}^d \right) - \psi \left( 0.5 \tilde{\nu}_{t;ci}^d \right) \right) \tag{7.44}$$

By using (7.44), the integral $L_{ci}^1$ (7.43) is computed as follows:

$$L_{ci}^1 \approx \sum_{d \in \mathbf{c}} \gamma_{t;d} \left[ H(\hat{\Theta}_t^d) \left( \ln \left( 0.5 \tilde{\Lambda}_{t;cj}^d \right) - \psi \left( 0.5 \tilde{\nu}_{t;cj}^d \right) \right) \right],$$

which proves (6.27).

## $L_{ci}^2$ (6.28)

By inserting (7.39) into (7.37) it is obtained:

$$L_{ci}^2 \approx \sum_{d \in \mathbf{c}} \text{Be}(\tilde{v}_t^d) \prod_{\substack{e \in \mathbf{c} \\ e \neq c}} \prod_{j=1}^{n_e} N(\tilde{v}_{t;ej}^d, \tilde{G}_{t;ej}^d) / P_e(\psi_e^c)$$

$$\times \int_{r_{ci}} \int_{\Theta_{ci}} \frac{1}{r_{ci}} \left( H_t(\hat{\Theta}_t^d) + \frac{\partial H_t(\hat{\Theta}_t^d)}{\partial r_{ci}} (r_{ci} - \hat{r}_{t;ci}^d) \right) r_{ci}^{-0.5(\nu_{ci} + |\psi_{ci}| + 2)} P(r_{ci}, \Theta_{ci} | \tilde{v}_{t;ci}^d, \tilde{V}_{t;ci}^d) d\Theta_{ci} dr_{ci} \tag{7.45}$$

Similarly to (7.44), the integral in (7.45) is derived as follows:

$$\int_{r_{ci}} \int_{\Theta_{ci}} \left[ \frac{1}{r_{ci}} \left( H_t(\hat{\Theta}_t^d) + \frac{\partial H_t(\hat{\Theta}_t^d)}{\partial r_{ci}} \left( r_{ci} - \hat{r}_{t;ci}^d \right) \right) r_{ci}^{-0.5(\nu_{ci}+|\psi_{ci}|+2)} \right.$$

$$\times \exp \left\{ -\frac{1}{2r_{ci}} \left( (\Theta_{ci} - \hat{\hat{\Theta}}_{t;ci}^d)^T \tilde{V}_{t;\psi ci}^d (\Theta_{ci} - \hat{\hat{\Theta}}_{t;ci}^d) + \tilde{\Lambda}_{t;ci}^d \right) \right\} d\Theta_{ci} dr_{ci} \bigg]$$

$$= \left| x = r_{ci}^{-0.5} (\tilde{V}_{t;\psi ci}^d)^{0.5} (\Theta_{ci} - \hat{\hat{\Theta}}_{t;ci}^d), \;\; dx = r_{ci}^{-0.5|\psi_{ci}|} \left| \tilde{V}_{t;\psi ci}^d \right|^{0.5} d\Theta_{ci} \right|$$

$$= \frac{(2\pi)^{0.5|\psi_{ci}|}}{\left| \tilde{V}_{t;\psi ci}^d \right|^{0.5}} \int_{r_{ci}} \frac{1}{r_{ci}} \left( H_t(\hat{\Theta}_t^d) + \frac{\partial H_t(\hat{\Theta}_t^d)}{\partial r_{ci}} \left( r_{ci} - \hat{r}_{t;ci}^d \right) \right) r_{ci}^{-0.5(\tilde{\nu}_{t;ci}^d+2)} \exp \left\{ -\frac{\tilde{\Lambda}_{t;ci}^d}{2r_{ci}} \right\} dr_{ci} \qquad (7.46)$$

$L_{ci}^2$ is then computed by using (7.46), (7.40) and (7.42) as follows:

$$L_{ci}^2 \approx \sum_{d \in \mathbf{c}} \gamma_{t;d} \left( H(\hat{\Theta}_t^d) \frac{\tilde{\nu}_{ci}^d}{\tilde{\Lambda}_{t;ci}} + \frac{\partial H(\hat{\Theta}_t^d)}{\partial r_{cj}} \frac{2}{\tilde{\nu}_{ci}^d + 2} \right),$$

which proves (6.28).

## $L_{ci}^3$ (6.29)

By inserting (7.39) into (7.38) it is obtained:

$$L_{ci}^3 \approx \sum_{d \in \mathbf{c}} \mathrm{Be}(\tilde{v}_t^d) \prod_{\substack{e \in \mathbf{c} \\ e \neq c}} \prod_{j=1}^{n_e} \mathrm{N}(\tilde{v}_{t;ej}^d, \tilde{G}_{t;ej}^d) / \mathrm{P}_e(\psi_e^{\mathbf{c}}) \int_{r_{ci}} \int_{\Theta_{ci}} \left[ \frac{1}{r_{ci}} \left( (\Theta_{ci} - \hat{\Theta}_{ci})^T V_{\psi ci} (\Theta_{ci} - \hat{\Theta}_{ci}) \right) \right.$$

$$\left( H_t(\hat{\Theta}_t^d) + \frac{\partial H_t(\hat{\Theta}_t^d)}{\partial r_{ci}} \left( r_{ci} - \hat{r}_{t;ci}^d \right) + \frac{\partial H_t(\hat{\Theta}_t^d)}{\partial \Theta_{ci}} \left( \Theta_{ci} - \hat{\Theta}_{t;ci}^d \right) \right) \mathrm{P}(r_{ci}, \Theta_{ci} | \tilde{v}_{t;ci}^d, \tilde{V}_{t;ci}^d) d\Theta_{ci} dr_{ci} \bigg]$$

$$(7.47)$$

The integration in (7.47) is done in the following steps:

$$\int_{r_{ci}} \int_{\Theta_{ci}} \left[ \frac{1}{r_{ci}} \left( (\Theta_{ci} - \hat{\Theta}_{ci})^T V_{\psi ci} (\Theta_{ci} - \hat{\Theta}_{ci}) \right) \left( H_t(\hat{\Theta}_t^d) + \frac{\partial H_t(\hat{\Theta}_t^d)}{\partial r_{ci}} \left( r_{ci} - \hat{r}_{t;ci}^d \right) + \frac{\partial H_t(\hat{\Theta}_t^d)}{\partial \Theta_{ci}} \left( \Theta_{ci} - \hat{\Theta}_{t;ci}^d \right) \right) \right.$$

$$r_{ci}^{-0.5(\nu_{ci}+|\psi_{ci}|+2)} \exp \left\{ -\frac{1}{2r_{ci}} \left( (\Theta_{ci} - \hat{\hat{\Theta}}_{t;ci}^d)^T \tilde{V}_{t;\psi ci}^d (\Theta_{ci} - \hat{\hat{\Theta}}_{t;ci}^d) + \tilde{\Lambda}_{t;ci}^d \right) \right\} d\Theta_{ci} dr_{ci} \bigg]$$

$$= \left| x = r_{ci}^{-0.5} (\tilde{V}_{t;\psi ci}^d)^{0.5} (\Theta_{ci} - \hat{\hat{\Theta}}_{t;ci}^d), \;\; dx = r_{ci}^{-0.5|\psi_{ci}|} \left| \tilde{V}_{t;\psi ci}^d \right|^{0.5} d\Theta_{ci} \right|$$

$$= \frac{1}{\left| \tilde{V}_{t;\psi ci}^d \right|^{0.5}} \int_{r_{ci}} \int_{\Theta_{ci}} \left[ \frac{1}{r_{ci}} \left( \left( \hat{\hat{\Theta}}_{t;ci}^d - \hat{\Theta}_{ci} + r_{ci}^{0.5} \left( \tilde{V}_{t;\psi ci}^d \right)^{-0.5} x \right)^T V_{\psi ci} \left( \hat{\hat{\Theta}}_{t;ci}^d - \hat{\Theta}_{ci} + r_{ci}^{0.5} \left( \tilde{V}_{t;\psi ci}^d \right)^{-0.5} x \right) \right) \right.$$

$$\left( H_t(\hat{\Theta}_t^d) + \frac{\partial H_t(\hat{\Theta}_t^d)}{\partial r_{ci}} \left( r_{ci} - \hat{r}_{t;ci}^d \right) + \frac{\partial H_t(\hat{\Theta}_t^d)}{\partial \Theta_{ci}} \left( r_{ci}^{0.5} \left( \tilde{V}_{t;\psi ci}^d \right)^{-0.5} x \right) \right)$$

$$r_{ci}^{-0.5(\nu_{ci}+2)} \exp \left\{ -\frac{x^T x}{2} \right\} \exp \left\{ -\frac{\tilde{\Lambda}_{t;ci}^d}{2r_{ci}} \right\} dx dr_{ci} \bigg] = \mathrm{I} + \mathrm{II}. \qquad (7.48)$$

The integrals I and II are computed as follows:

$$
\mathrm{I} = \frac{1}{\left|\tilde{V}^d_{t;\psi ci}\right|^{0.5}} \int_{\boldsymbol{r}_{ci}} \int_{\boldsymbol{\Theta}_{ci}} \left[ \frac{1}{r_{ci}} \left( \left( \hat{\tilde{\Theta}}^d_{t;ci} - \hat{\Theta}_{ci} + r^{0.5}_{ci} \left( \tilde{V}^d_{t;\psi ci} \right)^{-0.5} x \right)^T V_{\psi ci} \left( \hat{\tilde{\Theta}}^d_{t;ci} - \hat{\Theta}_{ci} + r^{0.5}_{ci} \left( \tilde{V}^d_{t;\psi ci} \right)^{-0.5} x \right) \right) \right.
$$

$$
\left. \left( \mathrm{H}_t(\hat{\Theta}^d_t) + \frac{\partial \mathrm{H}_t(\hat{\Theta}^d_t)}{\partial r_{ci}} \left( r_{ci} - \hat{r}^d_{t;ci} \right) \right) r^{-0.5(\nu_{ci}+2)}_{ci} \exp\left\{ -\frac{x^T x}{2} \right\} \exp\left\{ -\frac{\tilde{\Lambda}^d_{t;ci}}{2r_{ci}} \right\} dx dr_{ci} \right]
$$

$$
= \frac{(2\pi)^{0.5|\psi_{ci}|}}{\left|\tilde{V}^d_{t;\psi ci}\right|^{0.5}} \int_{\boldsymbol{r}_{ci}} \left[ \left( \frac{1}{r_{ci}} \left( \hat{\tilde{\Theta}}^d_{t;ci} - \hat{\Theta}_{ci} \right)^T V_{\psi ci} \left( \hat{\tilde{\Theta}}^d_{t;ci} - \hat{\Theta}_{ci} \right) + \mathrm{tr}\left( \left( \tilde{V}^d_{t;\psi ci} \right)^{-1} V_{\psi ci} \right) \right) \right.
$$

$$
\left. \left( \mathrm{H}_t(\hat{\Theta}^d_t) + \frac{\partial \mathrm{H}_t(\hat{\Theta}^d_t)}{\partial r_{ci}} \left( r_{ci} - \hat{r}^d_{t;ci} \right) \right) r^{-0.5(\nu_{ci}+2)}_{ci} \exp\left\{ -\frac{\tilde{\Lambda}^d_{t;ci}}{2r_{ci}} \right\} dr_{ci} \right]
$$

$$
\overset{(7.40),(7.42)}{=} \frac{(2\pi)^{0.5|\psi_{ci}|}}{\left|\tilde{V}^d_{t;\psi ci}\right|^{0.5}} \left( \frac{\tilde{\Lambda}^d_{t;ci}}{2} \right)^{-0.5\tilde{\nu}^d_{t;ci}} \Gamma\left( \frac{\tilde{\nu}^d_{t;ci}}{2} \right) \left\{ \mathrm{H}(\hat{\Theta}^d_t) \mathrm{tr}\left( \left( \tilde{V}^d_{\psi cj} \right)^{-1} V_{\psi cj} \right) \right.
$$

$$
\left. + \left( \hat{\tilde{\Theta}}^d_{t;cj} - \hat{\Theta}_{t;cj} \right)^T V_{\psi cj} \left( \hat{\tilde{\Theta}}^d_{t;cj} - \hat{\Theta}_{t;cj} \right) \left( \mathrm{H}(\hat{\Theta}^d_t) \frac{\tilde{\nu}^d_{cj}}{\tilde{\Lambda}^d_{t;cj}} + \frac{\partial \mathrm{H}(\hat{\Theta}^d_t)}{\partial r_{cj}} \frac{2}{\tilde{\nu}^d_{cj}+2} \right) \right\}, \qquad (7.49)
$$

$$
\mathrm{II} = \frac{1}{\left|\tilde{V}^d_{t;\psi ci}\right|^{0.5}} \int_{\boldsymbol{r}_{ci}} \int_{\boldsymbol{\Theta}_{ci}} \left[ \frac{1}{r_{ci}} \left( \left( \hat{\tilde{\Theta}}^d_{t;ci} - \hat{\Theta}_{ci} + r^{0.5}_{ci} \left( \tilde{V}^d_{t;\psi ci} \right)^{-0.5} x \right)^T V_{\psi ci} \left( \hat{\tilde{\Theta}}^d_{t;ci} - \hat{\Theta}_{ci} + r^{0.5}_{ci} \left( \tilde{V}^d_{t;\psi ci} \right)^{-0.5} x \right) \right) \right.
$$

$$
\left. \frac{\partial \mathrm{H}_t(\hat{\Theta}^d_t)}{\partial \Theta_{ci}} \left( r^{0.5}_{ci} \left( \tilde{V}^d_{t;\psi ci} \right)^{-0.5} x \right) r^{-0.5(\nu_{ci}+2)}_{ci} \exp\left\{ -\frac{x^T x}{2} \right\} \exp\left\{ -\frac{\tilde{\Lambda}^d_{t;ci}}{2r_{ci}} \right\} dx dr_{ci} \right]
$$

$$
= \frac{1}{\left|\tilde{V}^d_{t;\psi ci}\right|^{0.5}} \frac{\partial \mathrm{H}_t(\hat{\Theta}^d_t)}{\partial \Theta_{ci}} \int_{\boldsymbol{r}_{ci}} \int_{\boldsymbol{\Theta}_{ci}} \left[ \left( \tilde{V}^d_{t;\psi ci} \right)^{-0.5} x \left( x^T \left( \left( \tilde{V}^d_{t;\psi ci} \right)^{-0.5} \right)^T V_{\psi ci} \left( \hat{\tilde{\Theta}}^d_{t;ci} - \hat{\Theta}_{ci} \right) + \left( \hat{\tilde{\Theta}}^d_{t;ci} - \hat{\Theta}_{ci} \right)^T V_{\psi ci} \left( \tilde{V}^d_{t;\psi ci} \right)^{-0.5} x \right) \right.
$$

$$
\left. r^{-0.5(\nu_{ci}+2)}_{ci} \exp\left\{ -\frac{x^T x}{2} \right\} \exp\left\{ -\frac{\tilde{\Lambda}^d_{t;ci}}{2r_{ci}} \right\} dx dr_{ci} \right]
$$

$$
= \frac{(2\pi)^{0.5|\psi_{ci}|}}{\left|\tilde{V}^d_{t;\psi ci}\right|^{0.5}} \frac{\partial \mathrm{H}_t(\hat{\Theta}^d_t)}{\partial \Theta_{ci}} \left( \left( \tilde{V}^d_{t;\psi ci} \right)^{-1} \left( V_{\psi ci} + V^T_{\psi ci} \right) \left( \hat{\tilde{\Theta}}^d_{t;ci} - \hat{\Theta}_{ci} \right) \right) \int_{\boldsymbol{r}_{ci}} r^{-0.5(\nu_{ci}+2)}_{ci} \exp\left\{ -\frac{\tilde{\Lambda}^d_{t;ci}}{2r_{ci}} \right\} dr_{ci}
$$

$$
\overset{(7.40)}{=} \frac{(2\pi)^{0.5|\psi_{ci}|}}{\left|\tilde{V}^d_{t;\psi ci}\right|^{0.5}} \left( \frac{\tilde{\Lambda}^d_{t;ci}}{2} \right)^{-0.5\tilde{\nu}^d_{t;ci}} \Gamma\left( \frac{\tilde{\nu}^d_{t;ci}}{2} \right) \frac{\partial \mathrm{H}_t(\hat{\Theta}^d_t)}{\partial \Theta_{ci}} \left( \left( \tilde{V}^d_{t;\psi ci} \right)^{-1} \left( V_{\psi ci} + V^T_{\psi ci} \right) \left( \hat{\tilde{\Theta}}^d_{t;ci} - \hat{\Theta}_{ci} \right) \right). \qquad (7.50)
$$

By combining (7.47), (7.48), (7.49) and (7.50), $L_{ci}^3$ is obtained as follows:

$$L_{cj}^3 \approx \sum_{d \in \mathbf{c}} \gamma_{t;d} \Bigg[ H(\hat{\Theta}_t^d) \mathrm{tr}\left( \left( \tilde{V}_{\psi cj}^d \right)^{-1} V_{\psi cj} \right)$$
$$+ \left( \hat{\bar{\Theta}}_{t;cj}^d - \hat{\Theta}_{t;cj} \right)^T V_{\psi cj} \left( \hat{\bar{\Theta}}_{t;cj}^d - \hat{\Theta}_{t;cj} \right) \left( H(\hat{\Theta}_t^d) \frac{\tilde{v}_{cj}^d}{\tilde{\Lambda}_{t;cj}^d} + \frac{\partial H(\hat{\Theta}_t^d)}{\partial r_{cj}} \frac{2}{\tilde{v}_{cj}^d + 2} \right)$$
$$+ \frac{\partial H(\hat{\Theta}_t^d)}{\partial \Theta_{cj}} \left( \tilde{V}_{\psi cj}^d \right)^{-1} \left( V_{\psi cj} + V_{\psi cj}^T \right) \left( \hat{\bar{\Theta}}_{t;cj}^d - \hat{\Theta}_{t;cj} \right) \Bigg],$$

which proves (6.29).

# Bibliography

[1] J.M. Bernardo. *Expected information as expected utility*. The Annals of Statistics, 7(3): 686-690, 1979.

[2] M. Kárný and T.V. Guy. *On support of imperfect Bayesian participants*. In T.V. Guy, M. Kárný, and D.H. Wolpert, editors, Decision Making with Imperfect Decision Makers, volume 28. Springer, Berlin, 2012. Intelligent Systems Reference Library.

[3] M. Kárný. *Approximate Bayesian recursive estimation*. Information Sciences, 289:100-111, 2014.

[4] G. J. McLachlan and D. Peel. *Finite Mixture Models*. Wiley Series in Probability and Statistics. Wiley, New York, 2000.

[5] R.B. Nelsen. *An Introduction to Copulas*. Springer, New York, 1999.

[6] K. Dedecius, I. Nagy, M. Kárný, and L. Pavelková. *Parameter estimation with partial forgetting method*. In Proc. of the 15th IFAC Symposium on Identification and System Parameter Estimation - SYSID, 2009.

[7] O. Barndorff-Nielsen. *Information and Exponential Families in Statistical Theory*. Wiley, NY, 1978.

[8] R.E. Bellman, *Adaptive Control Processes*, Princeton University Press, NJ, 1961.

[9] L. Berec and M. Kárný. Identification of reality in Bayesian context. In K. Warwick and M. Kárný, editors, *Computer-Intensive Methods in Control and Signal Processing*, pages 181–193. Birkhäuser, 1997.

[10] J.O. Berger. *Statistical Decision Theory and Bayesian Analysis*. Springer, NY, 1985.

[11] C. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.

[12] T. Bohlin. *Interactive System Identification: Prospects and Pitfalls*. Springer, NY, 1991.

[13] P. Guan, M. Raginsky, and R. Willett. Online Markov decision processes with Kullback-Leibler control cost. In *American Control Conference*, pages 1388–1393. IEEE, 2012.

[14] I. Guyon, A. Saffari, G. Dror, and G. Cawley. Model selection: Beyond the Bayesian/frequentist divide. *Journal of Machine Learning Research*, 11:61–87, 2010.

[15] S. Haykin. *Neural Networks: A Comprehensive Foundation*. Macmillan, NY, 1994.

[16] G. Holmes and T.Y. Liu, editors. *Proceedings of 7th Asian Conference on Machine Learning (ACML2015), JMLR Workshop and Conference Proceedings*, volume 45, 2015.

[17] A.M. Jazwinski. *Stochastic Processes and Filtering Theory*. Academic Press, NY, 1970.

[18] M. Kárný. Recursive estimation of high-order Markov chains: Approximation by finite mixtures. *Infor. Sciences*, 326:188 – 201, 2016.

[19] M. Kárný, J. Böhm, T. V. Guy, L. Jirsa, I. Nagy, P. Nedoma, and L. Tesař. *Optimized Bayesian Dynamic Advising: Theory and Algorithms*. Springer, 2006.

[20] D.F. Kerridge. Inaccuracy and inference. *J. of the Royal Statistical Society*, B 23:284–294, 1961.

[21] R. Koopman. On distributions admitting a sufficient statistic. *Trans. of Am. Math. Society*, 39:399, 1936.

[22] W. Mason, J.W. Vaughan, and H. Wallach. Special issue: Computational social science and social computing. *Machine Learning*, 96:257–469, 2014.

[23] T. Mitchell. *Machine Learning*. McGraw Hill, 1997.

[24] P. Sadghi, R.A. Kennedy, P.B. Rapajic, and R. Shams. Finite-state Markov modeling of fading channels. *IEEE Signal Processing Magazine*, 57, 2008.

[25] L.J. Savage. *Foundations of Statistics*. Wiley, NY, 1954.

[26] A. Wald. *Statistical Decision Functions*. John Wiley, New York, London, 1950.

[27] M. Wiering and M. van Otterlo, editors. *Reinforcement Learning: State-of-the-Art*. Springer-Verlag, 2012.

[28] I. Mező. *Some infinite sums arising from the Weierstrass Product Theorem*. Applied Mathematics and Computation. 219: 9838–9846, 2013.

[29] A.A. Markov. *Extension of the limit theorems of probability theory to a sum of variables connected in a chain*. reprinted in Appendix B of: R. Howard. Dynamic Probabilistic Systems, volume 1: Markov Chains. John Wiley and Sons, 1971.

[30] M. Kárný, M.Ruman *Mixture-Ratio Modeling of Dynamic Environments*. Submitted to ECML-PKDD-2018.

[31] M. Ruman, F. Hůla, M. Kárný, and T.V. Guy. *Deliberation-aware responder in multi- proposer ultimatum game*. In Artificial Neural Networks and Machine Learning - Proceedings ICANN 2016, pages 230-237. Barcelona, 2016.

[32] V. Peterka. Bayesian system identification. In P. Eykhoff, editor, *Trends and Progress in System Identification*, pages 239–304. Pergamon Press, Oxford, 1981.

[33] Kárný M., R. Kulhavý (1988). *Structure determination of regression-type models for adaptive prediction and control*. In: Spall J. C. (ed.): Bayesian Analysis of Time Series and Dynamic Models. Marcel Dekker, New York.

[34] Colin A. Carter (1999). *Commodity futures markets: a survey*. In Australian Journal of Agricultural and Resource Economics, vol. 43, pages 209-247.