DOI: xxx/xxxx

# **<u>RESEARCH ARTICLE</u> Mixture Ratio Modelling of Dynamic Systems**

Miroslav Kárný, Marko Ruman

The Czech Academy of Sciences, Institute of Information Theory and Automation, Pod vodárenskou věží 4, 182 08 Prague 8, Czech Republic, http://www.utia.cas.cz/AS Email: school@utia.cas.cz, marko.ruman@gmail.com Correspondence: Miroslav Kárný, school@utia.cas.cz

#### Summary

Any knowledge extraction relies (possibly implicitly) on a hypothesis about the modelled-data dependence. The extracted knowledge ultimately serves to a decision making (DM). DM always faces uncertainty and this makes probabilistic modelling adequate. The inspected black-box modelling deals with "universal" approximators of the relevant probabilistic model. Finite mixtures with components in the exponential family are often exploited. Their attractiveness stems from their flexibility, the cluster interpretability of components and the existence of algorithms for processing high-dimensional data streams. They are even used in dynamic cases with mutually dependent data records while regression and auto-regression mixture components serve to the dependence modelling. These dynamic models, however, mostly assume data-independent component weights, i.e. memoryless transitions between dynamic mixture components. Such mixtures are not universal approximators of dynamic probabilistic models. Formally, this follows from the fact that the set of finite probabilistic mixtures is not closed with respect to the conditioning, which is the key estimation and predictive operation. The paper overcomes this drawback by using ratios of finite mixtures as universally approximating dynamic parametric models. The paper motivates them, elaborates their approximate Bayesian recursive estimation and reveals their application potential.

#### **KEYWORDS:**

universal approximation, black-box dynamic model, mixture model, approximate Bayesian estimation, data stream processing, forgetting, Kullback-Leibler divergence

## **1** | INTRODUCTION

The paper serves to *dynamic* decision making (DM) understood as a targeted choice among available actions (options) influencing dynamically evolving system. DM seen in this way strongly overlaps with machine learning<sup>1</sup>, signal processing,<sup>2</sup>, hypothesis testing<sup>3</sup>, classification<sup>4</sup>, knowledge sharing<sup>5</sup>, reinforcement learning<sup>6</sup>, control<sup>7</sup>, adaptive control<sup>?</sup>, etc. All these areas benefit from an enrichment of the domain knowledge by that hidden in recorded data. This "data mining" is an extremely broad domain with many results and applications, e.g.<sup>8</sup>. Even surveys specialise, for instance, to Internet of Things<sup>9</sup>, cyber-security<sup>10</sup>, etc.

*Paper Focus:* The paper contributes to clustering of data streams<sup>11</sup> with dynamically bonded data records<sup>12</sup>. Importance of such a processing is well seen on specific applications, for instance in geophysics<sup>13</sup> or on a stream clustering of independent data records<sup>14</sup>. Our paper provides universal approximators<sup>1</sup> of dynamic models and their recursive estimation.

<sup>&</sup>lt;sup>1</sup>We adopt this intuitive notion used in connection with neural networks<sup>15</sup>. Formally, it means that universal approximators form a dense subset of a set of functions they approximate. Their rigorous specification in various sets of approximated functions are, for instance, in <sup>16,17,18,19</sup>. They contain a range of other common references.

*DM View on Data Stream Processing:* A solution of a DM task leads to a strategy, i.e. a collection of decision rules mapping the evolving knowledge on  $actions^{20}$ . A good prescriptive DM theory should provide a strategy, which meets user's DM aims in the best possible way. The adopted Bayesian DM<sup>21</sup> is such a theory coping with incomplete knowledge, uncertainty and randomness of the system to which actions relate. Bayesian DM links DM consequences to the acquired knowledge and actions by conditional distributions. Here, conditional probability densities (pd<sup>2</sup>) describe them.

Actions are chosen sequentially. This enables to operate on gradually enriched knowledge, to learn better the system model in a single pass only. Adaptive control<sup>223</sup> fully relies on this feature. Bayesian data stream processing updates a statistic, ideally a sufficient one. Thus, it may serve as a knowledge, feature or cluster extractor suitable to all areas mentioned above.

*Addressed Modelling Problem:* An estimation is possible iff the inspected relations hold during the knowledge accumulation. The most general case relies on invariant mappings describing stochastic state-space model and leads to stochastic filtering<sup>24</sup>. This paper deals with a simpler case<sup>3</sup> and uses parametric, black-box, probabilistic models relating a few adjacent data records. It is a priori unknown, which parameter points to the best model. Bayesian paradigm offers the unambiguous deductive estimation via Bayes' rule<sup>26</sup>. It accumulates the knowledge into the posterior pd and redistributes this belief about the model adequacy<sup>27</sup>.

The modelling and thus DM quality depend on the model set over which the belief redistributes. *For a fixed condition*, any smooth pd, describing data-records dependence, can be arbitrarily-well approximated by a finite mixture of conditional pds (components)<sup>18</sup>. Components may have the same functional form, say Gaussian one, but differ in their parameters. Various versions of this ability to approximate "universally" were proved in neural-network context<sup>28</sup>. Zero-memory models with a fixed trivial condition are universally approximable by mixtures with constant component weights<sup>29</sup>. In the inspected dynamic case with a non-trivial condition, the component weights should depend on it, too. The common use of the condition-independent weights of dynamic components goes against the expansion logic and surely violate universal approximation property. Rare exceptions that considered condition-dependent weights<sup>30,31,32</sup> indicate how much is lost when using constant component weights.

The paper offers *ratios* of finite mixtures with components in the exponential family (EF)<sup>33</sup> as black-box,<sup>34</sup>, universally approximating, dynamic models<sup>15</sup>. It develops, inevitably approximate, but feasible Bayesian recursive estimation of models from this *extremely rich but yet unconsidered model set*. Note that the thesis<sup>35</sup> and the conference paper<sup>36</sup> consider it but they are direct predecessors of this paper.

*Layout:* The paper relies on an approximate Bayesian recursive estimation<sup>37</sup>. Its recall in Sec. 2 makes the paper selfcontaining and prepares notations. Sec. 3 justifies the mixture ratio models and Sec. 4 elaborates their estimation. Sec. 5 illustrates the theory. Sec. 6 outlines the gained application potential.

*Common Notation:* The bold symbol X denotes a set of Xs. It is a subset of either finite-dimensional real space or it consists of pds. Its exact specification is only given when needed. |X| marks cardinality of X. Random variables, their realisations and function arguments are undistinguished. The meaning follows from the context. Mappings, marked by san serif fonts, are taken as different if labels of their arguments differ. Mnemonic symbols prevail: M means a parametric model, J is a joint parametric pd of data vectors, P is posterior pd, O marks observations, A labels actions, etc. Decoration  $\tilde{}$  marks intermediate objects,  $\hat{}$  marks estimates and approximations.  $\equiv$  means defining equality.

## 2 | APPROXIMATE BAYESIAN RECURSIVE ESTIMATION

This section summarises the approximate recursive estimation published in  $^{37}$ . It serves us for handling of the novel model proposed in Sec. 3.

The inspected parametric model M consists of the conditional pds at discrete-time moments labelled by  $t \in t = \{1, 2, ..., |t|\}$ 

$$\mathsf{M}(O_t|A_t, D^{t-1}, \Theta), \ D_t \equiv (O_t, A_t), \ D^{t-1} \equiv (D_{t-1}, \dots, D_1, D_0).$$
(1)

The pd (1) relates the observation  $O_t \in O$  to the action  $A_t \in A$  and to the past data records  $D^{t-1}$  extended by a prior knowledge  $D_0$ . A time-invariant  $\Theta \in \Theta$  parameterises these pds. The employed decision rules modelled by pds  $R(A_t | D^{t-1}, \Theta)$  are unaware of the adequate  $\Theta$ . Thus, these decision rules meet natural conditions of control<sup>38</sup>

$$\mathsf{R}(A_t | D^{t-1}, \Theta) = \mathsf{R}(A_t | D^{t-1}).$$
<sup>(2)</sup>

<sup>&</sup>lt;sup>2</sup>Pd is a common abbreviation covering the usual ones for probability density and mass functions (pdf and pmf). Pds are evaluated with respect to a dominating, typically Lebesque's or counting, product measure <sup>22</sup>. Lebesque's notation  $\int \dots d\bullet$  is mostly used here.

<sup>&</sup>lt;sup>3</sup>The results also suit to stochastic state-space models, aka hidden Markov models<sup>25</sup>. The addressed universal approximation by mixture ratio, which is closed with respect to conditioning, is applicable to them. The more complex notation and inevitable technical details related to them would, however, mask the basic modelling idea.

The posterior pd  $P_{t-1}(\Theta) \equiv P(\Theta|A_t, D^{t-1}) = P(\Theta|D^{t-1})$  quantifies the knowledge about the parameter  $\Theta \in \Theta$  available at time t - 1. Bayes' rule<sup>26</sup> updates this knowledge by the data-record realisation  $D_t \equiv (O_t, A_t)$  to the posterior pd

$$\tilde{\mathsf{P}}_{t}(\Theta) \equiv \frac{\mathsf{M}(O_{t}|A_{t}, D^{t-1}, \Theta)\mathsf{P}_{t-1}(\Theta)}{\mathsf{F}(O_{t}|A_{t}, D^{t-1})} \propto \mathsf{M}(O_{t}|A_{t}, D^{t-1}, \Theta)\mathsf{P}_{t-1}(\Theta)$$
$$\mathsf{F}(O_{t}|A_{t}, D^{t-1}) \equiv \int_{\Theta} \mathsf{M}(O_{t}|A_{t}, D^{t-1}, \Theta)\mathsf{P}_{t-1}(\Theta)\mathrm{d}\Theta,$$
(3)

where the value of the forecasting pd F is unexpressed when using the proportionality  $\propto$ . The decision rule R cancels due to (2).

With a growing *t*, the analytic form of the posterior pd  $\tilde{P}_t$  generically becomes an excessively complex function of many variables  $\Theta \in \Theta$ . Then,  $P_t \in \mathbf{P}$ , with  $\mathbf{P}$  containing computationally feasible<sup>4</sup> pds, has to be evolved. Mostly, the posterior pd  $\tilde{P}_t \notin \mathbf{P}$  even if the pd  $P_{t-1} \in \mathbf{P}$ . To stay feasible, the pd  $\tilde{P}_t$  in (3) is to be projected on  $\mathbf{P}$ . Works<sup>39,40</sup> provide weak, generically met, conditions under which the minimiser  $\hat{P}_t \in \mathbf{P}$  of Kerridge's inaccuracy<sup>41</sup>

$$\mathsf{K}(\tilde{\mathsf{P}}_{t}||\hat{\mathsf{P}}_{t}) \equiv -\int_{\Theta} \tilde{\mathsf{P}}_{t}(\Theta) \ln\left(\hat{\mathsf{P}}_{t}(\Theta)\right) \mathrm{d}\Theta \tag{4}$$

is the adequate Bayesian projection (minimising an adequate expected utility) of the pd  $\tilde{P}_t$  on  $\hat{P}_t \in \mathbf{P}$ .

The projection *should not* serve as the prior pd in a further updating. Its use could cause a divergence of the projected pds from those projected optimally in the batch mode<sup>42</sup>. Forgetting, with a data-dependent factor  $\lambda_t \in [0, 1]$ , is the adequate countermeasure<sup>37</sup>. An application of the minimum *expected* Kullback-Leibler principle<sup>43</sup> shows it. It recommends to select

$$\mathsf{P}_{t} = \arg\min_{\mathsf{P}\in\mathsf{P}}[\lambda_{t}\mathsf{D}(\mathsf{P}||\hat{\mathsf{P}}_{t}) + (1-\lambda_{t})\mathsf{D}(\mathsf{P}||\mathsf{P}_{t-1})] \propto \hat{\mathsf{P}}_{t}^{\lambda_{t}}\mathsf{P}_{t-1}^{1-\lambda_{t}}.$$
(5)

The Kullback-Leibler divergence  $^{44}$  D(P|| $\hat{P}$ )

$$\mathsf{D}(\mathsf{P}||\hat{\mathsf{P}}) \equiv \int_{\Theta} \mathsf{P}(\Theta) \ln\left(\frac{\mathsf{P}(\Theta)}{\hat{\mathsf{P}}(\Theta)}\right) d\Theta = \int_{\Theta} \mathsf{P}(\Theta) \ln\left(\mathsf{P}(\Theta)\right) d\Theta + \mathsf{K}(\mathsf{P}||\hat{\mathsf{P}}) \tag{6}$$

is a shifted version of Kerridge's inaccuracy (4). The weight  $\lambda_t \in [0, 1]$  is the belief that  $\hat{P}_t$  is the best prior guess of the unknown  $P_t$  and  $(1 - \lambda_t)$  is the belief into  $P_{t-1}$ . The beliefs are "naturally" proportional to forecasting-pds values

$$\hat{\mathsf{F}}(O_t|A_t, D^{t-1}) \equiv \int_{\Theta} \mathsf{M}(O_t|A_t, D^{t-1}, \Theta) \hat{\mathsf{P}}_t(\Theta) \mathrm{d}\Theta \text{ and } \mathsf{F}(O_t|A_t, D^{t-1}), \text{ cf. (3).}$$
(7)

A priori there are no reasons to expect  $\hat{P}_t$  be better than  $P_{t-1}$  as the knowledge innovation brought by  $D_t$  may be spoiled by the projection of  $\tilde{P}_t$  on  $\hat{P}_t \in \mathbf{P}$ . This motivates the used equal prior belief into  $\hat{P}_t$  and  $P_{t-1}$ .

Bayes' rule, the projection via Kerridge inaccuracy, minimum expected Kullback-Leibler principle and Bayes' rule updating beliefs into  $\hat{P}_t$ ,  $P_{t-1}$  provide the updating of a feasible posterior pd  $P_{t-1}$  by  $D_t = (O_t, A_t)$ 

$$\tilde{\mathsf{P}}_{t}(\Theta) \propto \mathsf{M}(O_{t}|A_{t}, D^{t-1}, \Theta)\mathsf{P}_{t-1}(\Theta), \quad \hat{\mathsf{P}}_{t} \equiv \arg\min_{\hat{\mathsf{P}} \in \mathsf{P}} \mathsf{K}(\tilde{\mathsf{P}}_{t}||\hat{\mathsf{P}})$$

$$\mathsf{P}_{t}(\Theta) \propto \hat{\mathsf{P}}_{t}^{\lambda_{t}}(\Theta)\mathsf{P}_{t-1}^{1-\lambda_{t}}(\Theta), \quad \lambda_{t} \equiv \left[1 + \mathsf{F}(O_{t}|A_{t}, D^{t-1})/\hat{\mathsf{F}}(O_{t}|A_{t}, D^{t-1})\right]^{-1}, \quad \text{cf. (3),(7).}$$
(8)

#### Remarks

✓ It is worth stressing why the optimisation (4) runs over the second argument while (5) over the first argument. It corresponds with two conceptually different tasks. In the first one, an *approximation of a known* pd  $\tilde{P}_t$  is constructed. The complementary analyses in<sup>39</sup> and<sup>40</sup> recommend both the optimised functional and this order of arguments. The optimisation (5) *completes a partial knowledge about an unknown* pd  $P_t$ . Works<sup>45,40,43</sup> axiomatically justify that this completion should be done via minimum (expected) Kullback-Leibler principle that optimises over the first argument.

Note the referred works clarify the relation of the variational Bayes<sup>46</sup> and the expectation propagation<sup>47</sup> techniques.

 $\checkmark$  The employed log-convex P guarantees that P<sub>t</sub> stays in P. Otherwise, an extra projection has to follow the forgetting.

<sup>&</sup>lt;sup>4</sup>An intuitive understanding of this notion suffices. Sec. 3 provides example of such pds.

(9)

**a** 1

1 1

- ✓ Kerridge's inaccuracy K (4) has the same minimiser as Kulback-Leibler divergence but copes with pds P̃ having Dirac's constituents.
- ✓ The unavoidable projection of the pd  $\tilde{P}_t$  on the set **P** (8) is demanding but feasible, Sec. 3. The extra effort for evaluating  $\lambda_t$  is relatively small.
- ✓ The next implicit choice of the forgetting factor is possible and meaningful

$$\lambda_{t} = \left[1 + \mathsf{F}(O_{t}|A_{t}, D^{t-1})/\mathsf{F}_{\lambda_{t}}(O_{t}|A_{t}, D^{t-1})\right]^{-1}, \quad \mathsf{F}_{\lambda_{t}}(O_{t}|A_{t}, D^{t-1}) \equiv \int_{\Theta} \mathsf{M}(O_{t}|A_{t}, D^{t-1}, \Theta) \frac{\mathsf{P}_{t}^{\lambda_{t}}(\Theta)\mathsf{P}_{t-1}^{\lambda_{t}}(\Theta)}{\int_{\Theta} \hat{\mathsf{P}}_{t}^{\lambda_{t}}(\tilde{\Theta})\mathsf{P}_{t-1}^{1-\lambda_{t}}(\tilde{\Theta})\mathsf{d}\tilde{\Theta}} \mathsf{d}\Theta$$

It is here unused to avoid the related high computational costs.

## **3** | RATIO OF FINITE MIXTURES

This core section justifies a universal parametrisation of pds modelling dependent data records. They may be observed in a closed loop formed by a dynamic system and a randomised strategy that meets (2).

Markovian Modelling of a Joint PD: A joint parametric pd  $\tilde{J}(D^{|t|}|\Theta, D_0)$  of data-records sequences  $D^{|t|} = (D_t)_{t \in t}$  factorises<sup>38</sup>

parametric model decision rule

$$\tilde{\mathsf{J}}(D^{|t|}|\Theta, D_0) = \prod_{t \in t} \overbrace{\mathsf{M}(O_t|A_t, D^{t-1}, \Theta)} \overbrace{\mathsf{R}(A_t|D^{t-1})}^{\mathsf{M}(O_t|A_t, D^{t-1}, \Theta)}$$

Thus, the parametrisation and estimation only concern the system model.

Assumption 1 (Markov Time-Invariant Parametric Model). The system model is *time-invariant*, parameterised by a constant multivariate parameter  $\Theta \in \Theta$ . It is Markov model of an order  $n < \infty$ . It means

 $\mathsf{M}(O_t|A_t, D^{t-1}, \Theta) \equiv \mathsf{M}(O_t|\rho_t, \Theta)$ , with time invariant M. The regression vector, denoted  $\rho_t$ ,

(2)

 $\rho_t \text{ is a known function of } \rho_{t-1} \text{ and of } \begin{cases}
A_t, D_{t-1}, \dots, D_{t-n} & \text{if } 1 \leq n < \infty, \\
A_t \text{ or void} & \text{ for } n = 0.
\end{cases}$ 

Data entering the model M at time  $t \in t$  form the *data vector*  $\Psi_t \equiv (O_t, \rho_t) \in \Psi$ , extending the regression vector  $\rho_t$ .

The knowledge  $D_0$  provides the model structure and the regression vector  $\rho_0$ .

The Markov parametric model (9) is ratio of the joint pd  $J(\Psi|\Theta)$  and its marginal

$$\mathsf{M}(O_t|\rho_t,\Theta) = \frac{\mathsf{J}(O_t,\rho_t|\Theta)}{\int_{O}\mathsf{J}(O_t,\rho_t|\Theta)\mathsf{d}O_t} = \frac{\mathsf{J}(\Psi_t|\Theta)}{\int_{O}\mathsf{J}(O_t,\rho_t|\Theta)\mathsf{d}O_t}.$$
(10)

It is time invariant iff the joint pd  $J(\Psi_t | \Theta)$  is a time-invariant function multiplied by an arbitrary positive function of  $A_t$ ,  $D^{t-1}$ , e.g. any decision rule  $R(A_t | D^{t-1})$ . Thus, Assumption 1 practically makes the joint pd  $J(\Psi_t | \Theta)$  time-invariant.

#### Remark

✓ The Markov property is inevitable for the targeted feasibility of the recursive estimation. If not met naturally, it must be enforced via an approximation. The required time invariance can be relaxed by including time into  $\rho_t$  or by considering a time– and data– dependent unknown parameter  $\Theta$ . The latter case leads to the untreated hidden Markov models<sup>25</sup>.

Exponential Family: Our model exploits EF members described by the pd

$$\mathsf{M}(O_t|\rho_t,\Theta) \equiv \exp\left\langle \mathsf{B}(\Psi_t),\mathsf{C}(\Theta)\right\rangle, \ \Psi_t \in \Psi, \ \Theta \in \Theta, \ \text{where}$$
(11)

vector-valued functions  $B(\Psi_t)$ ,  $C(\Theta)$  have a finite dimension. The real-valued mapping  $\langle B(\Psi_t), C(\Theta) \rangle$  is *linear in*  $B(\Psi_t)$ -values.

Under (2), EF members posses conjugated (self-reproducing) pds<sup>26</sup>,

$$P_{t}(\Theta) \equiv P(\Theta|D^{t}) = P(\Theta|V_{t}) = \frac{\exp \langle V_{t}, C(\Theta) \rangle}{N(V_{t})}$$

$$N(V_{t}) \equiv \int_{\Theta} \exp \langle V_{t}, C(\Theta) \rangle \, d\Theta. \text{ The sufficient statistic } V_{t} \text{ evolves}$$

$$V_{t} = V_{t-1} + B(\Psi_{t}) \text{ for } t \in t \text{ while } V_{0} \text{ determines the prior pd.}$$

$$(12)$$

Remarks

- Definition (11) admits usual factors depending respectively on Ψ and Θ. It suffices to include constant entries into B and C. The support indicator is dropped for simplicity.
- $\checkmark$  The use of EF converts functional Bayes' rule into the algebraic evolution of the values of the sufficient statistic V<sub>t</sub>.
- ✓  $V_0$  in (12) gives the conjugated prior pd, stores a prior knowledge  $D_0$ , regularises the estimation and is to make  $N(V_0) < \infty$ .
- ✓ EF exhausts parametric models with  $\Theta$ -independent support, smooth in  $\Theta$  and having a finite-dimensional sufficient statistic  $V_t = V(D^t)^{48}$ .
- ✓ For truly dynamic parametric models with a non-constant regression vector, EF is quite narrow. It essentially contains normal linear-in-regression-coefficient models, for continuous observations, and Markov chain models, for discrete-valued observations and discrete-valued regression vectors. Only for them, the marginal pd of the regression vector, proportional to the normalisation in (10), depends on  $\Theta$  only (not on the regression vector  $\rho$ ). Static models in EF are much richer.

Universal Approximation: Smooth joint pds  $J(\Psi) = J(O, \rho)$  of data vectors  $\Psi \in \Psi$  can be universally approximated by a finite mixture of normal pds, to an arbitrary precision<sup>18</sup>. It holds even when the regression vectors  $\rho \in \rho$  are non-void. Thus, the joint pds of data vectors can be approximated by a finite mixture of pds from EF. This allows us to consider the joint parametric pd in (10) as the finite weighted sum of EF components, pds  $(J_c(\Psi|\Theta_c))_{c\in c}$  on  $\Psi$ ,

$$J(\Psi|\Theta) \equiv \sum_{c \in c} \alpha_c J_c(\Psi|\Theta_c) \equiv \sum_{c \in c} \alpha_c \exp \langle \mathsf{B}_c(\Psi), \mathsf{C}_c(\Theta_c) \rangle, \quad c \equiv \{1, \dots, |c|\}, \quad \alpha \equiv (\alpha_c)_{c \in c}$$
$$\alpha \in \alpha \equiv \left\{ \alpha_c \ge 0, \sum_{c \in c} \alpha_c = 1 \right\}, \quad \Theta \equiv (\alpha, (\Theta_c)_{c \in c}) \in \Theta \equiv \left( \alpha, (\Theta_c)_{c \in c} \right). \tag{13}$$

The insertion of (13) into (10) gives the non-standard parametric ratio model

$$\mathsf{M}(O_{t}|\rho_{t},\Theta) = \sum_{c\in c} \frac{\alpha_{c} \exp\left\langle\mathsf{B}_{c}(O_{t},\rho_{t}),\mathsf{C}_{c}(\Theta_{c})\right\rangle}{\sum_{\tilde{c}\in c} \alpha_{\tilde{c}} \int_{O} \exp\left\langle\mathsf{B}_{\tilde{c}}(O_{t},\rho_{t}),\mathsf{C}_{\tilde{c}}(\Theta_{\tilde{c}})\right\rangle \mathrm{d}O_{t}}$$

$$= \sum_{c\in c} \underbrace{\frac{\alpha_{c}\mathsf{N}_{c}(\rho_{t},\Theta_{c})}{\sum_{\tilde{c}\in c} \alpha_{\tilde{c}}\mathsf{N}_{\tilde{c}}(\rho_{t},\Theta_{\tilde{c}})}}_{\mathsf{W}_{c}(\rho_{t},\Theta_{c})} \underbrace{\frac{\exp\left\langle\mathsf{B}_{c}(\Psi_{t}),\mathsf{C}_{c}(\Theta_{c})\right\rangle}{\mathsf{N}_{c}(\rho_{t},\Theta_{c})}}_{\mathsf{M}_{c}(O_{t}|\rho_{t},\Theta_{c})} = \sum_{c\in c} \mathsf{w}_{c}(\rho_{t},\Theta)\mathsf{M}_{c}(O_{t}|\rho_{t},\Theta_{c}).$$

$$(14)$$

Assumption 1 and the *universal* approximation by finite mixtures for void  $\rho$  imply the quite strong property of the model (14):

The ratio (14) approximates any dynamic, Markov, time-invariant model of data.

#### Remarks

✓ The 2<sup>nd</sup> row of (14) represents the ratio model as a fully dynamic finite mixture. The component weights  $w_c(\rho_t, \Theta)$  depend on the regression vector  $\rho_t$  in *non-ambiguous way*, which needs *no extra parameter*.

✓ The components  $J_c(\Psi_t | \Theta_c) = \exp \langle B_c(\Psi_t), C_c(\Theta_c) \rangle$  in (13) are joint pds on the set  $\Psi$ . These pds may contain parameterindependent factors  $G_c$ 

$$\mathsf{J}_{c}(\Psi_{t}|\Theta_{c}) = \exp\left\langle\mathsf{B}_{c}(\Psi_{t}),\mathsf{C}_{c}(\Theta_{c})\right\rangle = \exp\left\langle\mathsf{B}_{c}(\Psi_{t;c}),\mathsf{C}_{c}(\Theta_{c})\right\rangle\mathsf{G}_{c}(\underline{\rho}_{t;c}).$$
(15)

The parametric factor models the data vector  $\Psi_{t;c} \equiv (O_t, \rho_{t;c})$ , where  $\rho_{t;c}$  is a sub-vector of  $\rho_t$ . The pd  $G_c(\rho_{t;c})$  is a parameterfree pd on the complement  $\rho_{t;c}$  of  $\rho_{t;c}$  to  $\rho_t$ . This mimics mixtures of principal component analysers<sup>49</sup> and diminishes the dimensionality curse<sup>50</sup>.

## 4 | ESTIMATION OF MIXTURE RATIOS

This section applies the approximate estimation (8) to the model (14). This yields the *universal and feasible probabilistic clustering of dynamic data streams*.

*Choice of the Set* **P** *of Feasible PDs:* The weights  $\alpha$  in (13) define the pd  $M(c_t = c | \Theta) = \alpha_c, c \in c$ , of a thought *unobserved* pointer,  $c_t \in c$ , to the active component  $J_{c_t}(\Psi_t | \Theta_{c_t})$  (15) "generating" the data vector  $\Psi_t$ . This model is from EF

$$\begin{split} \mathsf{M}(c_t|\Theta) &= \exp\left[\sum_{c \in c} \delta(c, c_t), \ln(\alpha_c)\right] = \exp\left\langle\delta(c_t), \ln(\alpha)\right\rangle \\ \delta(c, c_t) &\equiv \begin{cases} 1 & \text{if } c = c_t \\ 0 & \text{otherwise} \end{cases}, \ \delta(c_t) &\equiv [\delta(1, c_t), \dots, \delta(|c|, c_t)] \\ \ln(\alpha) &\equiv [\ln(\alpha_1), \dots, \ln(\alpha_{|c|})], \quad \left\langle\delta(c_t), \ln(\alpha)\right\rangle = \sum_{c \in c} \delta(c, c_t) \ln(\alpha_c). \end{split}$$

For the *observed* pointer  $c_t \in c$  to the active component, Dirichlet's pd<sup>51</sup> determined by a |c|-vector statistic  $v_t$ , is the self-reproducing pd

$$P_{t}(\alpha) \equiv P(\alpha|\mathbf{v}_{t}) \equiv \frac{\exp\left\langle \mathbf{v}_{t} - 1, \ln(\alpha)\right\rangle}{\mathsf{Be}(\mathbf{v}_{t})} \quad \text{with } \mathbf{v}_{t} = \mathbf{v}_{t-1} + \delta(c_{t}), \ \mathbf{v}_{0} > 0$$
$$\mathsf{Be}(\mathbf{v}) \equiv \frac{\prod_{c \in c} \Gamma(\mathbf{v}_{c})}{\Gamma\left(\sum_{c \in c} \mathbf{v}_{c}\right)}, \quad \Gamma(v) \equiv \int_{0}^{\infty} z^{v-1} \exp(-z) dz, \ v > 0.$$
(16)

For normal and Markov chain components or independent data vectors with components (15), the conjugate pds  $P_t(\Theta_c) = P(\Theta_c | D^t, c_t, \dots, c_1), c \in c$ , are

$$\mathsf{P}_{t}(\Theta_{c}) = \frac{\exp\left\langle\mathsf{V}_{t;c},\mathsf{C}_{c}(\Theta_{c})\right\rangle}{\mathsf{N}_{c}(\mathsf{V}_{t;c})}, \quad \mathsf{N}_{c}(\mathsf{V}_{t;c}) = \int_{\Theta_{c}} \exp\left\langle\mathsf{V}_{t;c},\mathsf{C}_{c}(\Theta_{c})\right\rangle \mathrm{d}\Theta_{c}.$$
(17)

They are "natural" candidates for creating the set **P** of feasible pds used in (8). They have to be given by *statistic*  $(v_t, (V_{t;c})_{c \in c})$  with values dependent on the observed data, not on the unobserved  $(c_t, \dots, c_1)$ .

The pds  $P_t(\alpha)$  (16),  $(P_t(\Theta_c))_{c \in c}$  (17) do not determine the joint pd  $P_t(\Theta)$ ,  $\Theta \equiv (\alpha, (\Theta_c)_{c \in c})$ , unambiguously<sup>52</sup>. Their product coupling is, however, preferable as it models the faced lack of information about mutual relations of the component parameters and their weights. This motivates the choice

$$\mathbf{P} \equiv \left\{ \mathsf{P}_{t}(\Theta) : \mathsf{P}_{t}(\Theta) = \mathsf{P}_{t}(\alpha) \prod_{c \in c} \mathsf{P}_{t}(\Theta_{c}) \text{ with factors given by (16), (17),} \right\}$$

which are given by the optional statistics values  $(\mathbf{v}_t, \mathbf{V}_t \equiv (\mathbf{V}_{t;c})_{c \in c})$ . (18)

*Evaluations of Kerridge's Inaccuracy* (4): With the chosen **P** (18), it remains to map (8) onto the updating of  $v_{t-1}$ ,  $V_{t-1} \equiv (V_{t-1;c})_{c \in c}$  to  $v_t$ ,  $V_t \equiv (V_{t;c})_{c \in c}$ .

Kerridge's inaccuracy of the pd  $\tilde{\mathsf{P}}_t$  (3) to a pd  $\hat{\mathsf{P}}_t \in \mathsf{P}$  (18), given by the statistic  $\hat{\mathsf{V}}_t \equiv (\hat{\mathsf{v}}_t, (\hat{\mathsf{V}}_{t;c})_{c \in c})$ , reads

$$K(\tilde{P}_{t}||\hat{P}_{t}) = \ln(\text{Be}(\hat{v}_{t})) + \sum_{c \in c} \ln(N_{c}(\hat{V}_{t;c}))$$

$$-\left\langle \hat{v}_{t} - 1, \int_{\Theta} \ln(\alpha)\tilde{P}_{t}(\Theta)d\Theta \right\rangle - \sum_{c \in c} \left\langle \hat{V}_{t;c}, \int_{\Theta} C_{c}(\Theta_{c})\tilde{P}_{t}(\Theta)d\Theta \right\rangle.$$
(19)

Bayes' rule (3) for the mixture ratio (14) with EF components (15) and a prior pd  $P_{t-1} \in P$  (18), gives

$$\tilde{\mathsf{P}}_{t}(\Theta) \propto \sum_{c \in c} \frac{\alpha_{c} \exp \langle \mathsf{B}_{c}(\Psi_{t}), \mathsf{C}_{c}(\Theta_{c}) \rangle \exp \langle \mathsf{v}_{t-1} - 1, \ln(\alpha) \rangle \prod_{\tilde{c} \in c} \exp \langle \mathsf{V}_{t-1;\tilde{c}}, \mathsf{C}_{\tilde{c}}(\Theta_{\tilde{c}}) \rangle}{\sum_{\tilde{c} \in c} \alpha_{\tilde{c}} \mathsf{N}_{\tilde{c}}(\rho_{t}, \Theta_{\tilde{c}})}$$

$$= \sum_{c \in c} \frac{\exp \langle \mathsf{v}_{t-1} + \delta(c) - 1, \ln(\alpha) \rangle}{\sum_{\tilde{c} \in c} \alpha_{\tilde{c}} \mathsf{N}_{\tilde{c}}(\rho_{t}, \Theta_{\tilde{c}})}$$

$$\times \prod_{\tilde{c} \in c} \exp \langle \mathsf{V}_{t-1;\tilde{c}} + \delta(c, \tilde{c}) \mathsf{B}_{c}(\Psi_{t;c}), \mathsf{C}_{\tilde{c}}(\Theta_{c}) \rangle \mathsf{G}_{c}(\underline{\rho}_{t;c}).$$
(20)

Further evaluations exploit auxiliary pds arising from summands in (20), cf. (16), (17), via an appropriate normalisation

$$\tilde{\mathsf{Q}}_{t;c}(\Theta) \equiv \frac{\exp\left\langle \mathsf{v}_{t-1} + \delta(c) - 1, \ln(\alpha) \right\rangle}{\beta_{t;c}} \prod_{\tilde{c} \in c} \exp\left\langle \mathsf{V}_{t-1;\tilde{c}} + \delta(c, \tilde{c})\mathsf{B}_{c}(\Psi_{t;c}), \mathsf{C}_{\tilde{c}}(\Theta_{\tilde{c}}) \right\rangle$$
$$\beta_{t;c} \equiv \mathsf{Be}(\mathsf{v}_{t-1} + \delta(c)) \prod_{\tilde{c} \in c} \mathsf{N}(\mathsf{V}_{t-1;\tilde{c}} + \delta(c, \tilde{c})\mathsf{B}_{c}(\Psi_{t;c})) / \mathsf{G}_{c}(\underline{\rho}_{-t;c}). \tag{21}$$

They independently model  $\alpha$  and  $(\Theta_c)_{c \in c}$ . The dependency enters  $\tilde{\mathsf{P}}_t$  via

$$H(\rho_{t},\Theta) \equiv \frac{1}{\sum_{\tilde{c}\in c} \alpha_{\tilde{c}} \mathsf{N}_{\tilde{c}}(\rho_{t},\Theta)} \text{ because } \tilde{\mathsf{P}}_{t}(\Theta) = \frac{\mathsf{H}(\rho_{t},\Theta)\sum_{c\in c} \beta_{t;c}\tilde{\mathsf{Q}}_{t;c}(\Theta)}{\mathsf{I}_{-1}(\Psi_{t},\mathsf{v}_{t-1},\mathsf{V}_{t-1})}$$
  
with 
$$\mathsf{I}_{-1}(\Psi_{t},\mathsf{v}_{t-1},\mathsf{V}_{t-1}) \equiv \int_{\Theta} \mathsf{H}(\rho_{t},\Theta)\sum_{c\in c} \beta_{t;c}\tilde{\mathsf{Q}}_{t;c}(\Theta)\mathsf{d}\Theta.$$
(22)

The determination of Kerridge's inaccuracy (19) needs to evaluate the integrals

$$I_{0c}(\Psi_{t}, \mathsf{v}_{t-1}, \mathsf{V}_{t-1}) \equiv \int_{\Theta} \ln(\alpha_{c}) \mathsf{H}(\rho_{t}, \Theta) \sum_{\tilde{c} \in c} \beta_{t;\tilde{c}} \tilde{\mathsf{Q}}_{t;\tilde{c}}(\Theta) \mathsf{d}\Theta, \quad I_{0} \equiv (I_{0c})_{c \in c},$$

$$I_{c}(\Psi_{t}, \mathsf{v}_{t-1}, \mathsf{V}_{t-1}) \equiv \int_{\Theta} \mathsf{C}_{c}(\Theta) \mathsf{H}(\rho_{t}, \Theta) \sum_{\tilde{c} \in c} \beta_{t;\tilde{c}} \tilde{\mathsf{Q}}_{t;\tilde{c}}(\Theta) \mathsf{d}\Theta, \quad c \in c.$$
(23)

Their insertion into Kerridge's inaccuracy (19) gives the best-projection statistic

$$\begin{split} \left[ \hat{\mathbf{v}}_{t}, (\hat{\mathbf{V}}_{t;c})_{c \in c} \right] &\in \operatorname{Arg} \min_{\hat{\mathbf{v}} \in \hat{\mathbf{v}}, (\hat{\mathbf{v}}_{c} \in \hat{\mathbf{V}}_{c})_{c \in c}} \left( \ln(\operatorname{\mathsf{Be}}(\hat{\mathbf{v}})) + \sum_{c \in c} \ln(\operatorname{\mathsf{N}}_{c}(\hat{\mathbf{V}}_{c})) \\ &- \left\langle \hat{\mathbf{v}} - 1, \frac{\operatorname{I}_{0}(\Psi_{t}, \mathsf{v}_{t-1}, \mathsf{V}_{t-1})}{\operatorname{I}_{-1}(\Psi_{t}, \mathsf{v}_{t-1}, \mathsf{V}_{t-1})} \right\rangle - \sum_{c \in c} \left\langle \hat{\mathbf{V}}_{c}, \frac{\operatorname{I}_{c}(\Psi_{t}, \mathsf{v}_{t-1}, \mathsf{V}_{t-1})}{\operatorname{I}_{-1}(\Psi_{t}, \mathsf{v}_{t-1}, \mathsf{V}_{t-1})} \right\rangle \right) \\ &= \left[ \operatorname{Arg} \min_{\hat{\mathbf{v}} \in \hat{\mathbf{v}}} \left( \ln(\operatorname{\mathsf{Be}}(\hat{\mathbf{v}})) - \left\langle \hat{\mathbf{v}} - 1, \frac{\operatorname{I}_{0}(\Psi_{t}, \mathsf{v}_{t-1}, \mathsf{V}_{t-1})}{\operatorname{I}_{-1}(\Psi_{t}, \mathsf{v}_{t-1}, \mathsf{V}_{t-1})} \right\rangle \right) \right) \\ &- \left( \operatorname{Arg} \min_{\hat{\mathbf{v}}_{c} \in \hat{\mathbf{V}}_{c}} \left( \ln(\operatorname{\mathsf{N}}_{c}(\hat{\mathbf{V}}_{c})) - \left\langle \hat{\mathbf{V}}_{c}, \frac{\operatorname{I}_{c}(\Psi_{t}, \mathsf{v}_{t-1}, \mathsf{V}_{t-1})}{\operatorname{I}_{-1}(\Psi_{t}, \mathsf{v}_{t-1}, \mathsf{V}_{t-1})} \right\rangle \right) \right) \right]_{c \in c} \right]. \end{split}$$

It needs (|c| + 1) independent minimisations. They are numerically simple. Parts of them have analytical solutions for normal and Markov chain models.

*Numerical Evaluations of* I's: The function  $H(\rho_t, \Theta)$  (22) depends on all involved parameters. It combines influences of respective components. A brute force evaluation of  $I_{-1}$ ,  $I_0$  and  $I_c$ ,  $c \in c$ , say by Monte Carlo (MC), is mostly too demanding in a typical stream processing. At the same time, the averaged functions in (22), (23) are smooth due to their construction. Moreover, their growth over their domains is well suppressed by generally light tails of the pds  $\tilde{Q}_{t,c}$  (21). Thus, we conjecture (and experiments support this conjecture) that the simplest approximation based on the first order Taylor expansion of  $H(\rho, \Theta)$  around expected values

$$\int_{\Theta} \alpha_{\tilde{c}} \tilde{\mathsf{Q}}_{t;c}(\Theta) \mathrm{d}\Theta, \qquad \int_{\Theta} \Theta_{\tilde{c}} \tilde{\mathsf{Q}}_{t;c}(\Theta) \mathrm{d}\Theta, \qquad \tilde{c}, c \in \boldsymbol{c},$$
(25)

suffices. In addition to (25), the approximation requires to evaluate  $\forall \tilde{c}, c \in c$ 

$$\int_{\Theta} \alpha \ln(\alpha_{\tilde{c}}) \tilde{\mathsf{Q}}_{t;c}(\Theta) \mathrm{d}\Theta, \quad \int_{\Theta} \mathsf{C}_{\tilde{c}}(\Theta_{\tilde{c}}) \tilde{\mathsf{Q}}_{t;c}(\Theta) \mathrm{d}\Theta, \quad \int_{\Theta} \Theta_{\tilde{c}} \mathsf{C}_{\tilde{c}}(\Theta_{\tilde{c}}) \tilde{\mathsf{Q}}_{t;c}(\Theta) \mathrm{d}\Theta.$$
(26)

#### Remarks

- ✓ The normalising constants Be, N (16), (17) of factors forming  $\tilde{Q}_{t;c}$  (21) can be analytically expressed for Normal-inverse-Wishart and Dirichlet's pds, which are conjugated to normal and Markov chain components<sup>35</sup>.
- ✓  $\tilde{c}$ th entries of expectations (25), (26) at time *t* coincide with those gained at time *t* 1 if  $\tilde{c} \neq c$ . This significantly reduces the computational load.
- The used approximation could be refined but often suffices. Its use allowed us to focus on the ratio model, on the key novelty brought.
- This estimation applied to standard mixtures reduces to the projection-based algorithm<sup>37</sup>, which is still the state-of-the-art recursive Bayesian estimator.

*Forgetting Handling:* It consists of the choice of the forgetting factor and its use. The formula for  $\lambda_t$  in (8) depends on the ratio of the forecasting pd divided by the forecasting pd constructed from the already updated posterior pd (both in the data realisation). Independent studies<sup>53</sup> imply that the forgetting should be applied component-wise. The overall forecasting quality then influences only the forgetting of component weights. It gives the forgetting factors  $\lambda_t$ ,  $(\lambda_{t,c})_{c \in c}$ , (8),

$$\begin{split} \lambda_{t}^{-1} &= 1 + \frac{\mathsf{Be}(\hat{v}_{t}) \prod_{\tilde{c} \in c} \mathsf{N}(\mathsf{V}_{t;\tilde{c}})}{\mathsf{Be}(\mathsf{v}_{t-1}) \prod_{\tilde{c} \in c} \mathsf{N}(\mathsf{V}_{t-1;\tilde{c}})} \frac{\mathsf{I}_{-1}(\Psi_{t}, \mathsf{v}_{t-1}, \mathsf{V}_{t-1})}{\mathsf{I}_{-1}(\Psi_{t}, \hat{v}_{t}, \hat{V}_{t})} \\ \lambda_{t;c}^{-1} &= 1 + \frac{\mathsf{N}(\mathsf{V}_{t-1;c} + \mathsf{B}_{c}(\Psi_{t;c}))\mathsf{N}(\hat{\mathsf{V}}_{t;c})}{\mathsf{N}(\mathsf{V}_{t-1;c})\mathsf{N}(\hat{\mathsf{V}}_{t;c} + \mathsf{B}_{c}(\Psi_{t;c}))} \frac{\mathsf{N}_{c}(\rho_{t}, \tilde{\Theta}_{t;c})}{\mathsf{N}_{c}(\rho_{t}, \tilde{\Theta}_{t-1;c})}, \ c \in c. \end{split}$$

There  $\mathbf{I}_{-1}(\Psi_t, \mathbf{v}_{t-1}, \mathsf{V}_{t-1})$  is approximately evaluated using an expansion around (25). The denominator  $\mathbf{I}_{-1}(\Psi_t, \hat{\mathbf{v}}_t, \hat{\mathbf{V}}_t)$  is computed according to the same formulae with  $\hat{\mathbf{v}}_t, \hat{\mathbf{V}}_t$  (24) replacing  $\mathbf{v}_{t-1}, \mathbf{V}_{t-1}$ .  $\tilde{\Theta}_{t-1;c}$ ,  $\tilde{\Theta}_{t;c}$  are the expected values of the component parameter  $\Theta_c$  computed as the *c*th part of (25) with the component statistics  $\mathbf{V}_{t-1;c}$  and  $\hat{\mathbf{V}}_{t;c}$ , respectively (14). The forgetting completes the updating (8) for the proposed mixture ratio model (14), (15)

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

## **5** | ILLUSTRATIVE EXAMPLES

This section deals both with the modelling and the stream estimation. The used toy examples serve to the aimed illustration of the theory. A systematic simulation and real-life studies are out of scope of this paper and will be published independently. Some experiments, including real-data application to commodity futures<sup>54</sup> are already available in<sup>35,36</sup>.

The first example illustrates the modelling potential of mixture ratios. The second one shows that the *standard mixture model is not good when the mixture ratio model is adequate while the mixture ratio is competitive even when the standard mixture model is adequate.* 

**On Mixture-Ratio Modelling Strength:** The simulated system E generates real scalar observations  $O_t$  with no actions and the regression vector  $\rho_t = O_{t-1}$ . E is the mixture of the joint normal pds  $\mathcal{N}_{\Psi}(\mu_c, \omega_c^{0.5}), \Psi_t \equiv (O_t, O_{t-1})$ , with expectations  $\mu_c$  and square roots  $\omega_c^{0.5}$  of precision matrices,  $c \in c \equiv \{1, 2\}$ ,

$$\mathsf{E}(\Psi_{t}) = \frac{1}{2} \mathcal{N}_{\Psi_{t}} \left( \underbrace{\begin{bmatrix} 1\\1 \end{bmatrix}}_{\mu_{1}}, \underbrace{\begin{bmatrix} 1&3\\0&0.5 \end{bmatrix}}_{\omega_{1}^{0.5}} \right) + \frac{1}{2} \mathcal{N}_{\Psi_{t}} \left( \underbrace{\begin{bmatrix} -1\\-1 \end{bmatrix}}_{\mu_{2}}, \underbrace{\begin{bmatrix} 1&-2\\0&0.5 \end{bmatrix}}_{\omega_{2}^{0.5}} \right).$$
(27)

The typical simulation results were sampled from  $E(O_t|O_{t-1})$  while differing only in initial values  $O_0 \in \{-0.8, 0, 3\}$ , are in Fig. 1. For comparison, the same mixture was simulated as the conditional one, i.e. with fixed component weights.

*Discussion:* The simulation results demonstrate the dynamic dependence of the components weight on the data realisation. This is the key feature of the truly dynamic model. It allows to model non-linear dynamic effects and unbalanced activations of respective components. The standard mixture model lacks both mentioned features of the mixture ratio model.

## **On Performance of the Standard and Ratio Mixture Models:**

Simulation Conditions: Scalar discrete observations  $O_t \in O = \{1, ..., |O|\}, |O| = 5$ , and  $2^{nd}$  order regression with  $\rho_t = (O_{t-1}, O_{t-2})$  are used. This defines data vectors (9)  $\Psi_t = (O_t, O_{t-1}, O_{t-2})$ . The *learnt* (estimated) models L and the simulated system models E are either the mixture ratio model M or the standard mixture S, both with |c| = 2,

$$E, L \in \{M, S\} = \{mixture ratio, standard mixture\}.$$
 (28)

This gives four combinations of the simulated and learnt models.

The structure of the joint pd  $J(\Psi_t | \Theta)$  defining the mixture ratio model M is

$$J(\Psi_{t}|\Theta_{M}) = \alpha_{M}J_{1}(O_{t}, O_{t-1}|\Theta_{M1})G_{1}(O_{t-2}) + (1 - \alpha_{M})J_{2}(O_{t}, O_{t-2}|\Theta_{M2})G_{2}(O_{t-1}).$$

The joint pds  $J_1(O_t, O_{t-1}|\Theta_{M1})$  and  $J_2(O_t, O_{t-2}|\Theta_{M2})$  are parameterised by their values  $\Theta_{M1}$ ,  $\Theta_{M2}$  assigned to possible pairs  $(O_t, O_{t-1})$ ,  $(O_t, O_{t-2})$ . The model parameter is  $\Theta_M = (\alpha_M, \Theta_{M1}, \Theta_{M2})$ . The parameter-free factors  $G_1(O_{t-2})$ ,  $G_2(O_{t-2})$  are uniform on O. They cancel in the mixture ratio M, which gets the form

$$\mathsf{M}(O_{t}|\rho_{t},\Theta_{\mathsf{M}}) = \frac{\alpha_{\mathsf{M}}\mathsf{J}_{1}(O_{t},O_{t-1}|\Theta_{\mathsf{M}1}) + (1-\alpha_{\mathsf{M}})\mathsf{J}_{2}(O_{t},O_{t-2}|\Theta_{\mathsf{M}2})}{\sum_{O_{t}\in O}\alpha_{\mathsf{M}}\mathsf{J}_{1}(O_{t},O_{t-1}|\Theta_{\mathsf{M}1}) + (1-\alpha_{\mathsf{M}})\mathsf{J}_{2}(O_{t},O_{t-2}|\Theta_{\mathsf{M}2})}.$$

The standard model  $S(O_t | \rho_t, \Theta_S) = \alpha_S S_1(O_t | O_{t-1}, \Theta_{S1}) + (1 - \alpha_S) S_2(O_t | O_{t-2}, \Theta_{S2})$  is parameterised by  $\Theta_S = (\alpha_S, \Theta_{S1}, \Theta_{S2})$ . Its components  $S_1(O_t | O_{t-1}, \Theta_{S1})$  and  $S_2(O_t | O_{t-2}, \Theta_{S2})$  are parameterised by probabilities  $\Theta_{S1}, \Theta_{S2}$  assigned to  $O_t$  when conditioned on  $O_{t-1}$  and  $O_{t-2}$ , respectively.

*MC Study over* E: It consists of 200 runs for  $t \in t$ , |t| = 500, with randomly generated parameters  $\Theta_E$  for both considered simulated model structures  $E \in \{M, S\}$  (28). In both cases, the mixture ratio M model and the standard mixture S model are learnt on the same observation sequences.

*MC Study with a Fixed*  $E = \underline{M}$ : The fixed system being the mixture ratio model  $\underline{M}$  given by the fixed parameter  $\underline{\Theta}_{M}$  having truly dynamic weights  $w_c(\rho_t)$  (14) is simulated. Its 200 runs differ in the noise realisations.

*Performance Evaluation:* At the time  $t \in t$ , the simulation and estimation dealt with the underlying model E,  $L \in \{M, S\}$  and assigned the following probabilities to observations

$$E_{t}(O_{t}) = \begin{cases} \underline{\mathsf{M}}_{t}(O_{t}) = \mathsf{M}(O_{t}|\rho_{t},\underline{\Theta}_{\mathsf{M}}) & \text{if } \mathsf{E} = \mathsf{M} \text{ with a fixed } \underline{\Theta}_{\mathsf{M}} \\ \underline{\mathsf{S}}_{t}(O_{t}) = \mathsf{S}(O_{t}|\rho_{t},\underline{\Theta}_{\mathsf{S}}) & \text{if } \mathsf{E} = \mathsf{S} & \text{with a fixed } \underline{\Theta}_{\mathsf{S}} \\ \mathsf{L}_{t}(O_{t}) = \mathsf{L}_{t}(O_{t}|O_{t-1},\ldots,O_{0}) = \int_{\Theta_{\mathsf{L}}} \mathsf{L}(O_{t}|\rho_{t},\Theta_{\mathsf{L}})\mathsf{P}(\Theta_{\mathsf{L}}|O_{t-1},\ldots,O_{0})\mathsf{d}\Theta_{\mathsf{L}}. \end{cases}$$
(29)

The corresponding Kullback-Leibler divergences were sequentially evaluated

$$D_{t}(\mathsf{E}||\mathsf{L}) = \sum_{\tau=1}^{t} \mathsf{D}(\mathsf{E}_{\tau}||\mathsf{L}_{\tau}), \quad \text{with} \quad \mathsf{D}(\mathsf{E}_{\tau}||\mathsf{L}_{\tau}) = \sum_{O_{\tau}\in O} \mathsf{E}_{\tau}(O_{\tau}) \ln\left(\frac{\mathsf{E}_{\tau}(O_{\tau})}{\mathsf{L}_{\tau}(O_{t})}\right)$$
$$\mathsf{D}(\mathsf{E}||\mathsf{L}) \equiv \mathsf{D}_{|t|}(\mathsf{E}||\mathsf{L}), \quad \text{and their differences} \quad \Delta = \mathsf{D}(\mathsf{E}||\mathsf{M}) - \mathsf{D}(\mathsf{E}||\mathsf{S}). \tag{30}$$

They quantified estimation quality within the comparable observation space.

*Results:* The simulation results are in Fig. 2. The left hand column contains histogram values  $D(E||L) = D_{|t|}(E||L), L \in \{M, S\}$ , at final time |t| = 500. The right hand column shows the corresponding sample cumulative distributions. The better outcomes have them shifted to the left. Fig. 3 shows the course of  $D_t(E||L), t \in t, L \in \{M, S\}$  for: (a)  $E = \underline{M}$ , given by a fixed  $\underline{\Theta}_{M}$ ; (b)  $E = \underline{S}$ , given by a fixed  $\underline{\Theta}_{S}$ ; (c)  $E = \underline{M}$  with  $\underline{\Theta}_{M}$  giving the dynamic weights  $w_c(\rho_t)$  (14). Tab. 1 shows sample statistics of increments  $\Delta$  (30). Negative values of sample means and medians reflect the dominance of the mixture ratio model.

*Discussion:* Tab. 1 and Figs. 2, 3 confirm the expectation that the mixture ratio M outperforms the standard mixture if the mixture ratio M is simulated. The improvement can be significant. The learnt mixture ratio copies the performance of the standard mixture S when S is simulated. It is even slightly better as it copes better with errors caused by the approximate estimation (8).

**TABLE 1** Sample statistics of differences of Kullback-Leibler divergences  $\Delta = D(E||M) - D(E||S)$  (30), where E is the predictor corresponding to the simulated system, to the learnt predictor with M, which corresponds to the mixture ratio, and S corresponding to the standard mixture. The 1st column shows statistics for the simulated system E = M with  $\Theta_M$  varied in 200 MC runs; the 2nd column shows results for the simulated system E = S with  $\Theta_S$  varied in 200 MC runs; the 3rd column reflects the case with the simulated E = M with a fixed  $\Theta_M$  causing the truly dynamic weights  $w_c(\rho_t)$  (14) (MC runs differ in noise).

Simulated Case	MC over $E = M$	MC over $E = S$	Fixed $E = \underline{M}$
Mean	-0.7001	-0.5729	-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

## 6 | POTENTIAL OF THE MIXTURE RATIO MODEL

The paper brings an important message that the ratios of finite mixtures model well dynamically related data. The next list explicates our gains.

- ✓ The mixture ratios universally approximate non-linear *dynamic* stochastic relations.
- ✓ The mixture ratios handle cases in which rare visits of active components are significant, cf. Fig. 1. This is vital in fault detection<sup>55</sup>, detection of non-standard fraud behaviours needed<sup>56</sup>, or in cyber-security applications<sup>10</sup>; everywhere, where outlier detection is faced<sup>57</sup>. These cases are hard dynamic versions of estimation with unbalanced data<sup>58</sup>.
- The mixture ratios suit to modelling of mixed (discrete and continuous valued) data. The components are the joint pds factorable into pds of continuous valued data, possibly conditioned on discrete ones, and pds relating discrete data.
- ✓ The mixture ratios may serve as a relatively universal dynamic feature extractor. Indeed, the approximately sufficient statistic collected during estimation, see Sec. 4, are the relevant features. It suffices to use Bayesian structure estimation on mixture ratios with low-dimensional components.
- ✓ The estimation of the mixture ratios provides the stationary joint pd of the data vector (13), which can directly serve for the design of adaptive decision strategies for an infinite decision horizon<sup>59</sup>. The estimated joint pd of data vector can be appropriately factorised and the current model of the stationary decision rule replaced by the rule, which makes the joint pd close to a desired stationary joint pd<sup>60</sup>.

In future, it is desirable:

- ✓ to elaborate numerically robust (factorised) procedures for normal models, sparse Markov-chains and mixed cases in the way mimic to<sup>50</sup>;<sup>35</sup> has made definite steps in this respect;
- $\checkmark$  to refine the data-dependent choice of forgetting factors;
- ✓ to tailor Bayesian structure estimation, ready for mixtures both with respect to suitable regression vectors  $\rho_c$  and the number of mixture components<sup>51</sup> to the mixture-ratio model;
- ✓ to perform real-life tests confirming that the theoretically higher approximating strength of the discussed mixture ratios, see Sec. 3, is mostly undiminished by the approximate estimation; the experiments with commodity futures trading, made in<sup>35</sup>, as well as other (yet unpublished) real-data processing, are more than promising.

The theoretical insight, experience from simulations, preliminary processing of real data, width and importance of the outlined contributions to estimation dynamic data relations make this effort worthwhile.

Data Availability The paper deals just with simulations. Authors will provide the experimental code upon a request.

*Conflict of interest* The authors declare that this is their original work and represents no potential conflict of interests. *Acknowledgement* This research was supported by MŠMT ČR LTC18075 and by EU-COST Action CA1622.

# References

- 1. Mitchell T. Machine Learning. McGraw Hill . 1997.
- Sadghi P, Kennedy R, Rapajic P, Shams R. Finite-State Markov Modeling of Fading Channels. *IEEE Signal Processing Magazine* 2008; 57. doi: 10.1109/MSP.2008.926683
- Guyon I, Saffari A, Dror G, Cawley G. Model Selection: Beyond the Bayesian/Frequentist Divide. *Journal of Machine Learning Research* 2010; 11: 61–87.
- 4. Bishop C. Pattern Recognition and Machine Learning. Springer . 2006.
- Mason W, Vaughan J, Wallach H. Special Issue: Computational social science and social computing. *Machine Learning* 2014; 96: 257–469.
- 6. Wiering M, Otterlo vM., eds. Reinforcement Learning: State-of-the-Art. Springer . 2012.
- Guan P, Raginsky M, Willett R. Online Markov Decision Processes with Kullback Leibler Control Cost. *IEEE Trans. on* AC 2014; 59(6): 1423–1438.
- 8. Holmes G, Liu T., eds., Proc. of 7th Asian Conf. on Machine Learning (ACML2015), JMLR Workshop and Conf. Proceedings. 45; 2015.
- 9. Tsai C, Lai C, Chiang M, Yang L. Data Mining for Internet of Things: A Survey. *IEEE Communications Surveys & Tutorials* 2014; 16(1): 77-95.
- 10. Buczak A, Guven E. A Survey of Data Mining and Machine Learning Methods for Cyber Security Intrusion Detection. *IEEE Communications Surveys Tutorials* 2016; 18(2): 1153-1176.
- 11. Nguyen H, Woon Y, Ng W. A survey on data stream clustering and classification. *Knowledge and Information Systems* 2015; 45: 535–569.
- Mirsky Y, Shapira B, Rokach L, Elovici Y. pcStream: A Stream Clustering Algorithm for Dynamically Detecting and Managing Temporal Contexts. In: Cao T, et al., eds. *Advances in Knowledge Discovery and Data Mining*. 9078. LNCS Springer, Cham. 2015 (pp. 230–237).
- 13. Appice A, A., Ciampi, Malerba D. Summarizing numeric spatial data streams by trend cluster discovery. *Data Mining and Knowledge Discovery* 2015; 29(1): 84–136.
- Silva J, Faria E, Barros R, Hruschka E, Carvalho A, Gama J. Data Stream Clustering: A Survey. ACM Computing Surveys 2014; 46. doi: 10.1145/2522968.2522981
- 15. Haykin S. Neural Networks: A Comprehensive Foundation. N.Y.: Macmillan . 1994.
- Kolmogorov A. On the representation of continuous functions of several variables by superposition of continuous functions of one variable and addition. *Doklady Akademii Nauk* 1957; 114: 953-956.
- Ferguson T. Bayesian density estimation by mixtures of normal distributions. In: Rizvi M, Rustagi J, Siegmund D., eds. *Recent Advances in Statistics*Ac. Press. 1983 (pp. 287âĂŞ-302).
- 18. McLachlan G, Peel D. Finite Mixture Models. Wiley Series in Probab. & Stat.N.Y.: Wiley . 2000.
- 19. Yarotsky D. Universal approximations of invariant maps by neural networks. arXiv:1804.10306v1 [cs.NE] 26 Apr (2018).
- 20. Wald A. Statistical Decision Functions. N.Y., London: J. Wiley . 1950.

- 21. Savage L. Foundations of Statistics. Wiley . 1954.
- 22. Rao M. Measure Theory and Integration. J. Wiley . 1987.
- 23. Mosca E. Optimal, Predictive, and Adaptive Control. Prentice Hall . 1994.
- 24. Jazwinski A. Stochastic Processes and Filtering Theory. Ac. Press . 1970.
- 25. Elliot R, Assoun L, Moore J. Hidden Markov Models. N.Y.: Springer-Verlag . 1995.
- 26. Berger J. Statistical Decision Theory and Bayesian Analysis. Springer . 1985.
- 27. Berec L, Kárný M. Identification of reality in Bayesian context. In: Warwick K, Kárný M., eds. *Computer-Intensive Methods in Control and Signal Processing*Birkhäuser. 1997 (pp. 181–193).
- 28. Park J, Sandberg I. Universal Approximation using Radial-Basis-Function Networks. Neural Comput. 1991; 3: 246-257.
- 29. McNicholas P. Mixture model-based classification. Boca Raton, London, N.Y.: CRC Press. 2017.
- 30. Catania L. Dynamic Adaptive Mixture Models. 2016. arXiv:1603.01308v1.
- Nagy I, Suzdaleva E, Kárný M, Mlynářová T. Bayesian estimation of dynamic finite mixtures. Int. J. Adapt Control Signal Process. 2011; 25(9): 765-787.
- 32. Frigessi A, Haug O, Rue H. A Dynamic Mixture Model for Unsupervised Tail Estimation without Threshold Selection. *Extremes* 2002; 5(3): 219 235.
- 33. Barndorff-Nielsen O. Information and Exponential Families in Statistical Theory. N.Y.: Wiley . 1978.
- 34. Bohlin T. Interactive System Identification: Prospects and Pitfalls. Springer . 1991.
- 35. Ruman M. Mixture Ratios for Decision Making. Master's thesis. FJFI, Czech Technical University. Prague: 2018.
- 36. Ruman M, Kárný M. Dynamic Mixture Ratio Model. In: Proc. ICCAIRO, IEEE; 2019: 92 99
- 37. Kárný M. Approximate Bayesian Recursive Estimation. Inf. Sci. 2014; 289: 100-111.
- 38. Peterka V. Bayesian system identification. In: Eykhoff P., ed. Trends and Progress in System IdentificationPerg. Press. 1981.
- 39. Bernardo J. Expected Information as Expected Utility. The An. of Stat. 1979; 7: 686-690.
- Kárný M, Guy T. On Support of Imperfect Bayesian Participants. In: Guy T, et al., eds. *Decision Making with Imperfect Decision Makers*. 28. Springer, Int. Syst. Ref. Lib. 2012 (pp. 29–56).
- 41. Kerridge D. Inaccuracy and inference. J. of the Royal Stat. Soc. 1961; B 23: 284–294.
- 42. Kulhavý R. A Bayes-closed approximation of recursive nonlinear estimation. *Int. J. Adapt Control Signal Process*. 1990; 4: 271–285.
- 43. Kárný M. Minimum Expected Relative Entropy Principle. In: Proc. of the 18th ECC, IFAC; 2020; Sankt Petersburg: 35-40.
- 44. Kullback S, Leibler R. On information and sufficiency. Ann Math Stat 1951; 22: 79-87.
- 45. Shore J, Johnson R. Axiomatic derivation of the principle of maximum entropy & the principle of minimum cross-entropy. *IEEE Tran. on Inf. Th.* 1980; 26(1): 26–37.
- 46. Šmídl V, Quinn A. The Variational Bayes Method in Signal Processing. Springer . 2005.
- 47. Minka T. A family of algorithms for approximate Bayesian inference. PhD thesis. MIT, 2001.
- 48. Koopman R. On distributions admitting a sufficient statistic. Trans. of Am. Math. Society 1936; 39: 399.
- 49. Tipping M, Bishop C. Probabilistic principal component analysis. J. of the Royal Society Series B 1999; 61: 611–622.

- 50. Kárný M. Recursive estimation of high-order Markov chains: Approx. by finite mixtures. Inf. Sci. 2016; 326: 188-201.
- 51. Kárný M, Böhm J, Guy T, et al. Optimized Bayesian Dynamic Advising: Theory and Algorithms. Springer . 2006.
- 52. Nelsen R. An Introduction to Copulas. N.Y.: Springer . 1999.
- 53. Dedecius K, Nagy I, Kárný M, Pavelková L. Parameter Estimation With Partial Forgetting Method. In: *Proc. of the 15th IFAC SYSID*, 2009.
- 54. Carter C. Commodity future markets: A survey. The Australian J. Agricultural and Resource Economics 1999; 43: 209-247.
- 55. Polycarpou M, Helmicki A. Automated fault detection and accommodation: A learning systems approach. *IEEE Trans. on Systems, Man, and Cybernetics* 1995; 25(11): 1447-1458.
- 56. Kou Y, Lu C, Sirwongwattana S, Huang Y. Survey of fraud detection techniques. In: *IEEE Intern. Conf. on Networking, Sensing and Control*, 2; 2004: 749-754
- 57. Hodge V, Austin J. A Survey of Outlier Detection Methodologies. Artificial Intelligence Review 2004; 22(2): 85-126.
- Bekkar M, Alitouche T. Imbalanced Data Learning Approaches Review. Int. J. of Data Mining & Knowledge Management Process 2013; 3(4): 15–33.
- 59. Kushner H. Introduction to Stochastic Control. Holt, Rinehart and Winston . 1971.
- 60. Kárný M, Kroupa T. Axiomatisation of fully probabilistic design. Inf. Sci. 2012; 186(1): 105-113.

# **AUTHORS' BIOGRAPHIES**



**Miroslav Kárný**. Ing. (MSc) Czech Technical University (CTU) Prague, 1973; CSc. (PhD) 1978, DrSc. (DSc) 1990, both the Institute of Information Theory and Automation, the Czechoslovak Academy of Sciences (UTIA CAS) employing him since 1973 in the Department of Adaptive Systems (AS). *Research:* conceptual, theoretical and algorithmic aspects of AS based on Bayesian dynamic decision making (DM) and its original fully probabilistic extension. *Teaching:* the advanced course on DM, CTU since 1991; supervision of 13 defended PhD students (+6 co-supervision) and numerous Bc, MSc theses and research projects. *Publications:* 1 monograph, 6 edited books,  $\approx 400$  works ( $\approx 10$  chapters,  $\approx 100$  articles); for the list after

1989 see http://www.utia.cz/people/karny.



**Marko Ruman**. Ing. (MSc) mathematical engineering at Faculty of Nuclear Sciences and Physical Engineering (FNSPE CTU) Prague, 2018; PhD Student 2018 – at FNSPE CTU, supervisor Dipl. Eng T.V. Guy, PhD. Since 2018 research assistant at AS UTIA CAS. *Research:* adaptive decision making, Markov decision processes, deep reinforcement learning, transferring knowledge between reinforcement learning tasks. *Theoretical skills:* advanced statistical and numerical methods, advanced functional analysis, dynamic decision making, Monte Carlo methods, classical physics, relational data bases. *Programming skills:* since 2006 active design of web-sites using HTML, CSS, PHP, MySQL and JavaSript, see e.g. www.ladywander.com;

also an active user of Matlab, C++, Python. *Languages:* fluent English, Slovak, basic Italian, French and German. *Publications:* co-author of 4 publications, see https://www.utia.cas.cz/node/3090.

How to cite this article: Kárný M., Ruman M., (2020), Mixture Ratio Modelling of Dynamic Systems, *Int J Adapt Control Signal Processing*, 2021; @@:@@-@@.



**FIGURE 1** Simulation of the system E (27). Its joint pd is shown in the middle. Trajectories in respective columns differ in the initial value of the regression vector  $\rho_1 = O_0 \in \{-0.8, 0, 3\}$ , respectively, while pseudo-random noise realisation is the same in all cases. The 1st row of figures shows the realised observations  $O_t$ ,  $t \in t \equiv \{1, ..., 300\}$ . It demonstrates the strong dependence of the observed trajectories on the initial condition  $\rho_1 = O_0$ . The 2nd row shows realised observations when the same components as in (27) are used as the conditional pds, i.e. when the standard mixture with constant component weights is simulated. The respective trajectories are almost identical. It has allowed us to replace the case corresponding  $\rho_1 = O_0 = 0$  by the simulated joint pd (27). The 3rd row shows evolution of dynamic weights  $(w_{t;1}(\rho_t))_{t \in t}$  (14) corresponding to the realised observations shown in the 1st row of figures. Again, the strong dependence on the initial  $\rho_1 = O_0$  is clearly seen.



FIGURE 2 Evaluation of Kullback-Leibler divergences D(E||L) of the observation predictors corresponding to the simulated systems E to the learnt predictors, see (30), in the Monte Carlo study. The left column shows histograms, the right one shows the empirical distribution functions on which the favourable shift to smaller divergence values is better seen. Colors distinguish the inspected learnt predictor  $L \in \{M, S\} = \{$ mixture ratio, standard mixture $\} = \{$ blue, red $\}$ . The 1st row deals with the simulated system E = M = mixture ratio, with its parameter  $\Theta_M$  varied in 200 MC runs; the 2nd row shows the results with the simulated system E = S = standard mixture, with its parameter  $\Theta_S$  varied in 200 MC runs; the 3rd row shows results for  $E = \underline{M}$  = mixture ratio, with a fixed parameter  $\underline{\Theta}_{M}$  causing truly dynamic evaluation of the weights  $w_{c}(\rho_{t})$  (14), there the MC study runs over 200 different noise realisations.

45

40

45

90

100



**FIGURE 3** Time evolutions of the Kullback-Leibler divergences  $D_t(E||L)$  (30) od predictors corresponding to the simulated system E and learnt predictors  $L \in \{M, S\} = \{$ full blue, dashed red $\} = ($ mixture ratio, standard mixture). The respective figures from left to right correspond to:  $E = \underline{M} =$ simulated mixture ratio with a randomly chosen  $\underline{\Theta}_M$ ;  $E = \underline{S} = mbox simulated standard mixture$  with a randomly chosen  $\underline{\Theta}_S$ ;  $E = \underline{M} =$ simulated mixture ratio with with the fixed  $\underline{\Theta}_M$  causing a truly dynamic weighting  $w_c(\rho_t)$  of respective components (14).