

# Mixture Ratio Modelling of Dynamic Relations

Miroslav Kárný · Marko Ruman

Received: date / Accepted: date

**Abstract** Any knowledge extraction from data sets implicitly or explicitly relies on hypothesis about relations between data records. The inspected black-box methodology uses “universally” approximating mapping describing a wide class of relations. The knowledge incompleteness, the ever present uncertainty and the ultimate use in a subsequent decision making (DM) single out probabilistic models as adequate. Among them, the use of finite mixtures with components in the exponential family (EF) dominate. Their dominance stems from their flexibility, cluster interpretability and the availability of a range sophisticated algorithms suitable for high-dimensional data streams. They are even used in dynamic situations with mutually dependent data records. The dependence is modelled by employing regression and auto-regression components. However, with a few specialised exceptions, these dynamic models assume data-independent weights of mixture components. Their use is illogical as it assumes independent (memory-less) transitions between dynamic mixture components. The restricted nature of such mixtures follows from the fact that the *set of finite probabilistic mixtures is not closed with respect to the conditioning*, which is the key learning operation. The paper overcomes this restriction by using *ratios of finite mixtures* as dynamic parametric models. It motivates them, elaborates their approximate Bayesian *stream* learning. The paper reveals their application potential stemming predominantly from their inherent ability to process unbalanced data sets. This is vital when a the portion of interesting data, reflecting for instance failures or attacks, is small.

**Keywords** Mixture Models · Approximate Bayesian Learning · Stream Learning · Kullback-Leibler Divergence

## 1 Introduction

The paper topic falls to a broad domain of decision making (DM) understood as a targeted choice among available options. This covers many domains including machine learning, Mitchell (1997), signal processing, Sadghi et al. (2008), filtering,

---

The Czech Academy of Sciences, Institute of Information Theory and Automation, POB 18, 182 08 Prague 8, Czech Republic, E-mail: school@utia.cas.cz, E-mail: marko.ruman@gmail.com

Jazwinski (1970), hypothesis testing, Guyon et al. (2010), classification and pattern recognition, Bishop (2006), knowledge sharing, Mason et al. (2014), reinforcement learning, Wiering and van Otterlo (2012), control, Guan et al. (2012), etc.

All these areas need to enrich the underlying domain knowledge by that hidden in, often excessive, amount of data records: they need data mining. Unsurprisingly, the amount of relevant theoretical, algorithmic and application results is excessive, see e.g. in Holmes and Liu (2015). This makes us to provide reference sample.

Nowadays, data mining is an extremely broad area with many results and applications. Even surveys specialise to Internet of Things, Tsai et al. (2014), cyber-security, Buczak and Guven (2016), etc.

*Paper Focus:* The paper contributes to data mining focused on clustering of data streams, Nguyen et al. (2015), with dynamically related data records, Mirsky et al. (2015). The contribution concerns modelling and conditioning. Importance of this data-mining type is seen on specific applications, e.g. in geophysics, Appice et al. (2015), or stream clustering of independent data records, Silva et al. (2014).

*DM Perspective on Knowledge Processing:* A solution of a DM problem leads to a strategy, a collection of decision rules mapping the evolving available knowledge on actions, Wald (1950). The strategy should meet DM aims in the best way under the faced circumstances. The adopted Bayesian DM, Savage (1954), proved to be an adequate methodology coping with incomplete knowledge, uncertainty and randomness of the environment to which actions relate. Bayesian paradigm links DM consequences to the acquired knowledge and the used actions by conditional distributions, here described by conditional probability densities (pd<sup>1</sup>).

Actions are chosen sequentially (in a single pass). This enables DM to enrich its knowledge, to learn a better environment model. This opportunity is exploited by adaptive controllers, Åström and Wittenmark (1989). The same processing of data streams serves as a statistic (knowledge, feature) extractor. The learning is possible iff the learnt relations practically do not change during the knowledge accumulation. The most general used structure of such relations leads to stochastic filtering, Jazwinski (1970). It is mostly applied to hidden Markov models, Elliot et al. (1995).

*The Addressed Modelling Problem:* The achievable modelling and thus DM quality are predetermined by the inspected model set. One-pass Bayesian learning of general hidden Markov models is possible but hard. The treatment of parametric models we focus on can be conceptually extended to them but it is out of scope of this paper. The considered relations of the observed data are described by parametric probabilistic models. Which parameter corresponds to the best model is a priori unknown. Bayesian learning offers the unambiguous deductive Bayes' rule, Berger (1985), which redistributes the belief (pd) about model adequacy, Berc and Kárný (1997). It accumulates the knowledge into the posterior pd.

As usual in the targeted DM, we rely on black-box models. They approximate joint pd of long sequences of observed data records by a stationary Markov model operating on a fixed-dimensional state. Focus on them is enabled by a common

<sup>1</sup> Pds are evaluated with respect to a dominating, typically Lebesgue or counting, product measure, Rao (1987). Lebesgue's notation  $\int \dots d\bullet$  is mostly used.

weak dependence of data realised at distant time moments and by the possibility to include time into the model state. Then, a state transition pd fully determines the needed joint pd for any given initial state.

For a fixed condition, any continuous transition pd can be arbitrarily-well approximated by a finite mixture of conditional pds (components). Components differ in their parameters but may have the same functional form, say Gaussian one. Various versions of this “universal approximation property” have been proved in neural-network context, Park and Sandberg (1991), as well as in probabilistic modelling, McLachlan and Peel (2000). Such expansions perfectly work and are widely used across domains for zero-memory models, McNicholas (2017).

By construction, the component weights should depend on the condition. The common use of the condition-independent component weights goes against the expansion logic. Rare ad hoc exceptions, Catania (2016); Frigessi et al. (2002), indicate how much is lost when considering constant component weights.

This paper goes beyond the common practice and offers *ratios* of finite mixtures with components in the exponential family (EF), Barndorff-Nielsen (1978), as black-box, Bohlin (1991), universally approximating models, Haykin (1994). It develops, inevitably approximate, but feasible sequential Bayesian learning of models from this *extremely rich but yet unconsidered model set*.

*Layout and Common Notation.* The paper relies on an approximate sequential Bayesian learning, Kárný (2014). Its recall in Sec. 2 prepares the used notation. Sec. 3 justifies the mixture ratio models in detail. Sec. 4 elaborates their learning. Sec. 5 numerically illustrates the theory. Sec. 6 outlines the model applicability.

Concerning notation,  $\mathbf{X}$  is a set of  $X$ s. It is either a non-empty subset of a finite-dimensional real space or a subset of pds.  $|\mathbf{X}|$  means cardinality of  $\mathbf{X}$ . Random variables, their realisations and function arguments are undistinguished. The context suffices to grasp the correct meaning. **San serif fonts** mark mappings, which are taken as different if their arguments differ, but mnemonic symbols are preferred. For instance,  $M$  denotes a parametric model,  $J$  is a joint parametric pd of data vectors,  $P$  is posterior pd,  $O$  marks observations,  $A$  labels actions, etc. Decoration  $\sim$  denotes intermediate objects and  $\hat{\phantom{x}}$  estimates.

## 2 Approximate Sequential Bayesian Learning

This section summarises the used approximate sequential learning, Kárný (2014), and introduces basic treated objects.

Parametric model  $M$  links observations  $O_t \in \mathcal{O}$  to actions  $A_t \in \mathcal{A}$  at discrete-time moments labelled by  $t \in \mathcal{t} = \{1, 2, \dots, |t|\}$ . Data records  $D_t = (O_t, A_t)$  are sequentially processed. Parametric model  $M$  is pd

$$M(O_t|A_t, D^{t-1}, \Theta), \quad D^{t-1} = (D_{t-1}, \dots, D_1, D_0) \quad (1)$$

that relates observation  $O_t$  to current action  $A_t$  and past data records  $D^{t-1}$ , extended by prior knowledge  $D_0$ . An unknown  $\Theta \in \Theta$  parameterises this pd.

The knowledge about parameter  $\Theta$  at time  $t-1$  is expressed by posterior pd  $P_{t-1}(\Theta) = P(\Theta|A_t, D^{t-1}) = P(\Theta|D^{t-1})$ . The last equality is implied by natural

conditions of control, Peterka (1981), see also Assumption 1 below. Bayes' rule updates this knowledge by the realisation of data record  $D_t$  to

$$\tilde{P}_t(\theta) = \frac{M(O_t|A_t, D^{t-1}, \theta)P_{t-1}(\theta)}{\int_{\Theta} M(O_t|A_t, D^{t-1}, \theta)P_{t-1}(\theta)d\theta} \propto M(O_t|A_t, D^{t-1}, \theta)P_{t-1}(\theta), \quad (2)$$

where  $\propto$  is proportionality. With 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>2</sup> pds, has to be evolved.

Generally, posterior pd  $\tilde{P}_t \notin \mathbf{P}$  even if pd  $P_{t-1} \in \mathbf{P}$ . To preserve feasibility, pd  $\tilde{P}_t$  in (2) is to be projected on  $\mathbf{P}$ . Work of Bernardo (1979) inspecting this approximation task provides weak conditions (met here) under which the minimiser  $\hat{P}_t \in \mathbf{P}$  of Kerridge's inaccuracy, Kerridge (1961),

$$K(\tilde{P}_t || \hat{P}_t) = - \int_{\Theta} \tilde{P}_t(\theta) \ln(\hat{P}_t(\theta)) d\theta \quad (3)$$

is the proper Bayesian projection of  $\tilde{P}_t$  on  $\hat{P}_t \in \mathbf{P}$ . The projection *should not* serve as the prior pd in a further updating: this may cause a divergence of the sequentially projected pds from those projected optimally in the batch mode, Kulhavý (1990). Forgetting, with a data-dependent factor  $\lambda_t \in [0, 1]$ , is the adequate countermeasure, Kárný (2014). It completes updating of feasible posterior pds

$$\begin{aligned} \tilde{P}_t(\theta) &\propto M(O_t|A_t, D^{t-1}, \theta)P_{t-1}(\theta), \quad \hat{P}_t = \arg \min_{\hat{P} \in \mathbf{P}} K(\tilde{P}_t || \hat{P}) \\ P_t(\theta) &\propto \hat{P}_t^{\lambda_t}(\theta)P_{t-1}^{1-\lambda_t}(\theta), \quad \lambda_t = \frac{[\int_{\Theta} M(O_t|A_t, D^{t-1}, \theta)P_{t-1}(\theta)d\theta]^2}{\int_{\Theta} M^2(O_t|A_t, D^{t-1}, \theta)P_{t-1}(\theta)d\theta}. \end{aligned} \quad (4)$$

The above plausible choice of forgetting factor  $\lambda_t$  is in Kárný et al. (2014).

The next assumption summarises conditions under which (4) is applicable. It adds reasons for them and serves for referencing.

**Assumption 1 (Conditions for Applicability of (4))**

- ✓ Parameter  $\theta \in \Theta$  is unknown to the employed decision rules, described by pds  $R(A_t|D^{t-1}, \theta)$ . Thus, they meet natural conditions of control, Peterka (1981),

$$R(A_t|D^{t-1}, \theta) = R(A_t|D^{t-1}) \Leftrightarrow P(\theta|A_t, D^{t-1}) = P(\theta|D^{t-1}). \quad (5)$$

Only such decision rules cancel in (2).

- ✓ Parametric model  $M(O_t|A_t, D^{t-1}, \theta)$  (1) has a  $\theta$ -independent support as justification of forgetting (4) needs this assumption, Kárný (2014).
- ✓ The set of feasible pds  $\mathbf{P}$  is log-convex, i.e. forgetting (4) applied to its members gives a pd in  $\mathbf{P}$ . Without this, an additional projection is needed.

**Remarks 1**

- ✓ Kerridge's inaccuracy  $K$  (3) is a shifted version of Kullback-Leibler divergence  $D(\tilde{P} || \hat{P})$ , Kullback and Leibler (1951)

$$D(\tilde{P} || \hat{P}) = \int_{\Theta} \tilde{P}(\theta) \ln\left(\frac{\tilde{P}(\theta)}{\hat{P}(\theta)}\right) d\theta = \int_{\Theta} \tilde{P}(\theta) \ln(\tilde{P}(\theta)) d\theta + K(\tilde{P} || \hat{P}).$$

It has the same minimiser but copes better with  $\tilde{P}$  containing Dirac's pds.

- ✓ The unavoidable projection of  $\tilde{P}_t$  on  $\mathbf{P}$  (4) is demanding but feasible, cf. Sec. 3. The extra effort for evaluating forgetting factor  $\lambda_t$  is small.

---

<sup>2</sup> An intuitive understanding of this notion suffices. Sec. 3 provides example of such pds.

### 3 Ratio of Finite Mixtures

This core section provides the proposed universal parametrisation of pds modelling data observed in the closed loop formed by a dynamic environment and a, possibly randomised and dynamic, decision strategy meeting (5).

*Markovian Modelling:* Joint parametric pd  $\tilde{J}(D^{|\mathbf{t}|}|\Theta, D_0)$  of data sequences  $D^{|\mathbf{t}|} = (D_t)_{t \in \mathbf{t}}$  factorises

$$\begin{aligned} \tilde{J}(D^{|\mathbf{t}|}|\Theta, D_0) &= \prod_{t \in \mathbf{t}} \overbrace{M(O_t|A_t, D^{t-1}, \Theta)}^{\text{parametric model}} \overbrace{R(A_t|D^{t-1}, \Theta)}^{\text{decision rule}} \\ &\stackrel{(5)}{=} \prod_{t \in \mathbf{t}} M(O_t|A_t, D^{t-1}, \Theta) R(A_t|D^{t-1}), \end{aligned}$$

i.e. the parametrisation and learning concern the parametric model.

**Assumption 2 (Markov Time-Invariant Parametric Model)** *The model is time-invariant, parameterised by a constant multivariate parameter  $\Theta \in \Theta$ . It is Markov model of an order  $n < \infty$ , which means*

$$\begin{aligned} M(O_t|A_t, D^{t-1}, \Theta) &= M(O_t|\psi_t, \Theta), \text{ with regression vector} \\ \psi_t &= \text{function of } \psi_{t-1} \ \& \ \begin{cases} A_t, D_{t-1}, \dots, D_{t-n} & \text{if } 1 \leq n < \infty \\ A_t \text{ or void} & \text{for } n = 0. \end{cases} \end{aligned} \quad (6)$$

Data entering model  $M$  at time  $t \in \mathbf{t}$  form data vector  $\Psi_t = (O_t, \psi_t) \in \Psi$ .

The prior knowledge  $D_0$  determines the model structure and also  $\psi_0$ .

Parametric model (6) is ratio of the joint pd  $J(\Psi|\Theta)$  and its marginal

$$M(O_t|\psi_t, \Theta) = \frac{J(O_t, \psi_t|\Theta)}{\int_{\mathcal{O}} J(O_t, \psi_t|\Theta) dO_t} = \frac{J(\Psi_t|\Theta)}{\int_{\mathcal{O}} J(O_t, \psi_t|\Theta) dO_t}. \quad (7)$$

It is time invariant iff joint pd  $J(\Psi_t|\Theta)$  is a time-invariant function multiplied by an arbitrary positive function of  $A_t, D^{t-1}, \Theta$ , e.g., a variable decision rule. Thus, Assumption 2 practically makes joint pd  $J(\Psi_t|\Theta)$  time-invariant.

As already discussed in Introduction, the Markov property is inevitable for the targeted feasibility of the sequential learning. If not met naturally, it must be enforced via an approximation. The required time invariance can be relaxed by inclusion of time into  $\psi_t$  or by considering a time- and data- dependent unknown parameter  $\Theta$ . The latter case corresponds with hidden Markov models, Elliot et al. (1995), which are out of scope of this paper.

*Exponential Family:* The proposed model exploits members of EF. Recall that a parametric Markov model  $\tilde{M}(O_t|\psi_t, \Theta)$  belongs to EF if it is described by the pd

$$\tilde{M}(O_t|\psi_t, \Theta) = \exp \langle B(\Psi_t), C(\Theta) \rangle, \quad \Psi_t \in \Psi, \Theta \in \Theta. \quad (8)$$

There, multivariate real-valued functions  $B(\Psi_t), C(\Theta)$  have a fixed finite dimension and real-valued mapping  $\langle B(\Psi_t), C(\Theta) \rangle$  is linear in  $B(\Psi_t)$ -values. Definition (8)

admits usual factors depending respectively on  $\Psi$  and  $\Theta$ . It suffices to include constant entries into  $\mathbf{B}$  and  $\mathbf{C}$ . The support indicator is dropped for simplicity.

Under natural conditions of control (5), exponential-family members possess conjugated (self-reproducing) pds, Berger (1985),

$$P_t(\Theta) = P(\Theta|D^t) = P(\Theta|\mathbf{V}_t) = \frac{\exp \langle \mathbf{V}_t, \mathbf{C}(\Theta) \rangle}{N(\mathbf{V}_t)} \quad (9)$$

$$N(\mathbf{V}_t) = \int_{\Theta} \exp \langle \mathbf{V}_t, \mathbf{C}(\Theta) \rangle d\Theta \quad \text{with updating } \mathbf{V}_t = \mathbf{V}_{t-1} + \mathbf{B}(\Psi_t).$$

Their use converts Bayes' rule, which is a functional recursion, into algebraic updating of values the fixed-dimensional sufficient statistic  $\mathbf{V}_t$ . Initial  $\mathbf{V}_0$  in (9) describes the conjugated prior pd quantifying prior knowledge  $D_0$ . It has to guarantee  $N(\mathbf{V}_0) < \infty$ . Such  $\mathbf{V}_0$  regularises the learning.

### Remarks 2

- ✓ Under Assumption 1, exponential family exhausts parametric models, smooth in  $\Theta$ , with a finite-dimensional sufficient statistic  $\mathbf{V}_t = \mathbf{V}(D^t)$ , Koopman (1936).
- ✓ 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 factor in (7), depends on  $\Theta$  only (not on regression vector). This makes them similar to the static case with  $\Psi_t = O_t$ , for which EF is much richer.

*Universal Approximation:* Joint pds  $J(\Psi) = J(O, \psi)$  of data vectors  $\Psi \in \mathbf{\Psi}$ , even with non-trivial regression vectors  $\psi \in \mathbf{\psi}$ , can be “universally” approximated by a finite mixture of normal pds, to an arbitrary precision, McLachlan and Peel (2000). Thus, it can be approximated by a finite mixture of pds from EF. This allows us to consider the joint parametric pd in (7) as the finite weighted sum of components, pds  $(J_c(\Psi_t|\Theta_c))_{c \in \mathbf{c}}$  on  $\mathbf{\Psi}$ , in EF

$$J(\Psi|\Theta) = \sum_{c \in \mathbf{c}} \alpha_c \overbrace{\exp \langle \mathbf{B}_c(\Psi), \mathbf{C}_c(\Theta_c) \rangle}^{J_c(\Psi_t|\Theta_c)}, \quad \mathbf{c} = \{1, \dots, |\mathbf{c}|\}, \quad \alpha = (\alpha_c)_{c \in \mathbf{c}}$$

$$\alpha \in \mathbf{\alpha} = \left\{ \alpha_c \geq 0, \sum_{c \in \mathbf{c}} \alpha_c = 1 \right\}, \quad \Theta = (\alpha_c, \Theta_c)_{c \in \mathbf{c}} \in \mathbf{\Theta} = (\mathbf{\alpha}, (\Theta_c)_{c \in \mathbf{c}}). \quad (10)$$

Insertion of (10) into (7) gives the parametric *ratio model*, which is needed in learning via Bayes' rule (2),

$$M(O_t|\psi_t, \Theta) = \sum_{c \in \mathbf{c}} \frac{\alpha_c \exp \langle \mathbf{B}_c(O_t, \psi_t), \mathbf{C}_c(\Theta_c) \rangle}{\sum_{\tilde{c} \in \mathbf{c}} \alpha_{\tilde{c}} \underbrace{\int_O \exp \langle \mathbf{B}_{\tilde{c}}(O_t, \psi_t), \mathbf{C}_{\tilde{c}}(\Theta_{\tilde{c}}) \rangle dO_t}_{W_{\tilde{c}}(\psi_t, \Theta_{\tilde{c}})}} \quad (11)$$

$$= \sum_{c \in \mathbf{c}} \underbrace{\frac{\alpha_c W_c(\psi_t, \Theta_c)}{\sum_{\tilde{c} \in \mathbf{c}} \alpha_{\tilde{c}} W_{\tilde{c}}(\psi_t, \Theta_{\tilde{c}})}}_{w_c(\psi_t, \Theta)} \underbrace{\frac{\exp \langle \mathbf{B}_c(\Psi_t), \mathbf{C}_c(\Theta_c) \rangle}{W_c(\psi_t, \Theta_c)}}_{M_c(O_t|\psi_t, \Theta_c)} = \sum_{c \in \mathbf{c}} w_c(\psi_t, \Theta) M_c(O_t|\psi_t, \Theta_c).$$

Assumption 2 and *universal approximation* by finite mixtures for void  $\psi$  imply:

Ratio (11) approximates any *dynamic*, Markov, time-invariant environment model.

### Remarks 3

- ✓ Let us stress that functions  $W_c(\psi_t, \Theta_c)$  result from the integration over  $\mathbf{O}$  only, not over  $\Psi$ . This often takes components  $(M_c(O_t|\psi_t, \Theta_c))_{c \in \mathbf{c}}$  out of EF.
- ✓ The second row of (11) shows that the model is a fully dynamic finite mixture. Both components and their weights  $w_c(\psi_t, \Theta)$  depend on  $\psi_t$ . The data-dependence is not arbitrary and needs no extra parameter.
- ✓ Components  $J_c(\Psi_t|\Theta_c) = \exp \langle B_c(\Psi_t), C_c(\Theta_c) \rangle$  in (10) are pds on  $\Psi$ . These pds may contain parameter-independent factors, i.e. pds

$$J_c(\Psi_t|\Theta_c) = \exp \langle B_c(\Psi_t), C_c(\Theta_c) \rangle = \exp \langle B_c(\Psi_{t;c}), C_c(\Theta_c) \rangle G_c(\psi_{\underline{t};c}) \quad (12)$$

model  $\Psi_{t;c} = (O_t, \psi_{t;c})$ . There  $\psi_{t;c}$  is a sub-vector of  $\psi_t$  and  $G_c(\psi_{\underline{t};c})$  is a parameter-free pd on its complement  $\psi_{\underline{t};c}$  to  $\psi_t$ . Their use mimics mixtures of principal component analysers, Tipping and Bishop (1999), and diminishes the dimensionality curse, Kárný (2016).

## 4 Learning with Mixture Ratios

This section applies learning (4) to mixture ratio model (11) and arrives at its feasible and justified sequential Bayesian learning.

*Choice of Set  $\mathbf{P}$  of Feasible PDs:* The component weights in (10) define pd  $M(c_t = c|\Theta) = \alpha_c$ ,  $c \in \mathbf{c}$ , of an *unobserved* pointer,  $c_t \in \mathbf{c}$ , to the active component  $J_{c_t}(\Psi_t|\Theta_{c_t})$  (12) “generating” data vector  $\Psi_t$ . This pointer model is from EF as

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

For the *observed* pointer  $c_t \in \mathbf{c}$  to the active component, Dirichlet’s pd, determined by  $|\mathbf{c}|$ -vector statistic  $\mathbf{v}_t$ , is the corresponding conjugated pd, Berger (1985),

$$\begin{aligned} P_t(\alpha) &= P(\alpha|\mathbf{v}_t) = \frac{\exp \langle \mathbf{v}_t - 1, \ln(\alpha) \rangle}{\text{Be}(\mathbf{v}_t)} \quad \text{with update } \mathbf{v}_t = \mathbf{v}_{t-1} + \delta(c_t), \quad \mathbf{v}_0 > 0 \\ \text{Be}(\mathbf{v}) &= \frac{\prod_{c \in \mathbf{c}} \Gamma(v_c)}{\Gamma(\sum_{c \in \mathbf{c}} v_c)}, \quad \Gamma(v) = \int_0^\infty z^{v-1} \exp(-z) dz, \quad v > 0. \end{aligned} \quad (13)$$

For normal and Markov chain components or independent data vectors with components (12), conjugate pds  $P_t(\Theta_c) = P(\Theta_c|D^t, c_t, \dots, c_1)$ ,  $c \in \mathbf{c}$ , are

$$P_t(\Theta_c) = \frac{\exp \langle \mathbf{V}_{t;c}, C_c(\Theta_c) \rangle}{N_c(\mathbf{V}_{t;c})}, \quad N_c(\mathbf{V}_{t;c}) = \int_{\Theta_c} \exp \langle \mathbf{V}_{t;c}, C_c(\Theta_c) \rangle d\Theta_c, \quad c \in \mathbf{c}. \quad (14)$$

They are “natural” candidates for creating set  $\mathbf{P}$  of feasible pds. They have to be given by *statistic*  $(\mathbf{v}_t, (\mathbf{V}_{t;c})_{c \in \mathbf{c}})$  with values given by the observed data, not by unobserved pointers  $(c_t, \dots, c_1)$ . Marginal pds  $P_t(\alpha)$  (13),  $(P_t(\Theta_c))_{c \in \mathbf{c}}$  (14) of  $P_t(\Theta)$ ,  $\Theta = (\alpha, (\Theta_c)_{c \in \mathbf{c}})$ , do not determine joint pd  $P_t(\Theta)$  unambiguously, Nelsen (1999). Under the faced lack of information about mutual relations of component parameters  $(\alpha, (\Theta_c)_{c \in \mathbf{c}})$ , the product combination of marginal pds is preferable. This motivates the relatively universal choice

$$\mathbf{P} = \left\{ P_t(\Theta) = P_t(\alpha) \prod_{c \in \mathbf{c}} P_t(\Theta_c) \right\}, \quad \text{with } \mathbf{V}_t = (\mathbf{v}_t, (\mathbf{V}_{t;c})_{c \in \mathbf{c}}), \quad \text{in (13), (14).} \quad (15)$$

*Evaluation of Kerridge's Inaccuracy (3):* With the chosen  $\mathbf{P}$ , it remains to convert updating (4) into the updating  $\mathbf{v}_{t-1}, \mathbf{V}_{t-1}$  to  $\mathbf{v}_t, \mathbf{V}_t$ . Kerridge's inaccuracy of pd  $\tilde{P}_t$  (2) to pd  $\hat{P}_t \in \mathbf{P}$  (15), given by  $\hat{\mathbf{V}}_t = (\hat{\mathbf{v}}_t, (\hat{\mathbf{V}}_{t;c})_{c \in \mathbf{c}})$  (13), (14), reads

$$\begin{aligned} K(\tilde{P}_t || \hat{P}_t) &= \ln(\text{Be}(\hat{\mathbf{v}}_t)) + \sum_{c \in \mathbf{c}} \ln(N_c(\hat{\mathbf{V}}_{t;c})) \\ &\quad - \left\langle \hat{\mathbf{v}}_t - 1, \int_{\Theta} \ln(\alpha) \tilde{P}_t(\Theta) d\Theta \right\rangle - \sum_{c \in \mathbf{c}} \left\langle \hat{\mathbf{V}}_{t;c}, \int_{\Theta} C_c(\Theta_c) \tilde{P}_t(\Theta) d\Theta \right\rangle. \end{aligned} \quad (16)$$

Bayes' rule (2) for (11), (12) and pd  $P_{t-1} \in \mathbf{P}$ , (13), (14), and (15), gives  $\tilde{P}_t(\Theta) \propto$

$$\begin{aligned} &\sum_{c \in \mathbf{c}} \frac{\alpha_c \exp \langle B_c(\Psi_t), C_c(\Theta_c) \rangle \exp \langle \mathbf{v}_{t-1} - 1, \ln(\alpha) \rangle \prod_{\tilde{c} \in \mathbf{c}} \exp \langle \mathbf{V}_{t-1;\tilde{c}}, C_{\tilde{c}}(\Theta_{\tilde{c}}) \rangle}{\sum_{\tilde{c} \in \mathbf{c}} \alpha_{\tilde{c}} W_{\tilde{c}}(\psi_t, \Theta_{\tilde{c}})} \\ &\sum_{c \in \mathbf{c}} \frac{\exp \langle \mathbf{v}_{t-1} + \delta(c) - 1, \ln(\alpha) \rangle \prod_{\tilde{c} \in \mathbf{c}} \exp \langle \mathbf{V}_{t-1;\tilde{c}} + \delta(c, \tilde{c}) B_c(\Psi_{t;c}), C_{\tilde{c}}(\Theta_{\tilde{c}}) \rangle G_c(\psi_{t;c})}{\sum_{\tilde{c} \in \mathbf{c}} \alpha_{\tilde{c}} W_{\tilde{c}}(\psi_t, \Theta_{\tilde{c}})}. \end{aligned} \quad (17)$$

Next steps exploit auxiliary pds arising from summands in (17), cf. (13), (14),

$$\begin{aligned} \tilde{Q}_{t;c}(\Theta) &= \frac{\exp \langle \mathbf{v}_{t-1} + \delta(c) - 1, \ln(\alpha) \rangle}{\beta_{t;c}} \prod_{\tilde{c} \in \mathbf{c}} \exp \langle \mathbf{V}_{t-1;\tilde{c}} + \delta(c, \tilde{c}) B_c(\Psi_{t;c}), C_{\tilde{c}}(\Theta_{\tilde{c}}) \rangle \\ \beta_{t;c} &= \text{Be}(\mathbf{v}_{t-1} + \delta(c)) \prod_{\tilde{c} \in \mathbf{c}} N(\mathbf{V}_{t-1;\tilde{c}} + \delta(c, \tilde{c}) B_c(\Psi_{t;c})) / G_c(\psi_{t;c}). \end{aligned} \quad (18)$$

They independently model  $\alpha$  and  $(\Theta_c)_{c \in \mathbf{c}}$ . Dependence is brought into  $\tilde{P}_t$  via

$$\begin{aligned} H(\psi_t, \Theta) &= \frac{1}{\sum_{\tilde{c} \in \mathbf{c}} \alpha_{\tilde{c}} W_{\tilde{c}}(\psi_t, \Theta)} \quad \text{as} \quad \tilde{P}_t(\Theta) = \frac{H(\psi_t, \Theta) \sum_{c \in \mathbf{c}} \beta_{t;c} \tilde{Q}_{t;c}(\Theta)}{I_{-1}(\Psi_t, \mathbf{v}_{t-1}, \mathbf{V}_{t-1})} \\ I_{-1}(\Psi_t, \mathbf{v}_{t-1}, \mathbf{V}_{t-1}) &= \int_{\Theta} H(\psi_t, \Theta) \sum_{c \in \mathbf{c}} \beta_{t;c} \tilde{Q}_{t;c}(\Theta) d\Theta. \end{aligned} \quad (19)$$

The determination of Kerridge's inaccuracy (16) needs to evaluate the integrals

$$\begin{aligned} I_{0c}(\Psi_t, \mathbf{v}_{t-1}, \mathbf{V}_{t-1}) &= \int_{\Theta} \ln(\alpha_c) H(\psi_t, \Theta) \sum_{\tilde{c} \in \mathbf{c}} \beta_{t;\tilde{c}} \tilde{Q}_{t;\tilde{c}}(\Theta) d\Theta, \quad I_0 = (I_{0c})_{c \in \mathbf{c}}, \\ I_c(\Psi_t, \mathbf{v}_{t-1}, \mathbf{V}_{t-1}) &= \int_{\Theta} C_c(\Theta) H(\psi_t, \Theta) \sum_{\tilde{c} \in \mathbf{c}} \beta_{t;\tilde{c}} \tilde{Q}_{t;\tilde{c}}(\Theta) d\Theta, \quad c \in \mathbf{c}. \end{aligned} \quad (20)$$



Their insertion into Kerridge's inaccuracy (16) gives the best-projection statistic

$$\begin{aligned}
\hat{\mathbf{V}}_t &= [\hat{\mathbf{v}}_t, (\hat{\mathbf{V}}_{t;c})_{c \in \mathbf{c}}] \in \text{Arg} \min_{\hat{\mathbf{v}} \in \hat{\mathbf{V}}, (\hat{\mathbf{V}} \in \hat{\mathbf{V}}_c)_{c \in \mathbf{c}}} \left( \ln(\text{Be}(\hat{\mathbf{v}})) + \sum_{c \in \mathbf{c}} \ln(\mathbf{N}_c(\hat{\mathbf{V}}_c)) \right. \\
&\quad \left. - \left\langle \hat{\mathbf{v}} - 1, \frac{\mathbf{l}_0(\Psi_t, \mathbf{v}_{t-1}, \mathbf{V}_{t-1})}{\mathbf{l}_{-1}(\Psi_t, \mathbf{v}_{t-1}, \mathbf{V}_{t-1})} \right\rangle - \sum_{c \in \mathbf{c}} \left\langle \hat{\mathbf{V}}_c, \frac{\mathbf{l}_c(\Psi_t, \mathbf{v}_{t-1}, \mathbf{V}_{t-1})}{\mathbf{l}_{-1}(\Psi_t, \mathbf{v}_{t-1}, \mathbf{V}_{t-1})} \right\rangle \right) \\
&= \left[ \text{Arg} \min_{\hat{\mathbf{v}} \in \hat{\mathbf{V}}} \left( \ln(\text{Be}(\hat{\mathbf{v}})) - \left\langle \hat{\mathbf{v}} - 1, \frac{\mathbf{l}_0(\Psi_t, \mathbf{v}_{t-1}, \mathbf{V}_{t-1})}{\mathbf{l}_{-1}(\Psi_t, \mathbf{v}_{t-1}, \mathbf{V}_{t-1})} \right\rangle \right), \right. \\
&\quad \left. \left( \text{Arg} \min_{\hat{\mathbf{V}} \in \hat{\mathbf{V}}_c} \left( \ln(\mathbf{N}_c(\hat{\mathbf{V}}_c)) - \left\langle \hat{\mathbf{V}}_c, \frac{\mathbf{l}_c(\Psi_t, \mathbf{v}_{t-1}, \mathbf{V}_{t-1})}{\mathbf{l}_{-1}(\Psi_t, \mathbf{v}_{t-1}, \mathbf{V}_{t-1})} \right\rangle \right) \right)_{c \in \mathbf{c}} \right]. \quad (21)
\end{aligned}$$

Thus,  $|\mathbf{c}| + 1$  independent minimisations are solved. They are numerically simple. Parts of them have analytical solutions for normal and Markov chain models.

*Numerical Evaluation of l's:* Function  $\mathbf{H}(\psi_t, \Theta)$  (19) depends on all involved parameters and combines influences of respective components. A brute force evaluation of  $\mathbf{l}_{-1}$ ,  $\mathbf{l}_0$  and  $\mathbf{l}_c$ ,  $c \in \mathbf{c}$ , say by Monte Carlo, is mostly too demanding in a typical on line learning. At the same time, the averaged functions in (19), (20) are by their construction smooth. Moreover, their growth over their domains is well suppressed by generally light tails of involved pds  $\tilde{\mathbf{Q}}_{t;c}$  (18). Thus, we conjecture that the simplest approximation based on the first order Taylor expansion of  $\mathbf{H}(\psi, \Theta)$  around expected values

$$\int_{\Theta} \alpha_{\tilde{c}} \tilde{\mathbf{Q}}_{t;c}(\Theta) d\Theta, \quad \int_{\Theta} \Theta_{\tilde{c}} \tilde{\mathbf{Q}}_{t;c}(\Theta) d\Theta, \quad \tilde{c}, c \in \mathbf{c}, \quad (22)$$

suffices. In addition to (22), such an approximation requires to evaluate

$$\int_{\Theta} \alpha \ln(\alpha_{\tilde{c}}) \tilde{\mathbf{Q}}_{t;c}(\Theta) d\Theta, \quad \int_{\Theta} \mathbf{C}_{\tilde{c}}(\Theta_{\tilde{c}}) \tilde{\mathbf{Q}}_{t;c}(\Theta) d\Theta, \quad \int_{\Theta} \Theta_{\tilde{c}} \mathbf{C}_{\tilde{c}}(\Theta_{\tilde{c}}) \tilde{\mathbf{Q}}_{t;c}(\Theta) d\Theta, \quad \tilde{c}, c \in \mathbf{c}. \quad (23)$$

#### Remarks 4

- ✓ *Normalising constants  $\text{Be}$ ,  $\mathbf{N}$  (13), (14) of factors forming  $\tilde{\mathbf{Q}}_{t;c}$  (18) can be analytically expressed for Normal-inverse-Wishart and Dirichlet's pds, which are conjugated to normal and Markov chain components, Ruman (2018).*
- ✓  *$\tilde{c}$ -th entry of expectations (22), (23) 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 advocated model, on the novelty brought.*

*Forgetting:* It consists of the choice of the forgetting factor and its use. The formula for  $\lambda_t$  in (4) arisen as ratio of the predictive pd divided by the predictive pd constructed from the already updated posterior pd (both in data realisation). The studies of partial forgetting, Dedecius et al. (2009), imply that forgetting should be applied component-wise. The overall prediction quality then influences only

forgetting of component weights. The value of  $l_{-1}$  (19) is proportional to the value of the predictive pd. This gives forgetting factors  $\lambda_t, (\lambda_{t;c})_{c \in \mathbf{c}}, (4)$ ,

$$\lambda_t = \frac{\text{Be}(\hat{\mathbf{v}}_t) \prod_{\tilde{c} \in \mathbf{c}} \text{N}(\hat{\mathbf{V}}_{t;\tilde{c}})}{\text{Be}(\mathbf{v}_{t-1}) \prod_{\tilde{c} \in \mathbf{c}} \text{N}(\mathbf{V}_{t-1;\tilde{c}})} \frac{l_{-1}(\Psi_t, \mathbf{v}_{t-1}, \mathbf{V}_{t-1})}{l_{-1}(\Psi_t, \hat{\mathbf{v}}_t, \hat{\mathbf{V}}_t)}$$

$$\lambda_{t;c} = \frac{\text{N}(\mathbf{V}_{t-1;c} + \mathbf{B}_c(\Psi_{t;c})) \text{N}(\hat{\mathbf{V}}_{t;c})}{\text{N}(\mathbf{V}_{t-1;c}) \text{N}(\hat{\mathbf{V}}_{t;c} + \mathbf{B}_c(\Psi_{t;c}))} \frac{W_c(\psi_t, \tilde{\Theta}_{t;c})}{W_c(\psi_t, \tilde{\Theta}_{t-1;c})}, \quad c \in \mathbf{c}.$$

There  $l_{-1}(\Psi_t, \mathbf{v}_{t-1}, \mathbf{V}_{t-1})$  is approximately evaluated using expansion around (22). Denominator  $l_{-1}(\Psi_t, \hat{\mathbf{v}}_t, \hat{\mathbf{V}}_t)$  is computed according to the same formulae with  $\hat{\mathbf{v}}_t, \hat{\mathbf{V}}_t$  (21) replacing  $\mathbf{v}_{t-1}, \mathbf{V}_{t-1}$ . Consequently,  $\tilde{\Theta}_{t-1;c}, \tilde{\Theta}_{t;c}$  are the expected values of component parameter  $\Theta_c$  computed as the  $c$ -th part of (22) with component statistics  $\mathbf{V}_{t-1;c}$  and  $\hat{\mathbf{V}}_{t;c}$ , respectively (11). The forgetting use completes updating (4) for proposed mixture ratio model (11), (12)

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

## 5 Illustrative Examples

This section illustrates the modelling theory and the corresponding learning. A systematic simulation and real-life studies are out of scope of this paper and will be published independently.

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. Additional experiments, including real-data application to commodity futures, Carter (1999), are in Ruman (2018).

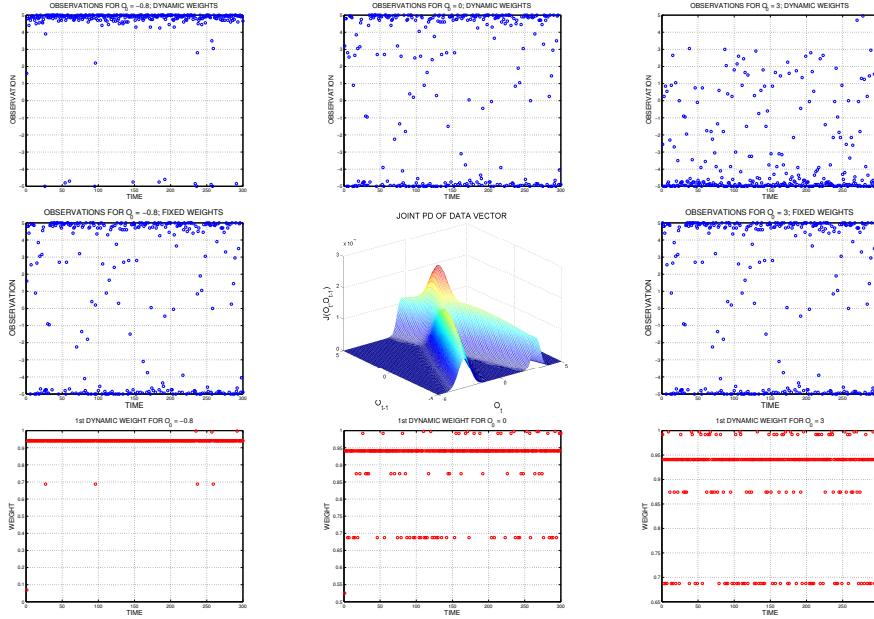
### 5.1 Mixture-Ratio Modelling Strength

Simulated environment model E, generating real scalar observations  $O_t$  with no actions and regression vector  $\psi_t = O_{t-1}$ , is the mixture of joint normal pds  $\mathcal{N}_{\Psi}(\mu_c, \omega_c^{0.5})$  with expectations  $\mu_c$  and square roots  $\omega_c^{0.5}$  of precision matrices,  $c \in \mathbf{c} = \{1, 2\}$ ,

$$\mathbb{E}(O_t, O_{t-1}) = \tag{24}$$

$$\frac{1}{2} \underbrace{\mathcal{N}_{O_t, O_{t-1}}}_{\Psi_t} \left( \underbrace{\begin{bmatrix} 1 \\ 1 \end{bmatrix}}_{\mu_1}, \underbrace{\begin{bmatrix} 1 & 3 \\ 0 & 0.5 \end{bmatrix}}_{\omega_1^{0.5}} \right) + \frac{1}{2} \underbrace{\mathcal{N}_{O_t, O_{t-1}}}_{\Psi_t} \left( \underbrace{\begin{bmatrix} -1 \\ -1 \end{bmatrix}}_{\mu_2}, \underbrace{\begin{bmatrix} 1 & -2 \\ 0 & 0.5 \end{bmatrix}}_{\omega_2^{0.5}} \right).$$

The typical simulation results, differing only in initial values  $O_0 \in \{-0.8, 0, 3\}$ , are shown in Fig. 1. It 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.



**Fig. 1** Simulation of environment (24). 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 realised observations when the same components as in (24) 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  $(w_{t,1}(\psi_t))_{t \in \mathcal{T}}$  (11) corresponding to the 1st-row.

## 5.2 Learning with Standard and Ratio Mixture Models

*Simulation Conditions:* Scalar observations  $O_t \in \mathcal{O} = \{1, \dots, |\mathcal{O}|\}$ ,  $|\mathcal{O}| = 5$ , and 2nd order regression vectors  $\psi_t = (O_{t-1}, O_{t-2})$  are considered. This defines data vectors (6)  $\Psi_t = (O_t, O_{t-1}, O_{t-2})$ . Learnt models  $\mathbf{L}$  and simulated environment  $\mathbf{E}$  are either mixture ratio model  $\mathbf{M}$  or standard mixture  $\mathbf{S}$  with two components  $|\mathcal{c}| = 2$

$$\mathbf{E}, \mathbf{L} \in \mathbf{L} = \{\mathbf{M}, \mathbf{S}\} = \{\text{mixture ratio, standard mixture}\}. \quad (25)$$

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

$$J(\Psi_t|\theta_M) = \alpha_M J_1(O_t, O_{t-1}|\theta_{M1}) G_1(O_{t-2}) + (1 - \alpha_M) J_2(O_t, O_{t-2}|\theta_{M2}) G_2(O_{t-1}).$$

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  $\theta_{M1}$ ,  $\theta_{M2}$  of possible pairs  $(O_t, O_{t-1})$ ,  $(O_t, O_{t-2})$ , and  $\theta_M = (\alpha_M, \theta_{M1}, \theta_{M2})$ .  $G_1(O_{t-2})$  and  $G_2(O_{t-1})$  are uniform pds on  $\mathcal{O}$ . This gives mixture ratio  $\mathbf{M}$

$$M(O_t|\psi_t, \theta_M) = \frac{\alpha_M J_1(O_t, O_{t-1}|\theta_{M1}) + (1 - \alpha_M) J_2(O_t, O_{t-2}|\theta_{M2})}{\sum_{O_t \in \mathcal{O}} \alpha_M J_1(O_t, O_{t-1}|\theta_{M1}) + (1 - \alpha_M) J_2(O_t, O_{t-2}|\theta_{M2})}.$$

Used standard mixture model  $\mathbf{S}$  is parameterised by  $\theta_{\mathbf{S}} = (\alpha_{\mathbf{S}}, \theta_{\mathbf{S}1}, \theta_{\mathbf{S}2})$

$$\mathbf{S}(O_t|\psi_t, \theta_{\mathbf{S}}) = \alpha_{\mathbf{S}}\mathbf{S}_1(O_t|O_{t-1}, \theta_{\mathbf{S}1}) + (1 - \alpha_{\mathbf{S}})\mathbf{S}_2(O_t|O_{t-2}, \theta_{\mathbf{S}2}),$$

where  $\mathbf{S}_1(O_t|O_{t-1}, \theta_{\mathbf{S}1})$  and  $\mathbf{S}_2(O_t|O_{t-2}, \theta_{\mathbf{S}2})$  are components parameterised by probabilities  $(\theta_{\mathbf{S}1}, \theta_{\mathbf{S}2})$  of  $O_t$  conditioned on  $O_{t-1}$  and  $O_{t-2}$ , respectively.

*Monte Carlo Study:* It consists of 200 runs for time steps  $t \in \mathbf{t}$ ,  $|\mathbf{t}| = 500$ , with randomly generated parameters  $\theta_{\mathbf{E}}$  for both model structures  $\mathbf{E} \in \mathbf{L} = \{\mathbf{M}, \mathbf{S}\}$  (25). In both cases, the mixture ratio and standard mixture models are sequentially learnt on the identical observations. An individual study is also considered with a fixed data-generating mixture ratio model given by parameter  $\underline{\theta}_{\mathbf{M}}$  with fixed  $\underline{\theta}_{\mathbf{M}}$  causing truly dynamic weights  $w_c(\psi_t)$  (11). These 200 runs differ in realisations.

*Evaluation:* At time  $t \in \mathbf{t}$ , the simulation and the sequential learning dealt with the next *predictors*, with an abused notation pointing to the underlying parametric models,

$$\begin{aligned} \mathbf{E}_t(O_t) &= \begin{cases} \underline{\mathbf{M}}_t(O_t) = \mathbf{M}(O_t|\psi_t, \underline{\theta}_{\mathbf{M}}) & \text{if } \mathbf{E} = \mathbf{M} \text{ with fixed parameter } \underline{\theta}_{\mathbf{M}} \\ \underline{\mathbf{S}}_t(O_t) = \mathbf{S}(O_t|\psi_t, \underline{\theta}_{\mathbf{S}}) & \text{if } \mathbf{E} = \mathbf{S} \text{ with fixed parameter } \underline{\theta}_{\mathbf{S}} \end{cases} \quad (26) \\ \mathbf{L}_t(O_t) &= \mathbf{L}_t(O_t|O_{t-1}, \dots, O_0) = \int_{\Theta_{\mathbf{L}}} \mathbf{L}(O_t|\psi_t, \theta_{\mathbf{L}}) \mathbf{P}(\theta_{\mathbf{L}}|O_{t-1}, \dots, O_0) d\theta_{\mathbf{L}}, \quad \mathbf{L} \in \mathbf{L}. \end{aligned}$$

The sequentially evaluated Kullback-Leibler divergences of environment  $\mathbf{E}$  and learnt predictors (26)

$$\begin{aligned} \mathbf{D}_t(\mathbf{E}||\mathbf{L}) &= \sum_{\tau=1}^t \mathbf{D}(\mathbf{E}_{\tau}||\mathbf{L}_{\tau}), \quad \text{with} \quad \mathbf{D}(\mathbf{E}_{\tau}||\mathbf{L}_{\tau}) = \sum_{O_{\tau} \in \mathbf{O}} \mathbf{E}_{\tau}(O_{\tau}) \ln \left( \frac{\mathbf{E}_{\tau}(O_{\tau})}{\mathbf{L}_{\tau}(O_{\tau})} \right) \\ \mathbf{D}(\mathbf{E}||\mathbf{L}) &= \mathbf{D}_{|\mathbf{t}|}(\mathbf{E}||\mathbf{L}), \quad \text{and their differences} \quad \Delta = \mathbf{D}(\mathbf{E}||\mathbf{M}) - \mathbf{D}(\mathbf{E}||\mathbf{S}) \quad (27) \end{aligned}$$

quantify learning quality in the comparable observation space.

*Results:* Results of Monte Carlo simulations are in Fig. 2. The left hand column contains histogram values  $\mathbf{D}(\mathbf{E}||\mathbf{L}) = \mathbf{D}_{|\mathbf{t}|}(\mathbf{E}||\mathbf{L})$ ,  $\mathbf{L} \in \mathbf{L}$ , at final time  $|\mathbf{t}| = 500$ . The right hand column shows the corresponding sample cumulative distributions. The better predictors have them shifted to the left. Fig. 3 shows course of  $\mathbf{D}_t(\mathbf{E}||\mathbf{L})$ ,  $t \in \mathbf{t}$ ,  $\mathbf{L} \in \mathbf{L} = \{\mathbf{M}, \mathbf{S}\}$  for: (a)  $\mathbf{E} = \underline{\mathbf{M}}$ , given by fixed  $\underline{\theta}_{\mathbf{M}}$ ; (b)  $\mathbf{E} = \underline{\mathbf{S}}$ , given by fixed  $\underline{\theta}_{\mathbf{S}}$ ; (c)  $\mathbf{E} = \underline{\mathbf{M}}$  with  $\underline{\theta}_{\mathbf{M}}$  giving dynamic weights  $w_c(\psi_t)$  (11).

Table 1 shows sample statistics of increments  $\Delta$  (27). Negative values mean the dominance of the mixture ratio model.

*Discussion:* Fig. s 2, 3 and Table 1 confirm the expectation that the mixture ratio outperforms the standard mixture if the mixture ratio is simulated. It leads to essentially same results when the standard mixture serves as the simulated environment. It is even slightly better in the latter case as it (probably) copes better with errors caused by approximate learning (4).

**Table 1** Sample statistics of  $\Delta$  (27): the 1st column  $E = \underline{M}$  with  $\underline{\Theta}_M$  varied in 200 Monte Carlo runs; the 2nd column  $E = \underline{S}$  with  $\underline{\Theta}_S$  varied in 200 Monte Carlo runs; the 3rd column  $E = \underline{M}$  with fixed  $\underline{\Theta}_M$  causing truly dynamic weights  $w_c(\psi_t)$  (11).

Simulated Case	Monte Carlo $E = \underline{M}$	Monte Carlo $E = \underline{S}$	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

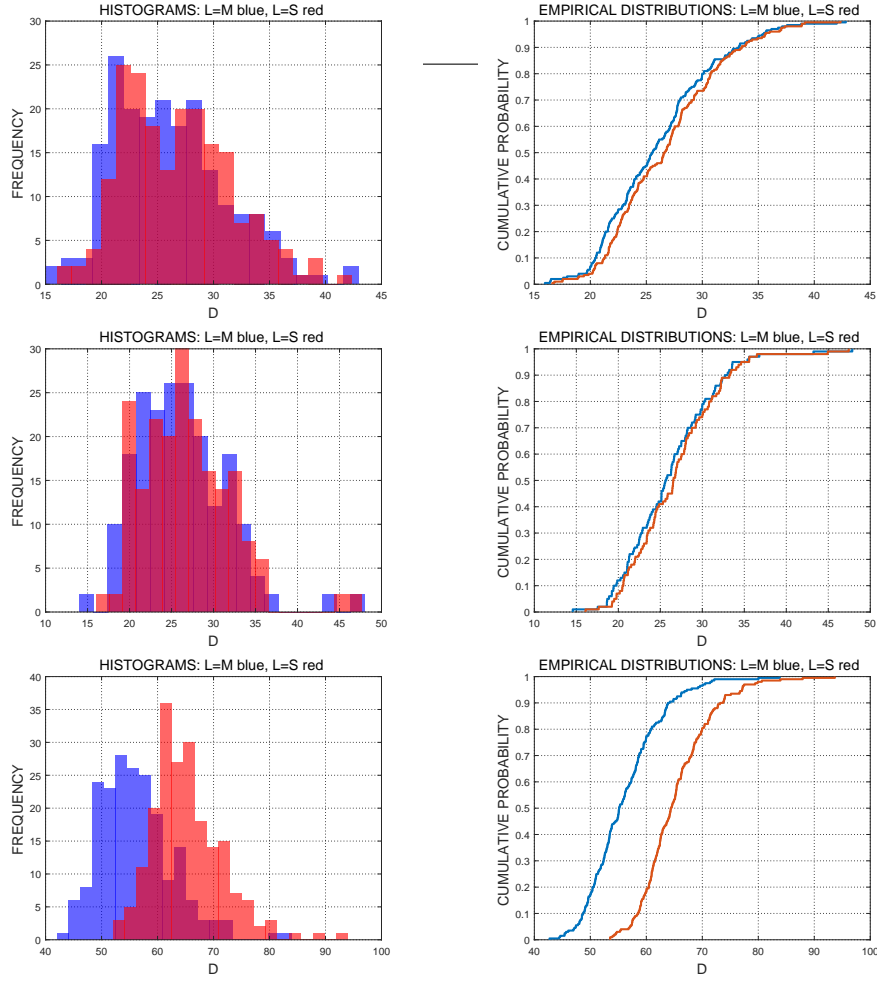
## 6 Potential of the Mixture Ratio Model

The paper brings an important message that the ratios of finite mixtures model well dynamic relations. The main promising features are as follows.

- ✓ The mixture ratios universally describe non-linear dynamic stochastic relations.
- ✓ The mixture ratios may serve as a relatively universal dynamic feature extractor. Indeed, the approximately sufficient statistic collected during learning, see Sec. 4, are the relevant features. Thus, it suffices to apply Bayesian structure estimation to mixture ratios with low-dimensional components.
- ✓ The mixture ratios handle cases in which rare visits of active components are significant, cf. Fig. 1. This is vital in fault detection, Polycarpou and Helmicki (1995), detection of non-standard fraud behaviours needed, Kou et al. (2004), or in cyber-security applications Buczak and Guven (2016), generally everywhere where outlier detection is faced, Hodge and Austin (2004). These cases are hard dynamic versions of learning with unbalanced data, Bekkar and Ali-touche (2013).
- ✓ The mixture ratios learn the stationary joint pd of the data vector (10), which can directly serve for the design of adaptive decision strategies for an infinite decision horizon, Kushner (1971). Indeed, the sequentially updated 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. This fits well to fully probabilistic design of decision strategies, Kárný and Kroupa (2012).
- ✓ The mixture ratios suit to modelling of mixed (discrete and continuous valued) data as the components are joint pds, which factorise into product of pds describing continuous valued data, possibly conditioned on discrete ones, and pds relating discrete values.

The paper indicates tractability all these cases, but it is necessary:

- ✓ to elaborate numerically robust (factorised) procedures for normal models, sparse Markov-chains and mixed cases in the way mimic to Kárný (2016); Ruman (2018) has made definite steps in this respect;
- ✓ to refine the data-dependent choice of forgetting factors;
- ✓ to tailor Bayesian structure estimation, ready for mixtures, Kárný et al. (2006), to the mixture-ratio model;
- ✓ to perform real-life tests confirming that the theoretically higher approximating strength of the discussed mixture ratios is mostly undiminished by the



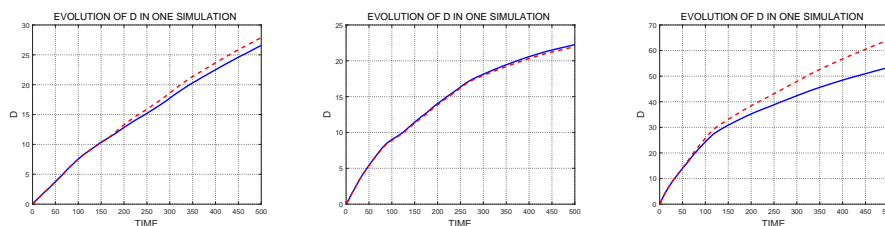
**Fig. 2** Histograms of  $D(E||L)$  ((27), left) and their empirical distribution function of (right),  $L \in \{M, S\} = \{\text{blue, red}\}$ . The 1st row  $E = \underline{M}$  with  $\underline{\theta}_M$  varied in 200 Monte Carlo runs; the 2nd row  $E = \underline{S}$  with  $\underline{\theta}_S$  varied in 200 Monte Carlo runs; the 3rd row for  $E = \underline{M}$  with fixed  $\underline{\theta}_M$  causing truly dynamic weights  $w_c(\psi_t)$  (11).

approximate learning; the experiments with commodity futures trading made in Ruman (2018), are more than promising.

The theoretical insight, experience from simulations and preliminary processing of real data, width and importance of possible contributions to machine learning and DM applications make this effort worthwhile.

#### Acknowledgement

This research is supported by GAČR GA16-09848S.



**Fig. 3** Time evolutions of  $D_t(E||L)$  (27),  $L \in \{M, S\} = \{\text{blue}, \text{red}\}$ . Left to right:  $E = \underline{M}$  with randomly chosen  $\underline{\Theta}_M$ ;  $E = \underline{S}$  with randomly chosen  $\underline{\Theta}_S$ ;  $E = \underline{M}$  with fixed  $\underline{\Theta}_M$  causing truly dynamic weights  $w_c(\psi_t)$  (11).

## References

- Appice A, Ciampi, Malerba D (2015) Summarizing numeric spatial data streams by trend cluster discovery. *Data Mining and Knowledge Discovery* 29(1):84–136
- Åström K, Wittenmark B (1989) *Adaptive Control*. Addison-Wesley
- Barndorff-Nielsen O (1978) *Information and Exponential Families in Statistical Theory*. Wiley, NY
- Bekkar M, Alitouche T (2013) Imbalanced data learning approaches review. *Int J of Data Mining & Knowledge Management Process* 3(4):15–33
- Berec L, Kárný M (1997) Identification of reality in Bayesian context. In: Warwick K, Kárný M (eds) *Computer-Intensive Methods in Control and Signal Processing*, Birkhäuser, pp 181–193
- Berger J (1985) *Statistical Decision Theory and Bayesian Analysis*. Springer, NY
- Bernardo J (1979) Expected information as expected utility. *The An of Stat* 7(3):686–690
- Bishop C (2006) *Pattern Recognition and Machine Learning*. Springer
- Bohlin T (1991) *Interactive System Identification: Prospects and Pitfalls*. Springer
- Buczak A, Guven E (2016) A survey of data mining and machine learning methods for cyber security intrusion detection. *IEEE Communications Surveys Tutorials* 18(2):1153–1176
- Carter C (1999) Commodity future markets: A survey. *The Australian Journal of Agricultural and Resource Economics* 43:209–247
- Catania L (2016) Dynamic adaptive mixture models. *ArXiv:1603.01308v1*
- Dedecius K, Nagy I, Kárný M, Pavelková L (2009) Parameter estimation with partial forgetting method. In: *Proc. of the 15th IFAC SYSID*
- Elliot R, Assoun L, Moore J (1995) *Hidden Markov Models*. Springer-Verlag, NY
- Frigessi A, Haug O, Rue H (2002) A dynamic mixture model for unsupervised tail estimation without threshold selection. *Extremes* 5(3):219 – 235
- Guan P, Raginsky M, Willett R (2012) Online Markov decision processes with Kullback-Leibler control cost. In: *Am. Control Conference, IEEE*, pp 1388–1393
- Guyon I, Saffari A, Dror G, Cawley G (2010) Model selection: Beyond the Bayesian/frequentist divide. *Journal of Machine Learning Research* 11:61–87
- Haykin S (1994) *Neural Networks: A Comprehensive Foundation*. Macmillan, NY
- Hodge V, Austin J (2004) A survey of outlier detection methodologies. *Artificial Intelligence Review* 22(2):85–126

- Holmes G, Liu T (eds) (2015) Proceedings of 7th Asian Conference on Machine Learning (ACML2015), JMLR Workshop and Conference Proceedings, vol 45
- Jazwinski A (1970) Stochastic Processes and Filtering Theory. Ac. Press, NY
- Kárný M (2014) Approximate Bayesian recursive estimation. *Infor Sciences* 289:100–111
- Kárný M (2016) Recursive estimation of high-order Markov chains: Approximation by finite mixtures. *Infor Sciences* 326:188–201
- Kárný M, Kroupa T (2012) Axiomatisation of fully probabilistic design. *Information Sciences* 186(1):105–113
- Kárný M, Böhm J, Guy TV, Jirsa L, Nagy I, Nedoma P, Tesař L (2006) Optimized Bayesian Dynamic Advising: Theory and Algorithms. Springer
- Kárný M, Macek K, Guy T (2014) Lazy fully probabilistic design of decision strategies. In: 11th Int. Symp. on Neural Net., no. 8866 in LNCS, pp 140–149
- Kerridge D (1961) Inaccuracy and inference. *J of the Royal Stat Society B* 23:284–294
- Koopman R (1936) On distributions admitting a sufficient statistic. *Trans of Am Math Society* 39:399
- Kou Y, Lu C, Sirwongwattana S, Huang Y (2004) Survey of fraud detection techniques. In: IEEE International Conference on Networking, Sensing and Control, 2004, vol 2, pp 749–754 Vol.2
- Kulhavý R (1990) A Bayes-closed approximation of recursive nonlinear estimation. *Int J Adaptive Control and Signal Proc* 4:271–285
- Kullback S, Leibler R (1951) On information and sufficiency. *Annals of Mathematical Statistics* 22:79–87
- Kushner H (1971) Introduction to Stochastic Control. Holt, Rinehart and Winston
- Mason W, Vaughan J, Wallach H (2014) Special issue: Computational social science and social computing. *Machine Learning* 96:257–469
- McLachlan G, Peel D (2000) Finite Mixture Models. Wiley Series in Probab. & Stat., Wiley, NY
- McNicholas P (2017) Mixture model-based classification. CRC Press, Boca Raton, London, New York
- Mirsky Y, Shapira B, Rokach L, Elovici Y (2015) 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, vol 9078, LNCS Springer, Cham, pp 230–237
- Mitchell T (1997) Machine Learning. McGraw Hill
- Nelsen R (1999) An Introduction to Copulas. Springer, NY
- Nguyen H, Woon Y, Ng W (2015) A survey on data stream clustering and classification. *Knowledge and Information Systems* 45:535–569, DOI <https://doi.org/10.1007/s1011>
- Park J, Sandberg I (1991) Universal approximation using radial-basis-function networks. *Neural Computations* 3:246–257
- Peterka V (1981) Bayesian system identification. In: Eykhoff P (ed) Trends and Progress in System Identification, Pergamon Press, Oxford, pp 239–304
- Polycarpou M, Helmicki A (1995) Automated fault detection and accommodation: a learning systems approach. *IEEE Transactions on Systems, Man, and Cybernetics* 25(11):1447–1458
- Rao M (1987) Measure Theory and Integration. John Wiley, NY



- Ruman F (2018) Mixture ratios for decision making. Master's thesis, FJFI, Czech Technical University, Prague
- Sadghi P, Kennedy R, Rapajic P, Shams R (2008) Finite-state Markov modeling of fading channels. *IEEE Signal Processing Magazine* 57, DOI 10.1109/MSP.2008.926683
- Savage L (1954) *Foundations of Statistics*. Wiley, NY
- Silva J, Faria E, Barros R, Hruschka E, Carvalho A, Gama J (2014) Data stream clustering: A survey. *ACM Computing Surveys* 46, DOI 10.1145/2522968.2522981
- Tipping M, Bishop C (1999) Probabilistic principal component analysis. *J of the Royal Society Series B – Stat Methodology* 61:611–622
- Tsai C, Lai C, Chiang M, Yang L (2014) Data mining for internet of things: A survey. *IEEE Communications Surveys & Tutorials* 16(1):77–95
- Wald A (1950) *Statistical Decision Functions*. John Wiley, NY, London
- Wiering M, van Otterlo M (eds) (2012) *Reinforcement Learning: State-of-the-Art*. Springer