



Approximate Recursive Bayesian Estimation of Dynamic Probabilistic Mixtures'/'FTCHV

Josef Andryšek

Department of Adaptive Systems
Institute of Information Theory and Automation
Academy of Sciences of the Czech Republic

Abstract. Majority of complex non-linear systems can be successfully modelled by a finite probabilistic mixture of linear models. The mixture model can be handled analytically, which is important for control of the system as well as for decision making. Quality of the model is a crucial requirement of all tasks of this type. The exact Bayesian methodology can not be used for estimation of this model, because complexity of the posterior distribution grows exponentially with number of data. Therefore, approximation techniques such as the quasi-Bayes algorithm must be used. This paper introduces a new estimation algorithm, which is based on minimization of Kullback-Leibler distance between the proper Bayesian posterior density and an approximate posterior density. The approximate posterior distribution is chosen from the exponential family in order to achieve numerically efficient estimation.

Keywords: parameter estimation, recursive estimation, probabilistic mixture, control of the complex system, system model

1 Introduction

The choice of a suitable model is essential for both control and decision making when dealing with complex systems. One way to face complexity is the principle of adaptivity, i.e. using models which evolve during their use. The demand for adaptivity of the model leads to the recursive estimation of its parameters, i.e. permanent updating of its parameter estimates by the new data. In other words, statistics describing estimates are corrected by newly acquired data. The model should be chosen from a sufficiently rich family of models to capture all properties of the modelled system. Naturally, computational cost associated with estimation of parameters of the model grows with complexity of the model. If the modelled system is non-linear, its model should be non-linear too. In this paper, we study finite probabilistic mixture of linear models. The finite mixtures provide a universal approximation of almost any probabilistic density function (Titterton, et al. 1985) and thus can be successfully used in modelling of complex systems. Invoking the principle of adaptivity, we seek an efficient recursive estimation of the mixture model parameters.

The resulting model can be then used both for control and decision making tasks. Universal algorithms for mixture-based control (Kárný, et al. 2003) were derived, but quality of the resulting control strategy strongly depends on the quality of the estimated model. Practical experience indicates that this is a weak element of adaptive control and that an improvement of the estimation part improves the overall control quality. Hence, we try to develop better estimation algorithms for the mixture model. The control algorithms (Kárný et al. 2003) as well as efficient structure estimation algorithms were derived using the Bayesian theory (Peterka 1981). The unknown model parameters are treated as random variables and all subsequent task are defined in terms of posterior distributions of the parameters rather than their point estimates.

The recursive Bayesian estimation evaluates the posterior distribution on parameters at time t as an update of the the posterior distribution at time $t - 1$ using the Bayes' rule and the data acquired at time t . The recursion starts at $t = 1$ with update of the prior distribution which must be chosen before the estimation starts. The posterior distribution obtained by the Bayes' rule may not be, however, analytically tractable and thus unsuitable for the next update.

In practice, mostly such prior distribution is used so that the posterior distribution in each estimation step has the same functional form as the prior distribution. Hence, just the sufficient statistics determining the posterior density are updated. Such a prior distribution is then known as conjugate with the observation model. For example, conjugate prior distribution is available for all models from the exponential family. If the conjugate prior does not exist the exact recursive estimation can not be achieved. In such a case, we seek approximate recursive estimation. This is the case of the probabilistic mixture model. Using the exact Bayesian update, the complexity of posterior density grows exponentially with number of the data samples. The quasi-Bayes





algorithm (Kárný, et al. 1998) (Kárný et al. 2003) or a modification of the EM algorithm (Titterton et al. 1985) are examples of approximate algorithms facing this problem.

This paper introduces a new approximate estimation method, which can be viewed as a generalization of quasi-Bayes algorithm. The basis of both approaches is finding the approximate posterior density in particular (well manipulable) class of densities.

The new algorithm finds the optimal projection of the correct Bayesian density into the selected class of densities. The projection is optimal in the sense of Kulback Leibler distance (Kullback & Leibler 1951). It should be mentioned that the Kullback Leibler distance is not symmetric. Algorithm presented in this paper minimizes the Kulback Leibler distance with the argument order, which conforms with Bayesian principles (Berec & Kárný 1997). An algorithm minimizing the Kullback Leibler distance with arguments in different order can be found in (Roberts & Penny 2002).

2 Notions and notations

x^* denotes the range of x , $x \in x^*$.

\hat{x} denotes the number of entries in the vector x .

\equiv means the equality by definition.

x_t is a quantity (vector) x at the discrete time labelled by $t \in t^* \equiv \{1, \dots, \hat{t}\}$.

$x_{i;t}$ is an i -th entry of the vector x_t . The semicolon in the subscript indicates that the symbol following it is the time index.

$x_{k;l;t}$ is a subvector of the vector x_t . $x_{k;l;t} = (x_{k;t}, \dots, x_{l;t})$.

$x(k_l) \equiv x_k, \dots, x_l$.

$x(t) \equiv x(1_t)$.

$x(t)$ is an empty sequence and reflects just the prior information if $t < 1$.

d is data array, d_t is data record at time t (vector with entries $(d_{1;t}, \dots, d_{\hat{d}_t;t})$).

Θ unknown parameter, finite-dimensional vector

f, π are the letters reserved for probability density functions(pdf).

$f(d_t | d(t-1), \Theta)$ means model of the system.

$f_c(d_t | d(t-1), \Theta_c)$ is component of the mixture.

$\pi_0(\Theta)$ denotes prior density of the unknown parameter Θ .

$\pi_t(\Theta | d(t)) \equiv \pi_t(\Theta | \mathcal{G}_t)$ means (approximate) posterior density of the parameter Θ determined by the sufficient statistic \mathcal{G}_t .

\propto is the proportion sign, $h \propto g$ means that function h equals to the function g up to the normalization. I.e.

$$\frac{h}{\int h} = \frac{g}{\int g}.$$

∂ is the model order.

$\mathcal{D}(\| \|)$ means the Kullback-Leibler distance (Kullback & Leibler 1951). This "distance" is familiarly used in Bayesian analysis as the measure how good the second pdf approximates the first pdf. For conciseness, the Kullback-Leibler distance is referred to as the KL distance. $\mathcal{D}\left(f \parallel g\right) = \int f \ln\left(\frac{f}{g}\right)$

$\Gamma(x)$ means gamma function, $\Gamma(x) = \int_0^{+\infty} t^{x-1} \exp(-t) dt$.

$\psi_0(x)$ is digamma function, $\psi_0(x) = \frac{\partial \ln \Gamma(x)}{\partial x}$.

δ denotes identity matrix. I.e. $\delta_{ij} = 1$ iff $i = j$, otherwise $\delta_{ij} = 0$.

Agreement 1 (Multimatrix, multivector) Multimatrix of type m, n

$$M = \begin{pmatrix} M_{11} & \cdots & M_{1n} \\ \vdots & \ddots & \vdots \\ M_{m1} & \cdots & M_{mn} \end{pmatrix}$$

is a mathematical object, where M_{ij} is either matrix or multimatrix. Hence matrix is a multimatrix. Multimatrix need not be a matrix. Definition of Multivector is analogical.

Agreement 2 (Multimatrix indexing) For M being a multimatrix of type m, n the following notation is used:

M_{ij} is ij -th entry of M .





$M_{\bullet j}$ is multimatrix $\begin{pmatrix} M_{1j} \\ \vdots \\ M_{mj} \end{pmatrix}$.

$M_{i\bullet}$ is multimatrix (M_{i1}, \dots, M_{in}) .

$M_{\bullet\bullet}$ means the same as M . We use this notation when we want to stress that M is a multimatrix (matrix).

Agreement 3 (Other Matrix notations) Let's M be a matrix of type m, n and c some scalar. Let's define the following operations:

$M \pm c$ is matrix of type m, n , $(M \pm c)_{ij} = M_{ij} \pm c$.

$\exp(M)$ is matrix of type m, n , $(\exp(M))_{ij} = \exp(M_{ij})$.

$\max M$ is scalar with maximal value of M .

3 Basic elements and tools

3.1 Recursive parameter estimation

The task of recursive parameter estimation is to determine the posterior density $\pi_t(\Theta|d(t))$ based on the knowledge of

- last posterior density $\pi_{t-1}(\Theta|d(t-1))$
- new data record d_t
- model of the system $f(d_t|d(t-1), \Theta)$ parameterized by unknown parameter Θ .

The algorithm starts from prior pdf $\pi_0(\Theta) \equiv \pi_0(\Theta|d(0))$. We assume existence of the sufficient statistic \mathcal{G}_t for posterior pdfs, i.e.

$$\pi_t(\Theta|d(t)) \equiv \pi_t(\Theta|\mathcal{G}_t).$$

Next, consider that the actual data record d_t doesn't depend on all historical data $d(t-1)$ but only on a subset $\phi_{t-1} = (d_{t-1}, d_{t-2}, \dots, d_{t-\delta})$. Hence,

$$f(d_t|d(t-1), \Theta) \equiv f(d_t|\phi_{t-1}, \Theta).$$

The standard Bayesian approach determines $\pi_t(\Theta|\mathcal{G}_t)$ as

$$\pi_t(\Theta|\mathcal{G}_t) \propto f(d_t|\phi_{t-1}, \Theta)\pi_{t-1}(\Theta|\mathcal{G}_{t-1}). \quad (1)$$

3.1.1 Recursive parameter estimation with conjugate pdf

The considered approach (1) can be effectively used in the case when $\pi_0(\Theta)$ is conjugate pdf to the system model $f(d_t|\phi_{t-1}, \Theta)$. In such a case, $\pi_t(\Theta|\mathcal{G}_t)$ has the same functional form as $\pi_0(\Theta)$. Hence, we can get

$$\pi_t(\Theta|\mathcal{G}_t) \equiv \pi(\Theta|\mathcal{G}_t), \quad \forall t.$$

When updating from $\pi(\Theta|\mathcal{G}_{t-1})$ to $\pi(\Theta|\mathcal{G}_t)$ it suffices to update the sufficient statistics: $(\mathcal{G}_{t-1}, d_t) \longrightarrow \mathcal{G}_t$.

3.1.2 Recursive parameter estimation without conjugate pdf

If the pdf conjugate to the system model doesn't exist, the dimension of sufficient statistic grows with number of data samples. Then, of course, complexity of π_t grows as well. In such a case we can proceed in the following way:

- we choose prior pdf in an arbitrary well manipulable functional form,
- we seek an approximate posterior pdf's of the same functional form,
- we set, in each step of estimation, the statistic determining the approximate posterior pdf in such a way that it is "closest" to the "correct Bayesian" pdf.

We need to specify what we mean by: "correct Bayesian" and "closest". Let's have the approximate posterior pdf $\pi(\Theta|\mathcal{G}_{t-1})$, which depends on the statistic \mathcal{G}_{t-1} . If we handle the approximate posterior pdf $\pi(\Theta|\mathcal{G}_{t-1})$





as the correct posterior pdf, the "correct Bayesian" posterior pdf in the next step $\hat{\pi}(\Theta|\mathcal{G}_{t-1}, d_t, \phi_{t-1})$ is (according to (1)) obtained as

$$\hat{\pi}(\Theta|\mathcal{G}_{t-1}, d_t, \phi_{t-1}) = \frac{f(d_t|\phi_{t-1}, \Theta)\pi(\Theta|\mathcal{G}_{t-1})}{\int f(d_t|\phi_{t-1}, \Theta)\pi(\Theta|\mathcal{G}_{t-1})d\Theta}.$$

The term "closest" means closest in sense of the KL distance. It means that we want to find \mathcal{G}_t so that

$$\mathcal{D}\left(\hat{\pi}(\Theta|\mathcal{G}_{t-1}, d_t, \phi_{t-1}) \parallel \pi(\Theta|\mathcal{G}_t)\right) \quad (2)$$

is minimized.

Remarks 1

1. Applicability of the presented algorithm strictly depends on the complexity of the KL distance. Except of trivial cases, it is usable only if the KL distance can be evaluated analytically.
2. The algorithm uses the approximate posterior pdf obtained in step $t - 1$ as the true posterior pdf in step t . This leads to error accumulation.

3.2 Dynamic probabilistic mixture

In this paper, we consider the parameterized model of the system in the form of a finite probabilistic mixture:

$$f(d_t|\phi_{t-1}, \Theta) \equiv \sum_{c \in c^*} \alpha_c f_c(d_t|\phi_{c;t-1}, \Theta_c), \quad c^* = \{1, \dots, \hat{c}\}, \quad \hat{c} < \infty, \quad \text{where} \quad (3)$$

$$f_c(d_t|\phi_{c;t-1}, \Theta_c) \equiv \text{c-th component given by component parameters } \Theta_c \text{ and the state}$$

$$\phi_{c;t-1} \equiv \text{subset of } \phi_{t-1}$$

$$\alpha_c \equiv \text{the probabilistic component weight}$$

$$\Theta \equiv \text{mixture parameter formed by the component weights and parameters}$$

$$\Theta \in \Theta^* \equiv \left\{ \{\Theta_c \in \Theta_c^*\}_{c \in c^*}, \alpha \equiv [\alpha_1, \dots, \alpha_{\hat{c}}] \in \alpha^* \equiv \left\{ \alpha_c \geq 0, \sum_{c \in c^*} \alpha_c = 1 \right\} \right\}.$$

Before fixing and refining nomenclature related to the mixture, we split the individual components into so called *factors* that provide flexibility of the parametric description.

Using the chain rule, the pdfs $f_c(d_t|\phi_{c;t-1}, \Theta_c)$ can be written as a product of pdfs of individual entries of d_t . Before applying the chain rule, entries of d_t can be permuted and some permutations may lead to parameterizations with less parameters. This motivates inclusion of permutations into the model description

$$d \rightarrow d_c \text{ with } d_{ic} = d_{j_{ic}}, \text{ where} \quad (4)$$

j_{ic} is i -th entry of the permuted indices $[1, \dots, \hat{d}]$. The assignment (4) is applied component-wise and together with the chain rule give

$$f_c(d_t|\phi_{c;t-1}, \Theta_c) = \prod_{i \in i^*} f_{ic}(d_{ic;t}|d_{(i+1)\hat{d};t}, \phi_{c;t-1}, \Theta_{ic}) \equiv \prod_{i \in i^*} f_{ic}(d_{ic;t}|\psi_{ic;t}, \Theta_{ic}). \quad (5)$$

The additional subscript i of the parameter Θ_{ic} indicates that only some entries of Θ_c may occur in i -th pdf (*factor*) in (5). Similarly, the *regression vector* $\psi_{ic;t}$ is generally a sub-vector of the vector

$$[d_{(i+1)\hat{d};t}, \phi'_{c;t-1}, 1]'. \quad (6)$$

Agreement 4 (Nomenclature related to mixtures)

Pdfs: The pdf $f_c(d_t|\phi_{c;t-1}, \Theta_c)$ in (3) is called parameterized component of a mixture and α_c is the weight of the c -th parameterized component.

The pdf $f_{ic}(d_{ic;t}|\psi_{ic;t}, \Theta_{ic})$ in (5) is called parameterized factor.

Data: The vector d_t containing data measured at time t is called data record.

The vector $\phi_{c;t-1}$ is the observable state of the parameterized component.

The parameterized factor is determined by regression vector $\psi_{ic;t}$ defined as a sub-selection of the vector

$$[d_{(i+1)\hat{d};t}, \phi'_{c;t-1}, 1]' \quad (6).$$

The coupling $\Psi_{ic;t} \equiv [d_{ic;t}, \psi'_{ic;t}]'$ is called data vector of the factor.



Remarks 2

1. We added the number 1 to the definition of the regression vector, because it helps us to effectively express the constant shifts in mean values of factors.
2. The adopted dynamic mixture model is not sufficiently general. The component weights should also depend on the state vector. The choice is driven by our inability to estimate this “natural” and more realistic model. See discussion in (He & Kárný 2003)

3.3 Form of the prior and the posterior pdf

According to the general hints in section 3.1.2 we need to choose the prior pdf in a form that is well manipulable, i.e. analytically tractable.

Agreement 5 (Considered forms of pdfs on Θ^*) The prior $\pi(\Theta) \equiv \pi(\Theta|d(0))$ and the posterior $\pi(\Theta|d(t)) \equiv \pi(\Theta|\mathcal{G}_t)$ are considered to be of the common form:

$$\pi(\Theta|\mathcal{G}_t) = Di_\alpha(\kappa_t) \prod_{i \in i^*, c \in c^*} \pi_{ic}(\Theta_{ic}|\mathcal{S}_{ic;t}), \quad t \in \{0\} \cup t^* \quad , \quad \text{where} \quad (7)$$

$$\mathcal{G}_t \equiv (\kappa_{\bullet;t}, \mathcal{S}_{\bullet;t}),$$

$$Di_\alpha(\kappa_t) \text{ is Dirichlet distribution, } Di_\alpha(\kappa_\bullet) \equiv \frac{\prod_{c \in c^*} \alpha_c^{\kappa_c - 1}}{\mathcal{B}(\kappa)}, \quad \mathcal{B}(\kappa) \equiv \frac{\prod_{c \in c^*} \Gamma(\kappa_c)}{\Gamma(\sum_{c \in c^*} \kappa_c)},$$

each pdf $\pi_{ic}(\Theta_{ic}|\mathcal{S}_{ic;t})$ is conjugate to the factor $f_{ic}(d_{ic;t}|\psi_{ic;t}, \Theta_{ic})$.

Parameters Θ_{ic} , $i \in i^* \equiv \{1, \dots, d\}$, $c \in c^*$, of the individual parameterized factors are mutually conditionally independent, and also, independent of the component weights α . The component weights have Dirichlet distribution $Di_\alpha(\kappa)$ with support on the probabilistic simplex α^* .

Remarks 3

1. The considered form of the posterior distribution restricts the class of mixture factors $f_{ic}(d_{ic;t}|\psi_{ic;t}, \Theta_{ic})$ to those having conjugate pdf.
2. Details about Dirichlet distribution $Di_\alpha(\kappa_t)$ can be found for example in (Andrýšek 2004).
3. The research report (Andrýšek 2004) contains all details and proofs, which were omitted in this paper.

3.4 Notations related to mixtures

In the sequel, we use the following elements: $i \in i^* \equiv \{1, \dots, d\}$, $c \in c^*$

$$\text{Factor prediction } \mathcal{I}_{ic;t} = \int f_{ic}(d_{ic;t}|\psi_{ic;t}, \Theta_{ic}) \pi_{ic}(\Theta_{ic}|\mathcal{S}_{ic;t-1}) d\Theta_{ic}$$

$$\text{Component prediction } \beta_{c;t} = \prod_{i=1}^{\hat{\Psi}} \mathcal{I}_{ic;t} \quad (8)$$

$$\text{Estimate of component weight } \hat{\alpha}_{c;t} = \frac{\kappa_{c;t}}{\sum_{c \in c^*} \kappa_{c;t}} \quad (9)$$

$$\text{QB weight of data } w_{c;t} = \frac{\hat{\alpha}_{c;t-1} \beta_{c;t}}{\sum_{c=1}^{\hat{c}} \hat{\alpha}_{c;t-1} \beta_{c;t}} \quad (10)$$

$$\text{“Correct” estimate of factor parameters } \pi_{ic}(\Theta_{ic}|\mathcal{S}_{ic;t}^U) = \frac{f_{ic}(d_{ic;t}|\psi_{ic;t}, \Theta_{ic}) \pi_{ic}(\Theta_{ic}|\mathcal{S}_{ic;t-1})}{\mathcal{I}_{ic;t}} \quad (11)$$

Remarks 4

1. The assumption of conjugacy of $\pi_{ic}(\Theta_{ic}|\mathcal{S}_{ic;t-1})$ to the factor $f_{ic}(d_{ic;t}|\psi_{ic;t}, \Theta_{ic})$ implies that $\frac{f_{ic}(d_{ic;t}|\psi_{ic;t}, \Theta_{ic}) \pi_{ic}(\Theta_{ic}|\mathcal{S}_{ic;t-1})}{\mathcal{I}_{ic;t}}$ has the same functional form as $\pi_{ic}(\Theta_{ic}|\mathcal{S}_{ic;t-1})$, and thus we need to evaluate only the statistic $\mathcal{S}_{ic;t}^U$.
2. As values of $\mathcal{I}_{ic;t}$ can be very close to zero, it is numerically advantageous to evaluate the weights $w_{\bullet;t}$ using $\mathcal{L}_{ic;t} = \ln \mathcal{I}_{ic;t} \cdot \mathcal{L}_{ic;t}$ can be computed directly without evaluating $\mathcal{I}_{ic;t}$.



Algorithm 1 $w_{\bullet;t} = \text{EVAL_WEIGHT}(\mathcal{L}_{\bullet;t}, \kappa_{\bullet;t-1})$

1. For each component c evaluate $\mathcal{H}_{c;t} = \ln \kappa_{c;t-1} + \sum_i \mathcal{L}_{ic;t}$
2. $\bar{\mathcal{H}}_{\bullet;t} = H_{\bullet;t} - \max H_{\bullet;t}$
3. $w_{\bullet;t} = \frac{\exp(\bar{\mathcal{H}}_{\bullet;t})}{\sum_c \exp(\bar{\mathcal{H}}_{\bullet;t})}$

Remarks 5 $w_{c;t}$ evaluated in this algorithm is the same as defined in (10):

$$\begin{aligned} w_{c;t} &= \frac{\exp(\mathcal{H}_{c;t} - \max \mathcal{H}_{\bullet;t})}{\sum (\exp(\mathcal{H}_{c;t} - \max \mathcal{H}_{\bullet;t}))} = \frac{\exp(\mathcal{H}_{c;t}) \exp(\max \mathcal{H}_{\bullet;t})}{\exp(\max \mathcal{H}_{\bullet;t}) \sum \exp(\mathcal{H}_{c;t})} = \\ &= \frac{\kappa_{c;t-1} \beta_{c;t}}{\sum \kappa_{c;t-1} \beta_{c;t}} = \frac{\sum_{\kappa_{c;t-1}}^{\kappa_{c;t-1} \beta_{c;t}}}{\sum_{\kappa_{c;t-1}}^{\kappa_{c;t-1} \beta_{c;t}}} = \frac{\hat{\alpha}_{c;t-1} \beta_{c;t}}{\sum \hat{\alpha}_{c;t-1} \beta_{c;t}}. \end{aligned}$$

4 Problem formulation and general solution

In this Section, we apply the approximation from section 3.1.2 to the introduced mixture model (3). We seek

the statistic \mathcal{G}_t that minimizes $\mathcal{D} \left(\hat{\pi}(\Theta | \mathcal{G}_{t-1}, \overbrace{d_t, \phi_{t-1}}^{\equiv \Psi_t}) \parallel \pi(\Theta | \mathcal{G}_t) \right)$, where

$$\begin{aligned} \hat{\pi}(\Theta | \mathcal{G}_{t-1}, \Psi_t) &= \frac{f(d_t | \phi_{t-1}, \Theta) \pi(\Theta | \mathcal{G}_{t-1})}{\int f(d_t | \phi_{t-1}, \Theta) \pi(\Theta | \mathcal{G}_{t-1}) d\Theta} \\ \pi(\Theta | \mathcal{G}_{t-1}) &= Di_{\alpha}(\kappa_{t-1}) \prod_{i=1, c=1}^{\hat{\Psi}, \hat{c}} \pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t-1}) \\ f(d_t | \phi_{t-1}, \Theta) &= \sum_{c=1}^{\hat{c}} \alpha_c \prod_{i=1}^{\hat{\Psi}} f_{ic}(d_{ic;t} | \psi_{ic;t}, \Theta_{ic}). \end{aligned}$$

In this case, the statistic \mathcal{G}_t consist of vector κ_t (of \hat{c} elements) and multimatrix $\mathcal{S}_{\bullet;t}$ of type $(\hat{\Psi}, \hat{c})$.

The next proposition summarizes the form of $\hat{\pi}(\Theta | \mathcal{G}_{t-1}, \Psi_t)$.

Proposition 1

$$\hat{\pi}(\Theta | \mathcal{G}_{t-1}, \Psi_t) = \sum_{c=1}^{\hat{c}} w_{c;t} Di_{\alpha}(\kappa_{t-1} + \delta_{\bullet,c}) \prod_{\substack{j,r=1 \\ r \neq c}}^{\hat{\Psi}, \hat{d}} \pi_{jr}(\Theta_{jr} | \mathcal{S}_{jr;t-1}) \prod_{j=1}^{\hat{\Psi}} \pi_{jc}(\Theta_{jc} | \mathcal{S}_{jc;t}^U). \quad (12)$$

Proposition 2 (Minimization of KL distance) For $\mathcal{G}_t \equiv \{\mathcal{S}_{\bullet;t}, \kappa_t\}$ minimizing

$$\mathcal{D} \left(\hat{\pi}(\Theta | \mathcal{G}_{t-1}, \Psi_t) \parallel \pi(\Theta | \mathcal{G}_t) \right),$$

it holds:

$$\begin{aligned} \kappa_t &\in \text{Arg min}_{\kappa_t} \left[\sum_{c=1}^{\hat{c}} w_{c;t} \mathcal{D} \left(Di_{\alpha}(\kappa_{t-1} + \delta_{\bullet,c}) \parallel Di_{\alpha}(\kappa_t) \right) \right] \\ \mathcal{S}_{ic;t} &\in \text{Arg min}_{\mathcal{S}_{ic;t}} \left[(1 - w_{c;t}) \mathcal{D} \left(\pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t-1}) \parallel \pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t}) \right) + w_{c;t} \mathcal{D} \left(\pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t}^U) \parallel \pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t}) \right) \right]. \end{aligned} \quad (13)$$

Remarks 6 The previous proposition split the overall problem into two subproblems. The subproblem (13) can be solved in general, as presented in section 4.2. Solution of the second subproblem depends on the choice of the system model. Solution for the Normal models is presented in section 5.



4.1 General algorithm

Following the proposition 2 we sketch the general algorithm of one mixture estimation step. We naturally suppose that $\Psi_{ic;t}$ can be obtained from $d(t)$.

Algorithm 2

Inputs - $\kappa_{\bullet;t-1}, \mathcal{S}_{\bullet\bullet;t-1}, \Psi_{\bullet\bullet;t}$

Outputs - $\kappa_{\bullet;t}, \mathcal{S}_{\bullet\bullet;t}$

1. For each factor ic evaluate $\mathcal{L}_{ic;t} = \ln \mathcal{I}_{ic;t}$
2. $w_{\bullet;t} = \text{EVAL_WEIGHT}(\mathcal{L}_{\bullet\bullet;t}, \kappa_{\bullet;t-1})$ (Algorithm 1)
3. $\kappa_t \in \text{Arg min}_{\kappa_t > 0} \left[\sum_{c=1}^{\hat{c}} w_{c;t} \mathcal{D} \left(Di_{\alpha}(\kappa_{t-1} + \delta_{\bullet c}) \parallel Di_{\alpha}(\kappa_t) \right) \right]$
4. For each factor ic evaluate $\mathcal{S}_{ic;t}^U$ so that $\pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t}^U) = \frac{\pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t-1}) f_{ic}(d_{ic;t} | \psi_{ic;t}, \Theta_{ic})}{\mathcal{I}_{ic}}$
5. For each factor ic evaluate $\mathcal{S}_{ic;t} \in \text{Arg min}_{\mathcal{S}_{ic;t}} \left[(1 - w_{c;t}) \mathcal{D} \left(\pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t-1}) \parallel \pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t}) \right) + w_{c;t} \mathcal{D} \left(\pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t}^U) \parallel \pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t}) \right) \right]$

Steps 1,4,5 depends on the specific choice of the system model, step 2 is solved, and step 3 is discussed in the next section.

4.2 Minimization with respect to κ_t

The following proposition converts the problem of KL distance minimization of κ -part to minimization of an algebraic expression.

Proposition 3 (Minimization with respect to κ_t)

For κ_t minimizing

$$\sum_{c=1}^{\hat{c}} w_{c;t} \mathcal{D} \left(Di_{\alpha}(\kappa_{t-1} + \delta_{\bullet c}) \parallel Di_{\alpha}(\kappa_t) \right)$$

it holds

$$\kappa_{\bullet;t} \in \text{Arg min} \left\{ \sum_{c=1}^{\hat{c}} \left[\ln \left(\Gamma(\kappa_{c;t}) \right) - \kappa_{c;t} \xi_{c;t} \right] - \ln \left(\Gamma \left(\sum_{c=1}^{\hat{c}} \kappa_{c;t} \right) \right) \right\}$$

where

$$\xi_{c;t} = \left(\psi_0(\kappa_{c;t-1}) + \frac{w_{c;t}}{\kappa_{c;t-1}} - \psi_0 \left(\sum_{c=1}^{\hat{c}} \kappa_{c;t-1} + 1 \right) \right).$$

Proposition 3 yields the following algorithm.

Algorithm 3 $\kappa_{\bullet;t} = \text{NEW_KAPPA}(w_{\bullet;t}, \kappa_{\bullet;t-1})$

1. For each component c evaluate $\xi_{c;t} = \psi_0(\kappa_{c;t-1}) + \frac{w_{c;t}}{\kappa_{c;t-1}} - \psi_0 \left(\sum_{c=1}^{\hat{c}} \kappa_{c;t-1} + 1 \right)$
2. $\kappa_{\bullet;t} \in \text{Arg min} \left\{ \sum_{j=1}^{\hat{c}} \left[\ln \left(\Gamma(\kappa_{j;t}) \right) - \kappa_{j;t} \xi_{j;t} \right] - \ln \left(\Gamma \left(\sum_{c=1}^{\hat{c}} \kappa_{c;t} \right) \right) \right\}$

Remarks 7

1. Minimization of the term (13) can be simply approximated by changing $\mathcal{D} \left(Di_{\alpha}(\kappa_1) \parallel Di_{\alpha}(\kappa_2) \right)$ into square of the Euclidean norm $\|\kappa_1 - \kappa_2\|^2$. The problem is then transformed into minimization of $\min_x \sum_c w_c \|x - x_c\|^2$ which has explicit solution: $x = \sum_c w_c x_c$. Applied to our case it yields $\kappa_t = \kappa_{t-1} + w_t$, which is identical to the solution obtained using the quasi-Bayes algorithm (Kárný et al. 1998).

2. The minimization problem in step 2 must be solved numerically or by suitable approximation. For detailed solution of this problem see (Nenuil 2004).

We have completed all steps which can be done on this general level. In the next parts of the paper, we are dealing with the special case of the factors.

5 Application to normal factors

In this section, we assume the parameterized factor to be dynamic Gaussian pdf with parameters $\Theta_{ic} \equiv (\theta_{ic}, r_{ic})$, where θ_{ic} is so called vector of regression coefficients and r_{ic} is noise variance of the factor.

$$f_{ic}(d_{ic;t}|\psi_{ic;t}, \Theta_{ic}) = N_{d_{ic;t}}(\theta'_{ic}\psi_{ic;t}, r_{ic}) = \frac{1}{\sqrt{2\pi r_{ic}}} \exp\left(-\frac{(d_{ic;t} - \theta'_{ic}\psi_{ic;t})^2}{2r_{ic}}\right)$$

We don't need to introduce a shift in the mean value, because the regression vector can contain number 1. See Remarks 2. The shifting constant is then placed to the corresponding place of the vector of regression coefficients.

The prior conjugate to this model is the Gauss inverse Wishart pdf with parameters $\mathcal{S}_{ic;t} = (\nu_{ic;t}, V_{ic;t})$, where $\nu_{ic;t}$ is scalar count of degrees of freedom and $V_{ic;t}$ is so called extended information matrix (symmetric, positive definite, of type $(\hat{\Psi}, \hat{\Psi})$).

$$\pi_{ic}(\Theta_{ic}|\mathcal{S}_{ic;t}) = GiW_{\theta_{ic}, r_{ic}}(V_{ic;t}, \nu_{ic;t}) \propto r_{ic}^{-0.5(\nu_{ic;t} + \hat{\psi}_{ic;t} + 2)} \exp\left\{-\frac{1}{2r_{ic}} \text{tr}(V_{ic;t}[-1, \theta'_{ic}]'[-1, \theta'_{ic}])\right\}$$

Note that the matrix V_{ic} can be equivalently manipulated through its $L'DL$ decomposition (i.e. with lower triangular matrix L_{ic} and diagonal matrix D_{ic} which fulfills the relation $V_{ic} = L'_{ic}D_{ic}L_{ic}$). Next, the matrices L_{ic} and D_{ic} can be equivalently expressed via matrix C_{ic} , vector $\hat{\theta}_{ic}$ and scalar ${}^l d_{ic}$.

Because all three representations described above are equivalent, we will not formally distinguish between them. If V_{ic} is a statistic of GiW factor, under the terms $L_{ic}, D_{ic}, \theta_{ic}, C_{ic}, {}^l d_{ic}$ we automatically mean the parts of corresponding representation of the matrix V_{ic} .

Now, we specify the steps 1,4,5 in the general algorithm 2 for Normal factors.

5.1 Evaluating $\mathcal{I}_{ic;t}$

\mathcal{I}_{ic} is defined as

$$\mathcal{I}_{ic;t} = \int f_{ic}(d_{ic;t}|\psi_{ic;t}, \Theta_{ic}) \pi_{ic}(\Theta_{ic}|\mathcal{S}_{ic;t-1}) d\Theta_{ic} = \int N_{d_{ic;t}}(\theta'_{ic}\psi_{ic;t}, r_{ic}) GiW_{\theta_{ic}, r_{ic}}(V_{ic;t-1}, \nu_{ic;t-1}) d\theta_{ic} dr_{ic}.$$

$\mathcal{I}_{ic;t}$ for normal factors is evaluated as:

$$\mathcal{I}_{ic;t} = \frac{\Gamma(0.5(\nu_{ic;t-1} + 1)) [{}^l d_{ic;t-1}(1 + \zeta_{ic;t})]^{-0.5}}{\sqrt{\pi} \Gamma(0.5\nu_{ic;t-1}) \left(1 + \frac{\hat{e}_{ic;t}^2}{{}^l d_{ic;t-1}(1 + \zeta_{ic;t})}\right)^{0.5(\nu_{ic;t-1} + 1)}} \quad (14)$$

where

$$\begin{aligned} \hat{e}_{ic;t} &\equiv d_{ic;t} - \hat{\theta}'_{ic;t-1}\psi_{ic;t} \equiv \text{prediction error} \\ \zeta_{ic;t} &\equiv \psi'_{ic;t} C_{ic;t-1} \psi_{ic;t} \end{aligned}$$

Remarks 8 According to remarks 4, we need to evaluate $\mathcal{L}_{ic;t} = \ln \mathcal{I}_{ic;t}$. It can be done efficiently via the product form of (14). The following algorithm summarizes this task. Recall that $\Psi_{ic;t} = [d_{ic;t}, \psi_{ic;t}]$.

Algorithm 4 (evaluation of $\mathcal{L}_{ic;t}$) $\mathcal{L}_{ic;t} = \text{FACNORM}(C_{ic;t-1}, \hat{\theta}_{ic;t-1}, {}^l d_{ic;t-1}, \nu_{ic;t-1}, \Psi_{ic;t})$

1. Evaluate $\zeta_{ic;t} = \psi'_{ic;t} C_{ic;t-1} \psi_{ic;t}$



2. Evaluate $\hat{e}_{ic;t} \equiv d_{ic;t} - \hat{\theta}'_{ic;t-1} \psi_{ic;t}$

3. Evaluate

$$\begin{aligned} \mathcal{L}_{ic;t} = \ln \mathcal{I}_{ic;t} = & \ln \Gamma(0.5(\nu_{ic;t-1} + 1)) - \ln \Gamma(0.5\nu_{ic;t-1}) - 0.5 \ln({}^{\text{ld}}D_{ic;t-1}) - 0.5 \ln(1 + \zeta_{ic;t}) - \\ & - 0.5(\nu_{ic;t-1} + 1) \ln \left(1 + \frac{\hat{e}_{ic;t}^2}{{}^{\text{ld}}D_{ic;t-1}(1 + \zeta_{ic;t})} \right) - 0.5 \ln(\pi) \end{aligned}$$

Remarks 9 Function $\ln \Gamma$ can be evaluated without computing Γ first (Abramowitz & Stegun 1972).

5.2 Evaluating $\mathcal{S}_{ic;t}^U$

$\mathcal{S}_{ic;t}^U \equiv [V_{ic}^U, \nu_{ic}^U]$ can be evaluated in the following way:

$$\begin{aligned} V_{ic;t}^U &= V_{ic;t-1} + \Psi_{ic;t} \Psi'_{ic;t} \\ \nu_{ic;t}^U &= \nu_{ic;t-1} + 1 \end{aligned} \quad (15)$$

The relation (15) can be rewritten in terms of $C, \hat{\theta}, {}^{\text{ld}}D$:

$$\begin{aligned} C_{ic;t}^U &= C_{ic;t-1} - \frac{1}{1 + \zeta_{ic;t}} z_{ic;t} z'_{ic;t}, \quad \hat{\theta}_{ic;t}^U = \hat{\theta}_{ic;t-1} + \frac{\hat{e}_{ic;t}}{1 + \zeta_{ic;t}} z_{ic;t} \\ {}^{\text{ld}}D_{ic;t}^U &= {}^{\text{ld}}D_{ic;t-1} + \frac{\hat{e}_{ic;t}^2}{1 + \zeta_{ic;t}} \quad z_{ic;t} = C_{ic;t-1} \psi_{ic;t} \end{aligned}$$

5.3 Minimizing the KL distance

According to the proposition 2, we need to minimize

$$(1 - w_{c;t}) \mathcal{D} \left(\pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t-1}) \parallel \pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t}) \right) + w_{c;t} \mathcal{D} \left(\pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t}^U) \parallel \pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t}) \right)$$

for each factor i within the component c . The minimization can be done factor-wise (see alg. 2), thus we can simplify the notation by considering one particular factor.

$$\begin{aligned} \mathcal{S}_{ic;t-1} &\equiv (V_{ic;t-1}, \nu_{ic;t-1}) \rightarrow (V, \nu) \\ \mathcal{S}_{ic;t} &\equiv (V_{ic;t}, \nu_{ic;t}) \rightarrow (V^\spadesuit, \nu^\spadesuit) \\ \mathcal{S}_{ic;t}^U &\equiv (V_{ic;t-1}^U, \nu_{ic;t-1}^U) \rightarrow (V^U, \nu^U) \\ w_{c;t} &\rightarrow w \\ \psi_{ic;t}, \Psi_{ic;t}, d_t &\rightarrow \psi, \Psi, d \end{aligned}$$

Thus, we minimize

$$\min_{V^\spadesuit, \nu^\spadesuit} \left\{ (1 - w) \mathcal{D} \left(GiW_{\theta,r}(V, \nu) \parallel GiW_{\theta,r}(V^\spadesuit, \nu^\spadesuit) \right) + w \mathcal{D} \left(GiW_{\theta,r}(V^U, \nu^U) \parallel GiW_{\theta,r}(V^\spadesuit, \nu^\spadesuit) \right) \right\}. \quad (16)$$

Remarks 10

1. It can be proven that this minimization task can be divided into two independent algebraic subproblems. First of them is minimization on two dimensional space $(\nu^\spadesuit, {}^{\text{ld}}D^\spadesuit)$, the second is minimization on multi-dimensional space $(\hat{\theta}^\spadesuit, C^\spadesuit)$. Both subproblems are solved in the following sections.
2. If we approximate $\mathcal{D} \left(GiW_{\theta,r}(V, \nu) \parallel GiW_{\theta,r}(V^\spadesuit, \nu^\spadesuit) \right)$ with $\|V - V^\spadesuit\|^2 + \|\nu - \nu^\spadesuit\|^2$, we can quickly achieve the result $V^\spadesuit = V + w\Psi\Psi'$, $\nu^\spadesuit = \nu + w$, which is exactly the same as the quasi-Bayes update (Kárný et al. 1998).



5.3.1 Searching for ${}^{\text{L}^d}D^{\blacklozenge}$ and ν^{\blacklozenge}

Proposition 4 For ν^{\blacklozenge} , ${}^{\text{L}^d}D^{\blacklozenge}$ minimizing (16) it holds:

$$\frac{\nu^{\blacklozenge}}{{}^{\text{L}^d}D^{\blacklozenge}} = (1-w)\frac{\nu}{{}^{\text{L}^D}D} + w\frac{\nu^U}{{}^{\text{L}^D}D^U}$$

$$\ln(0.5\nu^{\blacklozenge}) - \psi_0(0.5\nu^{\blacklozenge}) = \Upsilon, \text{ where}$$

$$\Upsilon \equiv (1-w)\left(\psi_0(0.5\nu) - \ln({}^{\text{L}^D}D)\right) + w\left(\psi_0(0.5\nu^U) - \ln({}^{\text{L}^D}D^U)\right) - \ln\left(0.5(1-w)\frac{\nu}{{}^{\text{L}^D}D} + 0.5w\frac{\nu^U}{{}^{\text{L}^D}D^U}\right)$$

Straightforward application of Proposition 4 yields the following algorithm. Recall that $\Psi = [d, \psi]$.

Algorithm 5 (Updating ${}^{\text{L}^d}D$ and ν) $({}^{\text{L}^d}D^{\blacklozenge}, \nu^{\blacklozenge}) = \text{UPDATE_DFM}(w, C, \nu, \hat{\theta}, {}^{\text{L}^D}D, \Psi)$

1. $\hat{e} = d - \hat{\theta}'\psi$, $\zeta = \psi' C \psi$
2. $\nu^U = \nu + 1$, ${}^{\text{L}^D}D^U = {}^{\text{L}^D}D + \frac{\hat{e}^2}{1+\zeta}$
3. $X^S = (1-w)\frac{\nu}{{}^{\text{L}^D}D} + w\frac{\nu^U}{{}^{\text{L}^D}D^U}$
4. $\Upsilon = (1-w)\left[\psi_0(0.5\nu) - \ln({}^{\text{L}^D}D)\right] + w\left[\psi_0(0.5\nu^U) - \ln({}^{\text{L}^D}D^U)\right] - \ln(0.5X^S)$
5. Solve the equation for ν^{\blacklozenge} : $\ln(0.5\nu^{\blacklozenge}) - \psi_0(0.5\nu^{\blacklozenge}) = \Upsilon$
6. ${}^{\text{L}^d}D^{\blacklozenge} = \frac{\nu^{\blacklozenge}}{X^S}$

Remarks 11

1. Step 5 must be solved numerically or using some suitable approximation.
For detail description of the numerical solution and for proof of unicity of the solution see (Nenutil 2004).
2. The proof of existence of solution fulfilling $\nu^{\blacklozenge} > 0$, ${}^{\text{L}^d}D^{\blacklozenge} > 0$ can be found in (Andrýšek 2004).

5.3.2 Searching for $\hat{\theta}^{\blacklozenge}$ and C^{\blacklozenge}

Proposition 5 For $\hat{\theta}^{\blacklozenge}$ and C^{\blacklozenge} minimizing (16) it holds:

$$C^{\blacklozenge} = C + w_c z z' \tag{17}$$

$$\hat{\theta}^{\blacklozenge} = \hat{\theta} + w_\theta z \tag{18}$$

where

$$z = C\psi, \quad \hat{e} = d - \hat{\theta}'\psi, \quad \zeta = \psi' C \psi$$

$$w_c = \left[\frac{\hat{e}^2}{(1+\zeta)^2} \frac{X X^U}{X + X^U} - \frac{w}{1+\zeta} \right], \quad w_\theta = \left[\frac{\hat{e}}{1+\zeta} \frac{X^U}{X + X^U} \right]$$

$$X = (1-w)\frac{\nu}{{}^{\text{L}^D}D}, \quad X^U = w\frac{\nu^U}{{}^{\text{L}^D}D^U}$$

Algorithm 6 (Updating $\hat{\theta}$ and C) $(C^{\blacklozenge}, \hat{\theta}^{\blacklozenge}) = \text{UPDATE_C}(w, C, \nu, \hat{\theta}, {}^{\text{L}^D}D, \Psi)$

1. $\hat{e} = d - \hat{\theta}'\psi$, $\zeta = \psi' C \psi$
2. $\nu^U = \nu + 1$, ${}^{\text{L}^D}D^U = {}^{\text{L}^D}D + \frac{\hat{e}^2}{1+\zeta}$
3. $X = (1-w)\frac{\nu}{{}^{\text{L}^D}D}$, $X^U = w\frac{\nu^U}{{}^{\text{L}^D}D^U}$
4. $z = C\psi$
5. $C^{\blacklozenge} = C + \left[\frac{\hat{e}^2}{(1+\zeta)^2} \frac{X X^U}{X + X^U} - \frac{w}{1+\zeta} \right] z z'$

$$6. \hat{\theta}^\spadesuit = \hat{\theta} + \left[\frac{\hat{e}}{1+\zeta} \frac{X^U}{X^S} \right] z$$

Remarks 12

1. The keystone of the previous algorithm is step 5. The formula is very simple, but its iterative use can cause numerical troubles. Therefore, in practice, we always work with matrix C in its $L'DL$ decomposition. Numerical stability of this operation is discussed in (Nenutil 2004).
2. It can be proven that C^\spadesuit obtained by this algorithm is positive definite.

6 Resulting PB algorithm

In this Section, we summarize all the elaborated parts into one consistent algorithm.

Algorithm 7 (PB)

Inputs - $\kappa_{\bullet;t-1}, C_{\bullet;t-1}, \hat{\theta}_{\bullet;t-1}, {}^l d D_{\bullet;t-1}, \nu_{\bullet;t-1}, \Psi_{\bullet;t}$
Outputs - $\kappa_{\bullet;t}, C_{\bullet;t}, \hat{\theta}_{\bullet;t}, {}^l d D_{\bullet;t}, \nu_{\bullet;t}$

1. For each factor ic : $\mathcal{L}_{ic;t} = \text{FACNORM}(C_{ic;t-1}, \hat{\theta}_{ic;t-1}, {}^l d D_{ic;t-1}, \nu_{ic;t-1}, \Psi_{ic;t})$. (algorithm 4)
2. Evaluate $w_{\bullet;t} = \text{EVAL_WEIGHT}(\mathcal{L}_{\bullet;t}, \kappa_{\bullet;t-1})$. (algorithm 1)
3. Evaluate $\kappa_{\bullet;t} = \text{NEW_KAPPA}(w_{\bullet;t}, \kappa_{\bullet;t-1})$. (algorithm 3)
4. For each factor ic : $({}^l d D_{ic;t}, \nu_{ic;t}) = \text{UPDATE_DFM}(w_{ic;t}, C_{ic;t-1}, \nu_{ic;t-1}, \hat{\theta}_{ic;t-1}, {}^l d D_{ic;t-1}, \Psi_{ic;t})$. (algorithm 5)
5. For each factor ic : $(C_{ic;t}, \hat{\theta}_{ic;t}) = \text{UPDATE_C}(w_{ic;t}, C_{ic;t-1}, \nu_{ic;t-1}, \hat{\theta}_{ic;t-1}, {}^l d D_{ic;t-1}, \Psi_{ic;t})$. (algorithm 6)

7 Comparison of PB and QB algorithms

In this Section, we compare the performance of the PB algorithm with the performance of the standard QB algorithm. The QB algorithm has been used extensively in real-life applications (Kárný, et al. 2003), and it is proven to be reliable and computationally efficient. Therefore, we study differences of the PB algorithm from the QB in terms of numerical properties and quality of estimation. The algorithms are based on different objective criteria for which they are optimal. Therefore, comparison of their behaviour is presented in a subjective way: arguing what seem to be more "rational".

In order to compare the analytical properties, we review the QB algorithm. Then, we investigate the differences between the two algorithms from analytical and computational point of view. Those findings are supported by experimental results.

7.1 The Quasi-Bayes algorithm

The general QB algorithm uses the following rule(see (Kárný et al. 1998)):

$$\kappa_t = \kappa_{t-1} + w_t$$

$$\pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t}) \propto [f_{ic}(d_{ic;t} | \psi_{ic;t}, \Theta_{ic})]^{w_{ic;t}} \pi_{ic}(\Theta_{ic} | \mathcal{S}_{ic;t-1})$$

Let's mark the statistics corresponding to the QB algorithm by the subscript Q . Application of the general algorithm to the case with Normal factors yields:

$$V_Q = V + w \Psi \Psi', \nu_Q = \nu + w, \kappa_{Q,\bullet;t} = \kappa_{\bullet;t-1} + w_{\bullet;t} \quad (19)$$

We would receive exactly this result, if we approximate the KL distances in the PB algorithm with squares of euclidian norms of the parameter difference (see remarks 7 and 10).

For better comparison of the QB algorithm with the PB algorithm, we rewrite the relations (19) in terms of $C, \hat{\theta}, {}^l d D$:

$$C_Q = C + w_{QC}zz' \quad (20)$$

$$\hat{\theta}_Q = \hat{\theta} + w_{Q\theta}z, \quad {}^L D_Q = {}^L D + \frac{w\hat{e}^2}{1+w\zeta} \quad (21)$$

where

$$z = C\psi, \quad \hat{e} = d - \hat{\theta}'\psi, \quad \zeta = \psi' C \psi$$

$$w_{QC} = \frac{-w}{1+w\zeta}, \quad w_{Q\theta} = \frac{w\hat{e}}{1+w\zeta}$$

7.2 Analytical comparison

Nature of both algorithms allows us to divide the analytical investigation into two parts. In the first part, we investigate the update of factors. This part is discussed next. In the second part, computing of the new component weights can be studied, see (Nenutil 2004).

7.2.1 Differences of the algorithms

Note that the expressions for the QB update (20), (21) are very similar to the expressions for the PB update (17),(18). Hence, it suffice to investigate differences between the pairs $(\nu^\spadesuit, \nu_Q), ({}^L D^\spadesuit, {}^L D_Q), (w_C, w_{QC}), (w_\theta, w_{Q\theta})$. This involves observation of 4 scalar variables, no matter what is the full dimension of the parameters.

We illustrates differences in behavior on the following examples. Consider the following situations:

a) $\nu = 165.39, {}^L D = 9.77, \hat{e} = -0.0140, \zeta = 0.59$

b) $\nu = 102.82, {}^L D = 1.14, \hat{e} = -0.7386, \zeta = 1.20$

The figures 1 and 2 shows the parameters $(\nu^\spadesuit, \nu_Q), ({}^L D^\spadesuit, {}^L D_Q), (w_C, w_{QC}), (w_\theta, w_{Q\theta})$ as functions of $w \in \ll 0, 1 >$. The parameters related to the PB algorithm are plotted with the thick line. It is clear, that values obtained using PB equals to those of QB for $w = 0, w = 1$.

7.2.2 Behaviour of the PB algorithm

In this Section, we study two particular factors and evaluate marginal distributions of their updates provided by both algorithms. For better comparison, we will also show the marginal pdf of the correct Bayesian update (12) which is a mixture of two GiW factors.

Consider the GiW factor $\pi(\Theta|\mathcal{S}) = GiW_{\theta,r}(V, \nu)$ and denote the associated densities as follows:

| | |
|---|---|
| trial update $\pi(\Theta \mathcal{S}^U) = GiW_{\theta,r}(V^U, \nu^U)$ | $V^U = V + \Psi\Psi', \nu^U = \nu + 1$ |
| QB update $\pi(\Theta \mathcal{S}_Q) = GiW_{\theta,r}(V_Q, \nu_Q)$ | $V_Q = V + w\Psi\Psi', \nu_Q = \nu + w$ |
| PB update $\pi(\Theta \mathcal{S}^\spadesuit) = GiW_{\theta,r}(V^\spadesuit, \nu^\spadesuit)$ | result of the algorithms 4 and 5 |
| correct update $\hat{\pi}(\Theta) = (1-w)\pi(\Theta \mathcal{S}) + w\pi(\Theta \mathcal{S}^U)$ | |

Consider the statistics V, ν of the GiW factor, updating weights w and actual data vectors of the factor Ψ , to be:

| | |
|--|--|
| a) | b) |
| $V = \begin{pmatrix} 1.16 & 0.12 \\ 0.12 & 0.83 \end{pmatrix}$ | $V = \begin{pmatrix} 1.96 & -1.47 \\ -1.47 & 6.07 \end{pmatrix}$ |
| $\nu = 102.82$ | $\nu = 108.06$ |
| $\Psi = (-0.59 \ 1)'$ | $\Psi = (-0.79 \ 1)'$ |
| $w = 0.43$ | $w = 0.39$ |

The figures 3 and 4 shows marginal pdfs of all discussed densities for both cases. From visual inspection of these figures, we can conclude that the PB algorithm can provide results significantly different from those of the QB algorithm. We also consider behavior of the PB algorithms as reasonable.

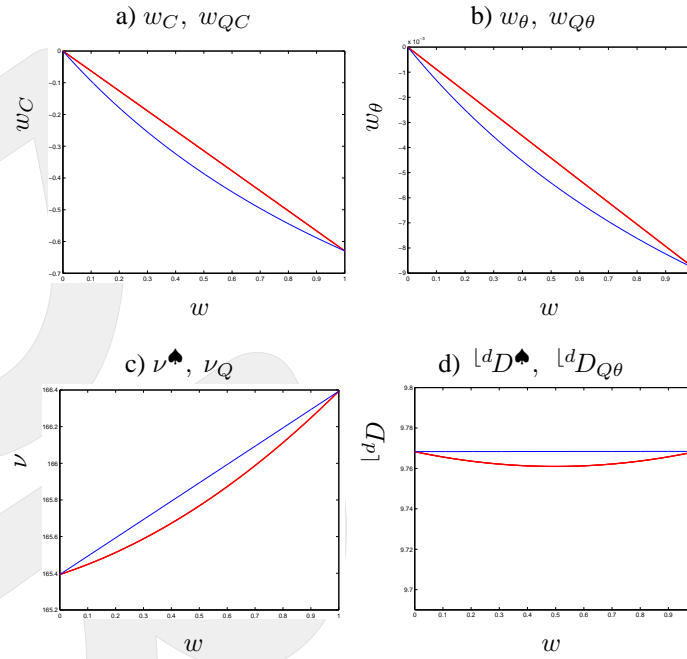


Figure 1. Similar behavior of the QB and PB algorithms for case a)

The figure shows the parameters (ν^*, ν_Q) , $(l^d D^*, l^d D_Q)$, (w_C, w_{QC}) , $(w_\theta, w_{Q\theta})$ as the functions of $w \in \langle 0, 1 \rangle$ for the case a) $\nu = 165.39$, $l^D D = 9.77$, $\hat{e} = -0.0140$, $\zeta = 0.59$. The parameters related to PB algorithm are plotted with the thick line. In this case the difference between the QB and PB algorithms is rather small.

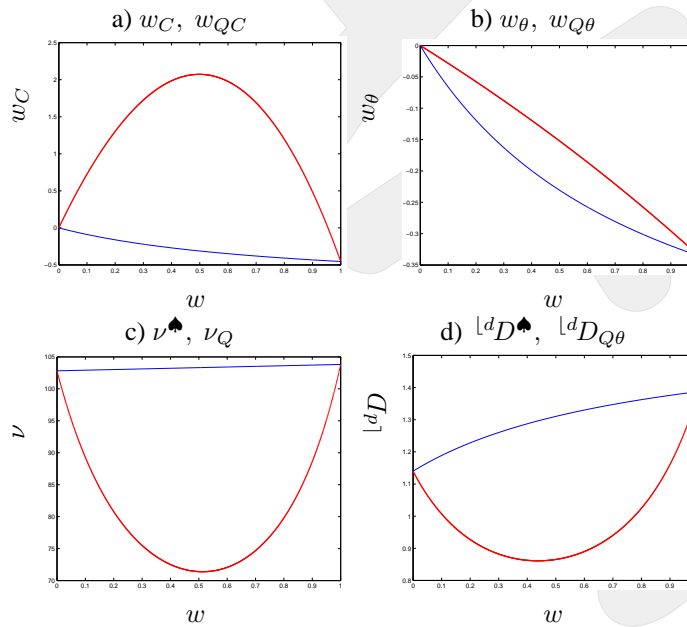


Figure 2. Different behavior of the QB and PB algorithms for case b)

The figure shows the parameters (ν^*, ν_Q) , $(l^d D^*, l^d D_Q)$, (w_C, w_{QC}) , $(w_\theta, w_{Q\theta})$ as the functions of $w \in \langle 0, 1 \rangle$ for the case b) $\nu = 102.82$, $l^D D = 1.14$, $\hat{e} = -0.7386$, $\zeta = 1.20$. The parameters related to the PB algorithm are plotted with the thick line. In this case, the difference between the QB and PB algorithms is significant.

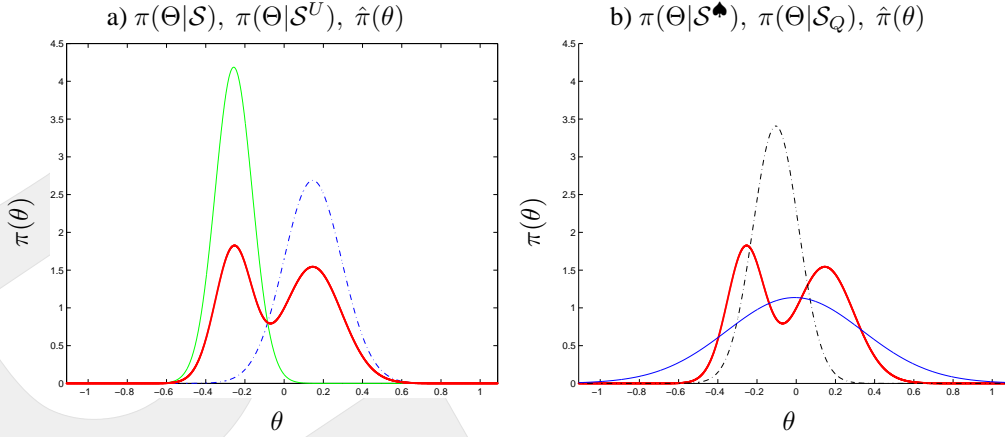


Figure 3. Marginal pdfs of the QB and PB updates for the case a)

The left part shows original factor (dashdot), its trial update (dotted) and the correct Bayesian update (thick), i.e. the mixture of the two mentioned factors. The right part shows how the QB update (dashdot) and the PB update (solid) approximates the correct Bayesian update (thick). It can be seen that the PB update is in this case flatter than the QB update which concentrates on smaller interval.

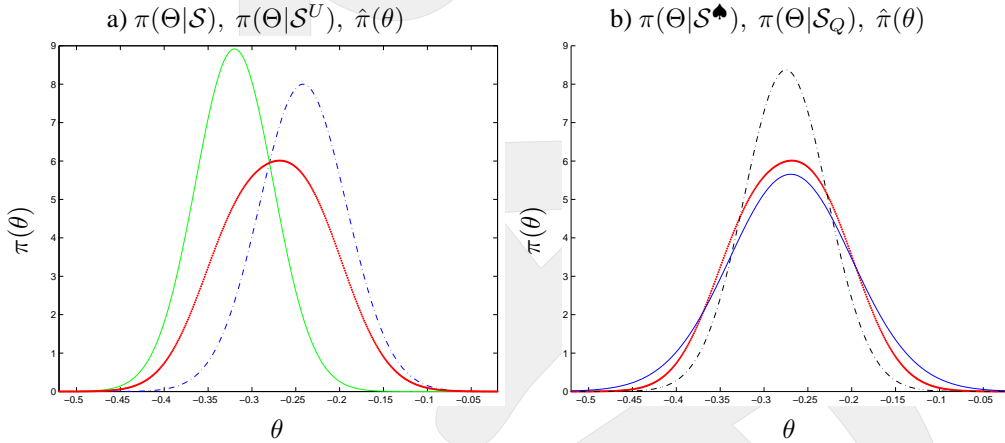


Figure 4. Marginal pdfs of the QB and PB updates for the case b)

The left part shows original factor (dashdot), its trial update (dotted) and the correct Bayesian update (thick), i.e. the mixture of the two mentioned factors. The right part shows how the QB update (dashdot) and the PB update (solid) approximates the correct Bayesian update (thick). It can be seen that the PB update in this case better approximates the correct pdf.

7.3 Experimental comparison

Intensive tests consisting of 1396 data sets were done. Data used for this test represent various types of systems (static, dynamic, multidimensional) and are part of standard testing procedure of new algorithms. As a quality measure, we used the likelihood (Kárný et al. 2003) of the estimated model. For each set, we evaluated a criterion h which is the difference between the likelihood obtained by the PB algorithm and the QB algorithm. (i.e $h > 0$ if the PB algorithm was better.) The table 1 shows the results. Mean value of h over all sets is 6.18.

7.4 Comparing of computational complexity

We compare all 5 steps of the PB algorithm (algorithm 7).

1. This step is needed in both algorithms.
2. This step is needed in both algorithms.

| condition | number of sets | percentage |
|---------------------|----------------|------------|
| $h > 0$ | 1125 | 80.6% |
| $h < 0$ | 271 | 19.4% |
| $\text{abs}(h) < 2$ | 1126 | 80.6% |
| $h > 2$ | 251 | 18.0% |
| $h < -2$ | 19 | 1.4% |

Table 1. Results of experimental comparison

The table shows some conditions for h and number of sets fulfilling each condition.

3. We have to find minimizer of a convex function with \hat{c} variables. There exist a good approximation of the starting point for iterative numerical algorithm, which warrants quick solution of this task (Nenutil 2004).
4. Solution of one-dimensional nonlinear equation must be found. However, a good approximation which always leads to solving the equation in a few steps was found (Nenutil 2004).
5. This step has the same complexity in both algorithms.

Addressing the previous considerations, we conclude that computational cost of numerical evaluation of the PB algorithm is comparable to the computational cost associated with the QB algorithms. Detailed case study of the computational costs of both algorithms can be found in (Nedoma & Andryšek 2004).

8 Conclusions

This work describes a novel and efficient algorithm for recursive estimation of finite probabilistic mixture. The algorithm has the potential of providing more accurate results than the well-established quasi-Bayes estimator. This improvement is important as mixtures represent a universal approximating tool for modelling of non-linear stochastic systems. Therefore, mixture models can be used to address complex control and decision-making problems in changing environments, such as multiple-participants decision making. Each participant (or group of participants) can be modelled by a component of the overall mixture model. All subsequent decision-making task can be easily formalized within the consistent formal framework of probabilistic mixture models. We believe, that the algorithms presented in this paper will be an important part of this framework.

9 Acknowledgments

This work was supported by AV ČR S1075351, GA ČR 102/03/0049, GA ČR 102/01/0608

References

- M. Abramowitz & I. Stegun (1972). *Handbook of mathematical functions*. Dover Publications, Inc., New York.
- J. Andryšek (2004). ‘Projection Based Estimation of Dynamic Probabilistic Mixtures’. Tech. Rep. 2098, ÚTIA AV ČR, Praha.
- L. Berc & M. Kárný (1997). ‘Identification of reality in Bayesian context’. In K. Warwick & M. Kárný (eds.), *Computer-Intensive Methods in Control and Signal Processing: Curse of Dimensionality*, pp. 181–193. Birkhäuser.
- L. He & M. Kárný (2003). ‘Bayesian Modelling and Estimation of ARMAX Model and Its Finite Mixtures’. Tech. Rep. 2080, ÚTIA AV ČR, Praha.
- M. Kárný, et al. (2003). ‘Productool background – theory, algorithms and software’. Tech. rep., UTIA AV CR. Draft of the report, 401 pp.
- M. Kárný, et al. (2003). ‘Mixture-based adaptive probabilistic control’. *International Journal of Adaptive Control and Signal Processing* **17**(2):119–132.
- M. Kárný, et al. (1998). ‘Quasi-Bayes estimation applied to normal mixture’. In J. Rojíček, M. Valečková,



- M. Kárný, & K. Warwick (eds.), *Preprints of the 3rd European IEEE Workshop on Computer-Intensive Methods in Control and Data Processing*, pp. 77–82, Praha. ÚTIA AV ČR.
- S. Kullback & R. Leibler (1951). 'On information and sufficiency'. *Annals of Mathematical Statistics* **22**:79–87.
- P. Nedoma & J. Andřýsek (2004). 'Benchmarks in Mixture Processing'. In M. Kárný, J. Kracík, & J. Andřýsek (eds.), *Proceedings of the International Workshop on Multiple Participant Decision Making 2004*, Praha. ÚTIA AV ČR.
- P. Nenutil (2004). 'Přibližné průběžné odhadování dynamických směsových modelů'. Tech. rep., ÚTIA AV ČR, Praha.
- V. Peterka (1981). 'Bayesian system identification'. In P. Eykhoff (ed.), *Trends and Progress in System Identification*, pp. 239–304. Pergamon Press, Oxford.
- S. J. Roberts & W. D. Penny (2002). 'Variational Bayes for Generalized Autoregressive Models'. *IEEE Transactions on Signal Processing* **50**(9):2245–2257.
- D. Titterington, et al. (1985). *Statistical Analysis of Finite Mixtures*. John Wiley & Sons, Chichester, New York, Brisbane, Toronto, Singapore. ISBN 0 471 90763 4.

