# Fusion of Probabilistic Unreliable Indirect Information into Estimation Serving to Decision Making

**Miroslav Kárný** · **František Hůla**

**Abstract** Bayesian decision making (DM) quantifies information by the probability density (pd) of treated variables. Gradual accumulation of information during acting increases the DM quality reachable by an agent exploiting it. The inspected accumulation way uses a parametric model forecasting observable DM outcomes and updates the posterior pd of its unknown parameter. In the thought multi-agent case, a neighbouring agent, moreover, provides a privately-designed pd forecasting the same observation. This pd may notably enrich the information of the focal agent. Bayes' rule is a unique deductive tool for a lossless compression of the information brought by the observations. It does not suit to processing of the forecasting pd. The paper extends solutions of this case. It: ▹ refines the Bayes'-rule-like use of the neighbour's forecasting pd ▹ deductively complements former solutions so that the learnable neighbour's reliability can be taken into account ▹ specialises the result to the exponential family, which shows the high potential of this information processing ▹ cares about exploiting population statistics.

## 1 INTRODUCTION

*Addressed problem:* This work contributes to the primary purpose of information fusion as stated in [12]: "...information fusion is the study of efficient methods for automatically or semi-automatically transforming information from different sources and different points in time into a representation that provides effective support for human or automated decision making".

The DM entity supported by this paper is called an agent and referred by "it" unless we want to stress the human involved. It fuses information about the relation of its
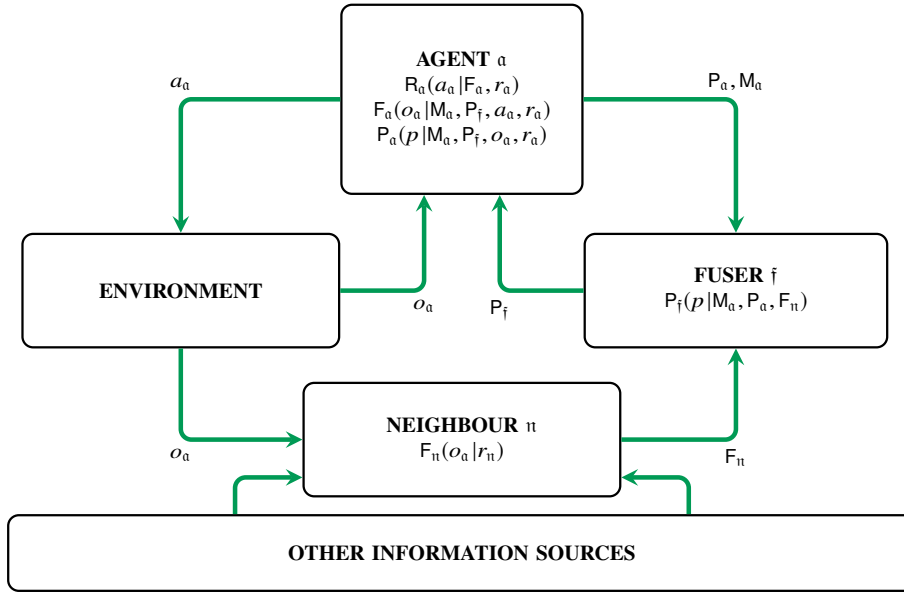
The Czech Academy of Sciences, Institute of Information Theory and Automation, 182 00 Prague 8, Czech Republic, E-mail: {school,hula}@utia.cas.cz

Fig. 1: The agent, $\mathfrak{a}$, learns a parameter, $p \in \boldsymbol{p}$, of the parametric model, $\mathsf{M}_\mathfrak{a}$, of the observation, $o_\mathfrak{a} \in \boldsymbol{o_\mathfrak{a}}$, for the given regressor, $r_\mathfrak{a} \in \boldsymbol{r_\mathfrak{a}}$. It evaluates the forecasting pd, $\mathsf{F}_\mathfrak{a}$, used for opting the action, $a_\mathfrak{a} \in \boldsymbol{a_\mathfrak{a}}$, by the decision rule, $\mathsf{R}_\mathfrak{a}$. $\mathsf{F}_\mathfrak{a}$ is expectation of $\mathsf{M}_\mathfrak{a}$ given by the posterior pd, $\mathsf{P}_\mathfrak{f}$, of $p \in \boldsymbol{p}$. The fuser, $\mathfrak{f}$, creates $\mathsf{P}_\mathfrak{f}$ from the model, $\mathsf{M}_\mathfrak{a}$, the agent's posterior pd, $\mathsf{P}_\mathfrak{a}$, updated by Bayes' rule, and from the forecasting pd, $\mathsf{F}_\mathfrak{n}$, of the observation, $o_\mathfrak{a} \in \boldsymbol{o_\mathfrak{a}}$. $\mathsf{F}_\mathfrak{n}$ is offered by the $\mathfrak{a}$'s neighbour, $\mathfrak{n}$, that gains $\mathsf{F}_\mathfrak{n}$ by using the common $o_\mathfrak{a}$, its regressor, $r_\mathfrak{n} \in \boldsymbol{r_\mathfrak{n}}$, and other private resources.

optional actions to observations. Figure 1 depicts the inspected fusion scenario (motivated below) with a detailed description of involved blocks, mappings and variables. It serves as a reference point and its individual parts are gradually explained.

The supported agent relies on a parametric model of the forecasted observation and on the posterior probability density[1] (pd) expressing the agent's information about the unknown model parameter. Bayes' rule [6] updates the posterior pd when the observed data record is available. The record consists of the observation, the applied action, and the used regressor made of past observations and actions. The agent also exploits a pd offered by its neighbour that forecasts the same observation. The neighbour may freely construct the offered forecasting pd in its private way and use its own information resources. The paper *proposes a fusing algorithm (fuser) that deductively exploits the offered forecasting pd*. The use of parametric models from the exponential family (EF [5]) provides a directly applicable fuser.

*Probabilistic modelling and complexity:* Among descriptions of uncertain information [22], the DM axiomatisation [41] makes us stay within the Bayesian paradigm,

---

[1] The existence of regular probability densities of inspected probabilistic measures with respect to Lebesgue's or counting measures is assumed [39].

where all variables are taken as random. For them, conditioning by the available information emerges as the proper information processing. Stochastic filtering [15] and Bayesian estimation [36] are common conditioning ways.

Complexity strongly limits the use of this prescriptive information processing. The complexity means an excessive need of resources for: ▹ the design of involved models ▹ algorithmic processing ▹ storing results; ▹ obtaining a sufficient amount of informative data; ▹ information-transmission, etc.

Complexity—irrespectively of its nature—makes agents to favourise simple models as linear regressions. The inevitably approximate expression of relations between the modelled variables does little harm when such models are learnt on-line. The local modelling then often suffices as the modelled dependencies are mostly smooth. The quality of such information processing can be notably enhanced by fusing the local statistic values with those gained in a preliminary off-line processing. On-line data processing and the statistics' fusion have to always be done within limited time and data budgets.

Generally, the complexity can only be counteracted by distributed—temporally, information-resources-wise or agents-wise—information processing. It surely needs an information fusion at some solution stage. The paper considers a specific but a widely applicable scenario outlined in Abstract and shown in Figure 1.

*Applications motivating the inspected scenario (the cases we met):* The probabilistic advisory system [18] models a closed loop, formed by a human agent and a world-part she cares about, by dynamic mixtures of pds. The system advises her on the actions promising the best forecast observation. Applications of this concept are described in [37]. The main one concerned rolling mill operators. The advices to them had to be refreshed on a cheap industrial hardware with clock rates of about 10 Hz while using a dynamic 2nd order model of more than ten variables. These circumstances limited the possibility to exploit databases containing around $10^6$ records. Thus, it makes sense to create a forecasting pd off-line and combine it with a simple parametric model updated on-line.

The same advisory system was tested for advising a therapeutic dose of radioactive iodine used for curing thyroid-gland cancer. The dose is to be personalised by using 3-4 measurements available for a specific patient. The patient-centric treatment must also exploit population statistic and physician forecasts. Obviously, any system supporting personalised medicine, as [21], should be able to combine sparse patient data with population statistics and experts' wisdom.

Recommendation systems in e-commerce represent another common use case [40]. We dealt with a version close to the advising to rolling mill operators. A customer is part of a retailer's loyalty programme that records the history of her shopping. Typically, her record reflects 3 to 10 shoppings per year in a specific category of goods. Lower units of years are recorded. The population data, describing past shopping of *all* other customers, can be processed off-line. The personalised on-line recommendations during a visit of the e-shop have to rely on the said, very sparse, personal data records in conjunction with the population model, all in tens of milliseconds.

*Other samples fitting the inspected scenario:* The fusion of filters' outcomes has many solutions and applications, e.g. in navigation [4] or robotics [3]. It has also been elaborated in the vein of this paper [9]. Fusions fitting our scenario that yield a soft cooperation of adaptive controllers exist, too [20].

*Layout and conventions:* Section 2 formalises the thought fusion problem. Section 3 solves it. Section 4 applies the solution to the EF of parametric models. Section 5 extends the gained results to handle the use of population statistics. Section 6 numerically illustrates the theory. Section 7 complements references made on the fly and contained in survey papers [31, 46] on information fusion by a few comments on the works related to ours.

*Throughout:*
✓ boldface $\boldsymbol{x}$ and **P** denote sets of possible $x$'s and P's;
✓ random variables, their values and realisations are formally undistinguished;
✓ models are pds described by san serif fonts similarly to other mappings;
✓ indices $\mathfrak{a}, \mathfrak{n}, \mathfrak{f}$ relate objects to their providers and (mostly implicitly) refer to the informations, $I_\mathfrak{a}$, $I_\mathfrak{n}$ and $I_\mathfrak{f}$, used by the agent, $\mathfrak{a}$, the neighbour, $\mathfrak{n}$, and the fuser, $\mathfrak{f}$;
✓ informations $I_\mathfrak{a}, I_\mathfrak{n}, I_\mathfrak{f}$ contain (often implicitly) the agent's action, $a_\mathfrak{a}$, and the regressors, $r_\mathfrak{a}, r_\mathfrak{n}, r_\mathfrak{f}$;
$\equiv$ stresses equality by the assignment;
$\propto$ is equality up to the normalising factor;
$^{\mathsf{opt}}$ indicates the optimal use of the richest information, $I_\mathfrak{f}$.

## 2 INFORMATION PROCESSING SCENARIO

A triple of information-handling agents is inspected, see Figure 1. The key one is the focal agent, $\mathfrak{a}$, to which the information processing serves. Its neighbour, $\mathfrak{n}$, serves as an additional information source. The fuser, $\mathfrak{f}$, represents the algorithm proposed here to help $\mathfrak{a}$ in using the probabilistic information provided by $\mathfrak{n}$.

*The agent,* $\mathfrak{a}$, uses its domain knowledge for selecting the parametric model, $\mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p)$ $\equiv \mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p, a_\mathfrak{a}, r_\mathfrak{a})$, of an observation, $o_\mathfrak{a} \in \boldsymbol{o}_\mathfrak{a}$. An unknown parameter, $p \in \boldsymbol{p}$, an agent-opted action, $a_\mathfrak{a}$, and the regressor, $r_\mathfrak{a}$, condition the model together with the implicitly present $a_\mathfrak{a}, r_\mathfrak{a}$. The agent quantifies its current information about the unknown parameter, $p \in \boldsymbol{p}$, by the posterior pd, $\mathsf{P}_\mathfrak{a}(p)$. Even when a new piece of information is obtained, $\mathfrak{a}$ *preserves the parametric model*, $\mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p)$, but replaces its posterior pd, $\mathsf{P}_\mathfrak{a}(p)$, by a new one, $\mathsf{P}_\mathfrak{f}(p)$, chosen from a given set of feasible posterior pds, **P**. Thus, the information, $I_\mathfrak{a}$, inherent to $\mathfrak{a}$ is

$$I_\mathfrak{a} \equiv (\mathsf{M}_\mathfrak{a}, \mathsf{P}_\mathfrak{a}, r_\mathfrak{a}, \mathbf{P}). \tag{1}$$

The assumption that $\mathfrak{a}$ is unwilling to change the parametric model, $\mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p)$, prevents induced, complexity-bringing, changes of the further information accumulation and of the decision-rule design. The replacement of $\mathsf{P}_\mathfrak{a}$ by $\mathsf{P}_\mathfrak{f} \in \mathbf{P}$ causes no such problem.

*The neighbour,* $\mathfrak{n}$, provides the agent, $\mathfrak{a}$, with extra information about the observation, $o_\mathfrak{a} \in \boldsymbol{o}_\mathfrak{a}$. $\mathfrak{n}$ offers the forecasting pd, $\mathsf{F}_\mathfrak{n}(o_\mathfrak{a}) \equiv \mathsf{F}_\mathfrak{n}(o_\mathfrak{a}|r_\mathfrak{n})$. Thus, the neighbour, $\mathfrak{n}$, transfers

a part of the information, $I_\mathfrak{n}$, via the fuser, $\mathfrak{f}$, (see below) to the agent, $\mathfrak{a}$. The neighbour possesses the information

$$I_\mathfrak{n} \equiv (\mathsf{F}_\mathfrak{n}, r_\mathfrak{n}, \dots). \tag{2}$$

The neighbour models the agent's observation in a way that may completely differ from that used by the agent. $\mathfrak{n}$ operates on its regressor, $r_\mathfrak{n}$, which is possibly unknown to $\mathfrak{a}$. The ellipsis in (2) indicates that the neighbour may have access to a richer or independent database. The neighbour, $\mathfrak{n}$, may have no clue about the existence of the environment model, $\mathsf{M}_\mathfrak{a}$, and thus about its unknown parameter, $p \in \boldsymbol{p}$.

The processing assumes that the forecasting pd, $\mathsf{F}_\mathfrak{n}$, models the observation, $o_\mathfrak{a}$ corresponding to the realised agent's action, $a_\mathfrak{a}$, and the regressor, $r_\mathfrak{a}$. Verifying this assumption is usually simple. For instance, when the correspondence of $a_\mathfrak{a}, r_\mathfrak{a}$ with $r_\mathfrak{n}$ follows from their simultaneous observation of the modelled environment.

*The fuser, $\mathfrak{f}$, is the algorithm designed in the paper.* It opts a posterior pd, $\mathsf{P}_\mathfrak{f} \in \mathbf{P}$, according to the fused information, $I_\mathfrak{f}$,

$$I_\mathfrak{f} \equiv (\mathsf{M}_\mathfrak{a}, \mathsf{P}_\mathfrak{a}, \mathsf{F}_\mathfrak{n}, r_\mathfrak{f}, \mathbf{P}). \tag{3}$$

$I_\mathfrak{f}$ consists of the environment model, $\mathsf{M}_\mathfrak{a}$, the agent's posterior pd, $\mathsf{P}_\mathfrak{a}$, of its parameter, $p \in \boldsymbol{p}$, the neighbour's forecasting pd, $\mathsf{F}_\mathfrak{n}$, the regressor, $r_\mathfrak{f}$, that unites $r_\mathfrak{a}$ and $r_\mathfrak{n}$, and the set of feasible posterior pds, $\mathbf{P}$.

The information, $I_\mathfrak{f}$ (3), does not determine the desired posterior pd, $\mathsf{P}_\mathfrak{f} \in \mathbf{P}$, uniquely. Thus, the fuser is to select such an improved posterior pd, $\mathsf{P}_\mathfrak{f} \in \mathbf{P}$, using the limited information, $I_\mathfrak{f}$. The choice is done by a static randomised DM strategy with a fusing rule[2]

$$\mathsf{S}_\mathfrak{f}(\mathsf{P}_\mathfrak{f}) \equiv \mathsf{S}_\mathfrak{f}(\mathsf{P}_\mathfrak{f} | I_\mathfrak{f}) = \mathsf{S}_\mathfrak{f}(\mathsf{P}_\mathfrak{f} | \mathsf{M}_\mathfrak{a}, \mathsf{P}_\mathfrak{a}, \mathsf{F}_\mathfrak{n}, r_\mathfrak{f}, \mathbf{P}), \;\; \mathsf{P}_\mathfrak{f} \in \mathbf{P}. \tag{4}$$

*The paper designs the optimal fusing rule, $\mathsf{S}^{\mathrm{opt}}$, of the form (4) for the used information-processing scenario. The design exploits a deductive, axiomatically justified, methodology.*


## 3 PROBLEM FORMULATION AND SOLUTION

The agent, $\mathfrak{a}$, allows the fuser, $\mathfrak{f}$, to modify its posterior pd, $\mathsf{P}_\mathfrak{a}$ – describing the parameter of the model, $\mathsf{M}_\mathfrak{a}$ – to the posterior pd, $\mathsf{P}_\mathfrak{f}$, offered by the fuser. The agent only allows the posterior pds from the set of feasible pds, $\mathbf{P}$. The agent is open to the replacement of $\mathsf{P}_\mathfrak{a}$ by a posterior pd, $\mathsf{P}_\mathfrak{f}$, from $\mathbf{P}$. It may even provide its rule, $\mathsf{S}_\mathfrak{a}$, which randomly selects a posterior pd from the set $\mathbf{P}$ using its limited information, $I_\mathfrak{a}$. This rule serves as a prior ansatz of the constructed fusing rule, $\mathsf{S}_\mathfrak{f}$, that exploits the richer fuser's information, $I_\mathfrak{f}$. The classical work [43] provides axiomatics that recommends to choose the optimal fusing rule $\mathsf{S}^{\mathrm{opt}}$ as the minimiser of the Kullback-Leibler

---

[2] Our manipulations assume discrete-valued modelled variables. The uncertain pds acting on them are probabilistic vectors and their distributions are then modelled without technicalities of the measure theory. The found solution is valid without this assumption.

divergence[3] (KLD, [27]) of $S_\mathfrak{f} \in \mathbf{S}_\mathfrak{f}$ to $S_\mathfrak{a}$. The work [19] extended this *minimum KLD principle* so that it copes with quite general sets of fusing rules, $\mathbf{S}_\mathfrak{f}$, over which the optimisation runs. The optimal fusing rule is thus

$$S^{opt} \in \text{Arg} \min_{S_\mathfrak{f} \in \mathbf{S}_\mathfrak{f}} D(S_\mathfrak{f} || S_\mathfrak{a}) \equiv \text{Arg} \min_{S_\mathfrak{f} \in \mathbf{S}_\mathfrak{f}} \int_{\mathbf{P}} S_\mathfrak{f}(P_\mathfrak{f}) \ln\left(\frac{S_\mathfrak{f}(P_\mathfrak{f})}{S_\mathfrak{a}(P_\mathfrak{f})}\right) dP_\mathfrak{f}. \tag{5}$$

The specification of the set, $\mathbf{S}_\mathfrak{f}$, containing suitable fusing rules, $S_\mathfrak{f}$ (4), determines the $S^{opt}$-choice (5). The proper choice has to reflect the richer information, $I_\mathfrak{f}$, available to the fuser compared to the agent's information, $I_\mathfrak{a}$. The next construction of $\mathbf{S}_\mathfrak{f}$ forms the core of our solution.

The used scenario, Section 2, implies that the parametric model[4], $M_\mathfrak{a}(o_\mathfrak{a}|p)$, is given. The optional posterior pds, $P_\mathfrak{f}(p)$, belong to the set of feasible pds, $\mathbf{P}$. The chain rule for pds [36] provides the joint pd of the observation, $o_\mathfrak{a}$, and the parameter, $p$,

$$J_{\mathfrak{a}\mathfrak{f}}(o_\mathfrak{a}, p|P_\mathfrak{f}) = M_\mathfrak{a}(o_\mathfrak{a}|p)P_\mathfrak{f}(p), \quad (o_\mathfrak{a}, p) \in \boldsymbol{o}_\mathfrak{a} \times \boldsymbol{p}. \tag{6}$$

The opted posterior pd, $P_\mathfrak{f}(p)$, and the neighbour's forecasting pd, $F_\mathfrak{n}(o_\mathfrak{a})$, also characterise a joint pd, $J_{\mathfrak{n}\mathfrak{f}}(o_\mathfrak{a}, p|P_\mathfrak{f})$, reflecting the neighbour's contribution to the fuser information, $I_\mathfrak{f}$. It is known that marginal pds, $F_\mathfrak{n}$ and $P_\mathfrak{f}$, do not determine the joint pd, $J_{\mathfrak{n}\mathfrak{f}}$, uniquely [33]. The conditions of the solved tasks make, however, the next product choice unambiguous

$$J_{\mathfrak{n}\mathfrak{f}}(o_\mathfrak{a}, p|P_\mathfrak{f}) = F_\mathfrak{n}(o_\mathfrak{a})P_\mathfrak{f}(p), \quad (o_\mathfrak{a}, p) \in \boldsymbol{o}_\mathfrak{a} \times \boldsymbol{p}. \tag{7}$$

It reflects that the neighbour forecasts the observation without information about the agent's parametric model, $M_\mathfrak{a}$. Moreover, the fuser models the parameter by the pd, $P_\mathfrak{f} \in \mathbf{P}$, chosen *before* seeing the observation, $o_\mathfrak{a}$. This makes $o_\mathfrak{a}, p$ independent.

The joint pds (of $o_\mathfrak{a}, p$) $J_{\mathfrak{n}\mathfrak{f}}$ and $J_{\mathfrak{a}\mathfrak{f}}$, comprise the more rich information, $I_\mathfrak{f}$ (3), than the joint pd, $J_\mathfrak{a}$, that only uses the agent's information, $I_\mathfrak{a}$ (1),

$$J_\mathfrak{a}(o_\mathfrak{a}, p|P_\mathfrak{a}) = M_\mathfrak{a}(o_\mathfrak{a}|p)P_\mathfrak{a}(p), \quad (o_\mathfrak{a}, p) \in \boldsymbol{o}_\mathfrak{a} \times \boldsymbol{p}. \tag{8}$$

Thus, the agent's joint pd, $J_\mathfrak{a}$, at most approximates the joint pds $J_{\mathfrak{a}\mathfrak{f}}$ (6) and $J_{\mathfrak{n}\mathfrak{f}}$ (7) *if they are given by a well-opted posterior pd*, $P_\mathfrak{f} \in \mathbf{P}$. The posterior pd is well opted if it really exploits the fuser information, $I_\mathfrak{f}$ (3). The use of this qualitative observation requires a quantitative expression of proximity of a pds pair. The works [7, 19] have shown that the approximation quality *is to be measured by the KLD of the approximated pd to its approximant* (unlike the popular variational Bayes method). This implies that good fusing rules, $S_\mathfrak{f}$ (4), make the expected values of $D(J_{\mathfrak{a}\mathfrak{f}}||J_\mathfrak{a})$, $D(J_{\mathfrak{n}\mathfrak{f}}||J_\mathfrak{a})$ small. This specifies the set, $\mathbf{S}_\mathfrak{f}$, of prospective fusing rules acting on $P_\mathfrak{f} \in \mathbf{P}$

$$\mathbf{S}_\mathfrak{f} \equiv \left\{ S(P_\mathfrak{f}) \equiv S(P_\mathfrak{f}|I_\mathfrak{f}) : \int_{\mathbf{P}} S(P_\mathfrak{f})D(J_{\mathfrak{a}\mathfrak{f}}||J_\mathfrak{a}) \, dP_\mathfrak{f} \leq b_\mathfrak{a} < \infty \right. \tag{9}$$

$$\left. \text{and} \quad \int_{\mathbf{P}} S(P_\mathfrak{f})D(J_{\mathfrak{n}\mathfrak{f}}||J_\mathfrak{a}) \, dP_\mathfrak{f} \leq b_\mathfrak{n} < \infty \right\},$$

---

[3] The work [43] calls the same functional "cross-entropy". The use of this term is often challenged so we stay with the name "Kulback-Leibler divergence".

[4] The agreed implicit conditioning on the agent's action, $a_\mathfrak{a}$, and its regressor, $r_\mathfrak{a}$, applies.

and completes the formulation of the optimisation task (5). The optional bounds, $b_\mathfrak{a}$ and $b_\mathfrak{n}$, parameterise the set (9) and ensure its non-emptiness. The next proposition provides the optimal fusing rule (5).

**Proposition 1 (Optimal Fusing)** *The optimal fusing rule,* $\mathsf{S}^{opt} \in \mathbf{S}_\mathfrak{f} \neq \emptyset$, *(5), (9), is* [5]

$$\mathsf{S}^{opt}(\mathsf{P}_\mathfrak{f}) \propto \mathsf{S}_\mathfrak{a}(\mathsf{P}_\mathfrak{f}) \exp\left[ -(\lambda_\mathfrak{a} + \lambda_\mathfrak{n}) \mathsf{D}(\mathsf{P}_\mathfrak{f} || \mathsf{P}^{opt}) \right] \tag{10}$$

$$\mathsf{P}^{opt}(p) \propto \mathsf{P}_\mathfrak{a}(p) \exp\left[ w \int_{o_\mathfrak{a}} \mathsf{F}_\mathfrak{n}(o_\mathfrak{a}) \ln(\mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p)) \, do_\mathfrak{a} \right], \;\; w \equiv \frac{\lambda_\mathfrak{n}}{\lambda_\mathfrak{a} + \lambda_\mathfrak{n}} \in [0,1].$$

*Kuhn-Tucker's multipliers [25], $\lambda_\mathfrak{a} \geq 0$, $\lambda_\mathfrak{n} \geq 0$, are chosen so that the constraints determining the set (9) are met.*

*The bounds $b_\mathfrak{a}$, $b_\mathfrak{n}$ in (9) are chosen so that at least one is active so that $\lambda_\mathfrak{a} + \lambda_\mathfrak{n} > 0$.*

*If the agent $\mathfrak{a}$ has no prior wish on the fusing rule and identifies $\mathsf{S}_\mathfrak{a} = \mathsf{S}_\mathfrak{f}$ [18] then the optimal fusing rule is deterministic and concentrates on $\mathsf{P}^{opt}(p)$ (10).*

*Proof:* The optimised KLD is a strictly convex functional on the convex set (9). Thus, a unique minimum exists and can be found by minimising the Kuhn-Tucker's functional. It is given by the non-negative Kuhn-Tucker's multipliers, $\lambda_\mathfrak{a}, \lambda_\mathfrak{n}$, chosen so that the constraints in (9) are met. In the next expression of this functional, the KLDs of the involved joint pds $\mathsf{J}_{\mathfrak{af}}, \mathsf{J}_{\mathfrak{nf}}, \mathsf{J}_\mathfrak{a}$, (6), (7), (8), are explicitly written. Also, the functional arguments are re-arranged, the normalisation of pds and Fubini's theorem on multiple integrations [39] are exploited

$$\int_\mathbf{P} \mathsf{S}_\mathfrak{f}(\mathsf{P}_\mathfrak{f}) \Bigg\{ \ln\left(\frac{\mathsf{S}_\mathfrak{f}(\mathsf{P}_\mathfrak{f})}{\mathsf{S}_\mathfrak{a}(\mathsf{P}_\mathfrak{f})}\right) + \lambda_\mathfrak{a} \int_{o_\mathfrak{a} \times p} \mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p) \mathsf{P}_\mathfrak{f}(p) \ln\left(\frac{\mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p)\mathsf{P}_\mathfrak{f}(p)}{\mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p)\mathsf{P}_\mathfrak{a}(p)}\right) do_\mathfrak{a} dp$$

$$+ \;\; \lambda_\mathfrak{n} \int_{o_\mathfrak{a} \times p} \mathsf{F}_\mathfrak{n}(o_\mathfrak{a}) \mathsf{P}_\mathfrak{f}(p) \ln\left(\frac{\mathsf{F}_\mathfrak{n}(o_\mathfrak{a})\mathsf{P}_\mathfrak{f}(p)}{\mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p)\mathsf{P}_\mathfrak{a}(p)}\right) do_\mathfrak{a} dp \Bigg\} d\mathsf{P}_\mathfrak{f}$$

$$= \;\; \int_\mathbf{P} \mathsf{S}_\mathfrak{f}(\mathsf{P}_\mathfrak{f}) \Bigg\{ \ln\left(\frac{\mathsf{S}_\mathfrak{f}(\mathsf{P}_\mathfrak{f})}{\mathsf{S}_\mathfrak{a}(\mathsf{P}_\mathfrak{f})}\right) + (\lambda_\mathfrak{a} + \lambda_\mathfrak{n})$$

$$\times \;\; \int_p \mathsf{P}_\mathfrak{f}(p) \Bigg[ \ln\left(\frac{\mathsf{P}_\mathfrak{f}(p)}{\mathsf{P}_\mathfrak{a}(p)}\right) + \frac{\lambda_\mathfrak{n}}{\lambda_\mathfrak{a}+\lambda_\mathfrak{n}} \underbrace{\int_{o_\mathfrak{a}} \mathsf{F}_\mathfrak{n}(o_\mathfrak{a}) \ln\left(\frac{\mathsf{F}_\mathfrak{n}(o_\mathfrak{a})}{\mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p)}\right) do_\mathfrak{a}}_{-\int_{o_\mathfrak{a}} \mathsf{F}_\mathfrak{n} \ln\left(\mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p)\right) do_\mathfrak{a} + a\, constant} \Bigg] dp \Bigg\} d\mathsf{P}_\mathfrak{f}$$

$$\overset{(10)}{=} \int_\mathbf{P} \mathsf{S}_\mathfrak{f}(\mathsf{P}_\mathfrak{f}) \Bigg\{ \ln\left(\frac{\mathsf{S}_\mathfrak{f}(\mathsf{P}_\mathfrak{f})}{\mathsf{S}_\mathfrak{a}(\mathsf{P}_\mathfrak{f})}\right) + (\lambda_\mathfrak{a}+\lambda_\mathfrak{n}) \underbrace{\int_p \mathsf{P}_\mathfrak{f}(p) \ln\left(\frac{\mathsf{P}_\mathfrak{f}(p)}{\mathsf{P}^{opt}(p)}\right) dp}_{\mathsf{D}(\mathsf{P}_\mathfrak{f} || \mathsf{P}^{opt})} \Bigg\} d\mathsf{P}_\mathfrak{f}$$

$$+ \;\; a\; constant = \mathsf{D}(\mathsf{S}_\mathfrak{f} || \mathsf{S}^{opt}) + another\; constant.$$

---

[5] It uses the implicit conditioning $\mathsf{F}_\mathfrak{n}(o_\mathfrak{a}) = \mathsf{F}_\mathfrak{n}(o_\mathfrak{a}|a_\mathfrak{a}, r_\mathfrak{a})$, $\mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p) = \mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p, a_\mathfrak{a}, r_\mathfrak{a})$. The assumed neighbour, see Section 2, implies the relevance of $a_\mathfrak{a}, r_\mathfrak{a}$ in the forecasting pd $\mathsf{F}_\mathfrak{n}$.

The last expression gives the optimal fusing rule (10) as the KLD reaches its minimum for identical arguments. The weight $w$ is $w = \lambda_{\mathfrak{n}}/(\lambda_{\mathfrak{n}} + \lambda_{\mathfrak{a}}) \in [0,1]$ for arbitrary $\lambda_{\mathfrak{a}}, \lambda_{\mathfrak{n}} \geq 0$. The assumption that at least one bound in (9) is active implies that $\lambda_{\mathfrak{a}} + \lambda_{\mathfrak{n}} > 0$.

For $S_{\mathfrak{a}}(P_{\mathfrak{f}}) = S_{\mathfrak{f}}(P_{\mathfrak{f}})$, the optimised functional is linear in the optimised fusing rule, $S_{\mathfrak{f}} \in \mathbf{S}_{\mathfrak{f}}$, and the basic lemma of stochastic control [2] applies.                                        □

## Discussion 1

✓ *The weight, $w$ (10), is zero for $\lambda_{\mathfrak{n}} = 0$. Then, the information about the neighbour's forecasting pd, $F_{\mathfrak{n}}$, does not change the agent's parameter description, $P_{\mathfrak{a}}(p)$.*

✓ *The multiplier $\lambda_{\mathfrak{n}}$ is zero (giving $w = 0$) if the corresponding proximity bound, $b_{\mathfrak{n}}$ (9), is chosen so large that it is not reached. A large bound means that the fuser assigns a negligible relevance to the forecasting pd offered by the neighbour.*

✓ *The weight $w \to 1$ for large values of $\lambda_{\mathfrak{n}}$ reflecting tight bound $b_{\mathfrak{n}}$. It expresses a high importance assigned to the neighbour by the fuser.*

✓ *Altogether, the weight (10) can be safely interpreted as the reliability assigned by the fuser to the neighbour.*

✓ *The value $\lambda_{\mathfrak{a}} + \lambda_{\mathfrak{n}}$ in (10) controls the closeness of the randomly sampled $P_{\mathfrak{f}} \sim S^{\mathsf{opt}}$ to $P^{\mathsf{opt}}$. Thus, both the relative trust weight, $w$ (10), and the individual values of Kuhn-Tucker multipliers, $\lambda_{\mathfrak{a}}, \lambda_{\mathfrak{n}}$, play a significant role in the proposed fusion.*

✓ *The individual multipliers reflect the individual tightness of bounds $b_{\mathfrak{a}}, b_{\mathfrak{n}}$ in the set (9). We assumed that at least one is tight, so that $\lambda_{\mathfrak{a}} + \lambda_{\mathfrak{n}} > 0$.*

✓ *The used formulation of the fusion task provides the top randomisation level that enriches the solution space. Indeed, the randomisation of the fusing rules $S_{\mathfrak{f}}$ (4) adds no flexibility [1].*

## 4 APPLICATION TO EXPONENTIAL FAMILY

The complexity curse is the main reason for the information processing requiring a fusion, see Section 1. The same reason motivates the wide-spread use of parametric models from the EF. The EF includes the vast majority of models that admit a sufficient statistic of a fixed finite dimension [23]. This property allows to convert the functional Bayes' rule into the algebraic recursion exactly. This motivates the presented specialisation of Proposition 1 to parametric models from the EF.

A member of the EF forecasts observation $o_{\mathfrak{a}} \in \boldsymbol{o}_{\mathfrak{a}}$ by the parametric model

$$M_{\mathfrak{a}}(o_{\mathfrak{a}}|p) = M_{\mathfrak{a}}(o_{\mathfrak{a}}|p, a_{\mathfrak{a}}, r_{\mathfrak{a}}) \equiv \exp \langle A(o_{\mathfrak{a}}, a_{\mathfrak{a}}, r_{\mathfrak{a}}), B(p) \rangle. \tag{11}$$

There, $\langle A, B \rangle$ is the scalar product of finite-dimensional, real values of the known functions, $A(o_{\mathfrak{a}}, a_{\mathfrak{a}}, r_{\mathfrak{a}})$ and $B(p)$. Here, the finite-dimensional action, $a_{\mathfrak{a}}$, and the regressor, $r_{\mathfrak{a}}$, forming a part of information $I_{\mathfrak{a}}$ (1), are explicitly referred to. The EF members have the conjugated (self-reproducing) prior pd [6]

$$P_{\mathfrak{a}}(p) \equiv P_{\mathfrak{a}}(p|V_{\mathfrak{a}}) \propto \exp \langle V_{\mathfrak{a}}, B(p) \rangle, \quad p \in \boldsymbol{p}. \tag{12}$$

It is given by a finite-dimensional real array, $V_{\mathfrak{a}}$, for which the scalar product $\langle V_{\mathfrak{a}}, B(p) \rangle$ makes sense and for which the function (12) is normalisable to a pd.

The next proposition just specialises Proposition 1 to the EF.

**Proposition 2 (Optimal Fusing in EF)** *Let the parametric model* $\mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p,a_\mathfrak{a},r_\mathfrak{a})$ *(11) be given. Let the agent,* $\mathfrak{a}$, *employ the conjugated prior pd,* $\mathsf{P}_\mathfrak{a}(p|V_\mathfrak{a})$ *(12). Then, the optimal fusing rule (5) within the set (9) has the general form (10) given by the conjugated pd,* $\mathsf{P}^{\mathsf{opt}}(p|V^{\mathsf{opt}}) \propto \exp\left\langle V^{\mathsf{opt}},\mathsf{B}(p)\right\rangle$, *with*

$$V^{\mathsf{opt}} = V_\mathfrak{a} + w\int_{o_\mathfrak{a}} \mathsf{F}_\mathfrak{n}(o_\mathfrak{a}|a_\mathfrak{a},r_\mathfrak{a})\mathsf{A}(o_\mathfrak{a},a_\mathfrak{a},r_\mathfrak{a})\,\mathrm{d}o_\mathfrak{a}, \quad w \in [0,1]. \tag{13}$$

**Discussion 2**

✓ *Bayes' rule updating the conjugated pd (12) by the data record,* $o_\mathfrak{a},a_\mathfrak{a},r_\mathfrak{a}$, *provides the posterior pd of this functional form with* $V_\mathfrak{a}$ *replaced by the statistic value*

$$V^{\mathsf{opt}} = V_\mathfrak{a} + \mathsf{A}(o_\mathfrak{a},a_\mathfrak{a},r_\mathfrak{a}). \tag{14}$$

*This offers an insight into our information fusion. The conjugated pd,* $\mathsf{P}^{\mathsf{opt}}(p|V)$, *delimiting the optimal fuser of the functional form (12), is given by the statistic (13). Thus, instead of incrementing* $V_\mathfrak{a}$ *by the value of* $\mathsf{A}$ *in the yet unavailable observation, the fuser increments* $V_\mathfrak{a}$ *by the expectation of* $\mathsf{A}(o_\mathfrak{a},a_\mathfrak{a},r_\mathfrak{a})$ *with respect to the forecasting pd,* $\mathsf{F}_\mathfrak{n}(o_\mathfrak{a}|a_\mathfrak{a},r_\mathfrak{a})$. *The expected value is weighted by the trust weight,* $w \in [0,1]$, *assigned by the fuser to the information-offering neighbour.*

✓ *If the neighbour supplies an* $\mathsf{F}_\mathfrak{n}(o_\mathfrak{a}|a_\mathfrak{a},r_\mathfrak{a})$, *which is fully concentrated on a crisp observation, then the standard Bayesian updating recovers whenever the fuser fully trusts the neighbour and sets* $w = 1$.

✓ *A crisp observation and* $w < 1$ *provide a weighted Bayes' rule handling an unreliable likelihood arising due to, for instance, its approximate evaluation.*

**Example 1 (Fusion Supports Markov Decision Processes)** *A Markov, action-dependent, environment model is a key ingredient of widely-used Markov decision processes [32]. This example shows that our theory may enhance its learning.*

*The parametric Markov model with discrete-valued observable state,* $o_\mathfrak{a} \in o_\mathfrak{a}$, *and action,* $a \in a$, *is the key EF member. It is parameterised by an array of transition probabilities,* $p$. *Its entry* $p_{o|a,r}$ *is the probability of the next environment state,* $o \in o_\mathfrak{a}$, *if the action* $a \in a_\mathfrak{a}$ *is chosen and the environment is in a state defining its regressor* $r \in r_\mathfrak{a} = o_\mathfrak{a}$. *Formally,*

$$\mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p) = \mathsf{M}(o_\mathfrak{a}|p,a_\mathfrak{a},r_\mathfrak{a}) \equiv p_{o_\mathfrak{a}|a_\mathfrak{a},r_\mathfrak{a}} = \prod_{(o,a,r)\in(o_\mathfrak{a},a_\mathfrak{a},r_\mathfrak{a})} p_{o|a,r}^{\delta[(o,a,r),(o_\mathfrak{a},a_\mathfrak{a},r_\mathfrak{a})]}$$

$$= \exp\left[\sum_{(o,a,r)\in(o_\mathfrak{a},a_\mathfrak{a},r_\mathfrak{a})} \delta[(o,a,r),(o_\mathfrak{a},a_\mathfrak{a},r_\mathfrak{a})]\ln(p_{o|a,r})\right], \tag{15}$$

*where* $\delta[\bullet,\star] \equiv 1$ *if* $\bullet = \star$, $\delta[\bullet,\star] \equiv 0$ *if* $\bullet \neq \star$ *is Kronecker's delta.*

*The last form in (15) shows that* $\mathsf{M}_\mathfrak{a}$ *is from the EF with* $\mathsf{A}_{o|a,r}(o_\mathfrak{a},a_\mathfrak{a},r_\mathfrak{a}) \equiv \delta[(o,a,r),(o_\mathfrak{a},a_\mathfrak{a},r_\mathfrak{a})]$, $\mathsf{B}_{o|a,r}(p) \equiv \ln(p_{o|a,r})$, *and* $\langle\mathsf{A},\mathsf{B}\rangle \equiv \sum_{(o,a,r)} \mathsf{A}_{o|a,r}\mathsf{B}_{o|a,r}$.

*The conjugated prior pd has the finite-dimensional sufficient statistic with entries* $V_{\mathfrak{a},o|a,r} > 0$, $o \in o_\mathfrak{a}$, $a \in a_\mathfrak{a}$, $r \in r_\mathfrak{a}$. *This conjugated prior is Dirichlet's pd [18]*

$$\mathsf{P}_\mathfrak{a}(p|V_\mathfrak{a}) \propto \prod_{(o,a,r)\in(o_\mathfrak{a},a_\mathfrak{a},r_\mathfrak{a})} p_{o|a,r}^{V_{\mathfrak{a},o|a,r}-1}. \tag{16}$$

*The Bayesian updating (14) with an observed triple $o_\mathfrak{a}, a_\mathfrak{a}, r_\mathfrak{a}$ reduces to counting*

$$V_{\mathfrak{a}, o_\mathfrak{a} | a_\mathfrak{a}, r_\mathfrak{a}} = V_{\mathfrak{a}, o_\mathfrak{a} | a_\mathfrak{a}, r_\mathfrak{a}} + 1. \tag{17}$$

*The optimal fusing (13) for the given $a_\mathfrak{a}$, $r_\mathfrak{a}$ reduces to the updating*

$$V^{\mathsf{opt}}_{o | a_\mathfrak{a}, r_\mathfrak{a}} = V_{\mathfrak{a}, o | a_\mathfrak{a}, r_\mathfrak{a}} + w \mathsf{F}_\mathfrak{n}(o | a_\mathfrak{a}, r_\mathfrak{a}), \quad \forall o \in \boldsymbol{o}_\mathfrak{a}. \tag{18}$$

**Discussion 3**

- ✓ *If the neighbour supplies a crisp observation $o_\mathfrak{a}$ with $\mathsf{F}_\mathfrak{n}(o_\mathfrak{a} | a_\mathfrak{a}, r_\mathfrak{a}) = 1$ then (18) reduces to (17) if the fuser fully trusts the neighbour and sets $w = 1$.*
- ✓ *The standard Bayesian estimation and fusion of the neighbour's information may run whenever the data record $(o_\mathfrak{a}, a_\mathfrak{a}, r_\mathfrak{a})$ or the forecasting pd $\mathsf{F}_\mathfrak{n}(o | a_\mathfrak{a}, r_\mathfrak{a})$, $o \in \boldsymbol{o}_\mathfrak{a}$, are available. No prior updating schedule is needed.*
- ✓ *The estimation and fusion provide the values, $V_\mathfrak{x}$, $\mathfrak{x} \in \{\mathfrak{a}, \mathfrak{n}, \mathfrak{f}\}$. They allow to forecast the observation $o_\mathfrak{a} \in \boldsymbol{o}_\mathfrak{a}$ for any action $a_\mathfrak{a} \in \boldsymbol{a_\mathfrak{a}}$ and any regressor $r_\mathfrak{a} \in \boldsymbol{r_\mathfrak{a}} = \boldsymbol{o_\mathfrak{a}}$, [18],*

$$\mathsf{F}_\mathfrak{x}(o_\mathfrak{a} | a_\mathfrak{a}, r_\mathfrak{a}, V_\mathfrak{x}) = \frac{V_{\mathfrak{x}, o_\mathfrak{a} | a_\mathfrak{a}, r_\mathfrak{a}}}{\sum_{o \in \boldsymbol{o}_\mathfrak{a}} V_{\mathfrak{x}, o | a_\mathfrak{a}, r_\mathfrak{a}}}, \quad \mathfrak{x} \in \{\mathfrak{a}, \mathfrak{n}, \mathfrak{f}\}. \tag{19}$$

- ✓ *The formula (19) provides a way in which the neighbour may obtain the forecasting pd $\mathsf{F}_\mathfrak{n}$. It simply collects its statistic $V_\mathfrak{n}$ on an other data set than the processed one and uses its version of (19). For instance, it processes the data set concerning the whole population of modelled environments. The neighbour's array, $V_{\mathfrak{n}, o | a, r}$, $(o, a, r) \in (\boldsymbol{o}_\mathfrak{a}, \boldsymbol{a}_\mathfrak{a}, \boldsymbol{r}_\mathfrak{a})$, may be a sub-array of an off-line-collected $V_{\mathfrak{n}, o | a, r}$ with $o \in \boldsymbol{o}_\mathfrak{n} \supseteq \boldsymbol{o}_\mathfrak{a}$, $a \in \boldsymbol{a}_\mathfrak{n} \supseteq \boldsymbol{a}_\mathfrak{a}$, $r \in \boldsymbol{r}_\mathfrak{n} \supseteq \boldsymbol{r}_\mathfrak{a}$.*
- ✓ *The result (18) has a great appeal. Probabilities $\mathsf{F}_\mathfrak{n}(o | a_\mathfrak{a}, r_\mathfrak{a})$, $\forall o \in \boldsymbol{o}_\mathfrak{a}$ and given $a_\mathfrak{a}, r_\mathfrak{a}$, replace the values of $\delta[(o, a_\mathfrak{a}, r_\mathfrak{a}), (o_\mathfrak{a}, a_\mathfrak{a}, r_\mathfrak{a})]$ (15) that are unknown when the observation $o_\mathfrak{a}$ is yet unavailable. This replacement coincides with the heuristic called quasi-Bayes estimation [14, 44].*

**Example 2 (Fusion for Linear Gaussian Models)** *A linear-in-regression coefficients, $\theta$, Gaussian model of a real vector, $o_\mathfrak{a}$, with a constant conditional covariance, $\rho$, is another prominent example of the (dynamic) EF parameterised by $p \equiv (\theta, \rho)$. If ▷ the action and regressor fill a column vector with unity at its end ▷ I is a unit matrix of the size of $o_\mathfrak{a}$ ▷ ′ is transposition, then the form (11) of an EF member is gained with*

$$\langle \mathsf{A}, \mathsf{B} \rangle = \mathrm{tr}[\mathsf{A}'\mathsf{B}], \qquad \mathsf{A}(o_\mathfrak{a}, a_\mathfrak{a}, r_\mathfrak{a}) = \begin{bmatrix} o_\mathfrak{a} \\ a_\mathfrak{a} \\ r_\mathfrak{a} \end{bmatrix} [o_\mathfrak{a}', a_\mathfrak{a}', r_\mathfrak{a}']$$

$$p = (\theta, \rho), \qquad \mathsf{B}(p) = -0.5 \begin{bmatrix} -\mathsf{I} & \theta \\ 0 & 1 \end{bmatrix}' \begin{bmatrix} \rho^{-1} & 0 \\ 0 & \ln(|\rho|) \end{bmatrix} \begin{bmatrix} -\mathsf{I} & \theta \\ 0 & 1 \end{bmatrix}.$$

*The sufficient statistic is a positive-definite, extended information matrix $V_\mathfrak{a}$. It determines the conjugated Gauss-inverse-Wishart pd of the unknown parameter. The*

*Bayesian updating (14) reduces to the recursive least squares, for details see [36]. The optimal fusing (13) reads*

$$V^{\text{opt}} = V_{\mathfrak{a}} + w \int_{o_{\mathfrak{a}}} \mathsf{F}_{\mathfrak{n}}(o_{\mathfrak{a}}|a_{\mathfrak{a}},r_{\mathfrak{a}}) \begin{bmatrix} o_{\mathfrak{a}} \\ a_{\mathfrak{a}} \\ r_{\mathfrak{a}} \end{bmatrix} [o'_{\mathfrak{a}}, a'_{\mathfrak{a}}, r'_{\mathfrak{a}}] \, \mathrm{d}o_{\mathfrak{a}} \tag{20}$$

$$= V_{\mathfrak{a}} + w \left\{ \begin{bmatrix} \bar{o}_{\mathfrak{n}} \\ a_{\mathfrak{a}} \\ r_{\mathfrak{a}} \end{bmatrix} [\bar{o}'_{\mathfrak{n}}, a'_{\mathfrak{a}}, r'_{\mathfrak{a}}] + \begin{bmatrix} cov_{\mathfrak{n}}(o_{\mathfrak{a}}|a_{\mathfrak{a}},r_{\mathfrak{a}}) & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \right\}.$$

*Thus, the fuser updates $V_{\mathfrak{a}}$ by the dyad with the unavailable $o_{\mathfrak{a}}$ replaced by its expectation $\bar{o}_{\mathfrak{n}} = \int_{o_{\mathfrak{a}}} o_{\mathfrak{a}} \mathsf{F}_{\mathfrak{n}}(o_{\mathfrak{a}}|a_{\mathfrak{a}},r_{\mathfrak{a}}) \, \mathrm{d}o_{\mathfrak{a}}$. The corresponding sub-matrix of $V_{\mathfrak{a}}$ is, moreover, increased by the covariance $cov_{\mathfrak{n}}(o_{\mathfrak{a}}|a_{\mathfrak{a}},r_{\mathfrak{a}})$ given by $\mathsf{F}_{\mathfrak{n}}(o_{\mathfrak{a}}|a_{\mathfrak{a}},r_{\mathfrak{a}})$. This result appeared in [38] with the heuristically motivated weight, w.*

*If $\mathfrak{n}$ supplies a crisp observation $o_{\mathfrak{a}}$ with $cov_{\mathfrak{n}}(o_{\mathfrak{a}}|a_{\mathfrak{a}},r_{\mathfrak{a}}) = 0$ then $\bar{o}_{\mathfrak{n}} = o_{\mathfrak{a}}$ and the formula (20) reduces to the standard Bayesian updating with the processed data record weighted by $\sqrt{w}$.*

## 5 AN EXTENDED FORMULATION AND ITS SOLUTION

The proposed fusing treats $\mathsf{F}_{\mathfrak{n}}$ as a single data record. It is obvious when considering Propositions 1, 2 with this forecasting pd concentrated on a crisp observation. At the same time, it is clear that information contents of a few and many data records may differ substantially even if they lead to the same forecasting pd. In other words, the sufficient statistic representing the data records is their sample pd $\tilde{F}_{\mathfrak{n}}$ *and* their (effective) number $\nu_{\mathfrak{n}}$.

This section respects the recalled fact under an additional assumption: the neighbour, $\mathfrak{n}$, provides the forecasting pd, $\mathsf{F}_{\mathfrak{n}}(o_{\mathfrak{a}}|a_{\mathfrak{a}},r_{\mathfrak{a}})$, *relevant for the realised action, $a_{\mathfrak{a}}$, and the realised regressor, $r_{\mathfrak{a}}$, together with the effective number, $\nu_{\mathfrak{n}} \geq 1$, of data records $d_{\mathfrak{n}} = (o_k, a_{\mathfrak{a}}, r_{\mathfrak{a}})_{k=1}^{\nu_{\mathfrak{n}}}$ that led to the forecasting pd $\mathsf{F}_{\mathfrak{n}}(o_{\mathfrak{a}}|a_{\mathfrak{a}},r_{\mathfrak{a}})$.*

The effective number $\nu_{\mathfrak{n}}$ may coincide with the real number of processed data records or may be lower due to the use of a weighted Bayes' rule, see [26] and Discussion 2. It can also be quite subjective to express the number of fictitious data records the neighbour used for the construction of the forecasting pd, $\mathsf{F}_{\mathfrak{n}}$.

Let us imagine that $\nu_{\mathfrak{n}}$ data records $d_{\mathfrak{n}}$ with *common* $a_{\mathfrak{a}}, r_{\mathfrak{a}}$ were fed into Bayes' rule with the agent's parametric model $\mathsf{M}_{\mathfrak{a}}(o_{\mathfrak{a}}|p) = \mathsf{M}_{\mathfrak{a}}(o_{\mathfrak{a}}|p,a_{\mathfrak{a}},r_{\mathfrak{a}})$ and the pd $\mathsf{P}_{\mathfrak{a}}(p)$. It gives [24]

$$\tilde{\mathsf{P}}_{\mathfrak{n}}(p|d_{\mathfrak{n}}) \propto \mathsf{P}_{\mathfrak{a}}(p) \exp \left[ \nu_{\mathfrak{n}} \int_{o_{\mathfrak{a}}} \tilde{\mathsf{F}}_{\mathfrak{n}}(o_{\mathfrak{a}}|a_{\mathfrak{a}},r_{\mathfrak{a}}) \ln(\mathsf{M}_{\mathfrak{a}}(o_{\mathfrak{a}}|p)) \, \mathrm{d}o_{\mathfrak{a}} \right] \tag{21}$$

$$\tilde{\mathsf{F}}_{\mathfrak{n}}(o_{\mathfrak{a}}|a_{\mathfrak{a}},r_{\mathfrak{a}}) \equiv \frac{1}{\nu_{\mathfrak{n}}} \sum_{k=1}^{\nu_{\mathfrak{n}}} \delta[o_{\mathfrak{a}}, o_k].$$

A fair neighbour believes that its forecasting pd $\mathsf{F}_{\mathfrak{n}}$ models reality well and is close to the sample pd $\tilde{\mathsf{F}}_{\mathfrak{n}}$ in (21) obtained from $\nu_{\mathfrak{n}}$ data records. This motivates us to *assign*[6]

---

[6] Let us stress that the neighbour, $\mathfrak{n}$, is generally unaware of the model, $\mathsf{M}_{\mathfrak{a}}$, and its parameter, $p$.

to the neighbour's information, $I_\mathfrak{n} = (F_\mathfrak{n}, \nu_\mathfrak{n})$, the posterior pd

$$P_\mathfrak{n}(p) \propto P_\mathfrak{a}(p) \exp\left[\nu_\mathfrak{n} \int_{o_\mathfrak{a}} F_\mathfrak{n}(o_\mathfrak{a}) \ln(M_\mathfrak{a}(o_\mathfrak{a}|p)) \, do_\mathfrak{a}\right]. \tag{22}$$

The *assignment* (22) leads to the alternative, comparing to (9), specification of the set of prospective fusing rules[7]

$$\mathbf{S}_\mathfrak{f} = \left\{ S(P_\mathfrak{f}) = S(P_\mathfrak{f}|I_\mathfrak{f}) : \int_{\mathbf{P}_\mathfrak{f}} S(P_\mathfrak{f}) D(P_\mathfrak{f}||P_\mathfrak{a}) \, dP_\mathfrak{f} \le b_\mathfrak{a} < \infty \tag{23}$$

$$\text{and} \quad \int_{\mathbf{P}_\mathfrak{f}} S(P_\mathfrak{f}) D(P_\mathfrak{f}||P_\mathfrak{n}) \, dP_\mathfrak{f} \le b_\mathfrak{n} < \infty \right\}.$$

The choice (23) gives the next solution of the task (5).

**Proposition 3 (Optimal Fusing with $\nu_\mathfrak{n} \ge 1$)** *The optimal fusing rule, $S^\text{opt} \in \mathbf{S}_\mathfrak{f} \ne \emptyset$, (5), (22), (23), has the form*

$$S^\text{opt}(P_\mathfrak{f}) \propto S_\mathfrak{a}(P_\mathfrak{f}) \exp\left[-(\lambda_\mathfrak{a} + \lambda_\mathfrak{n})D(P_\mathfrak{f}||P^\text{opt})\right] \tag{24}$$

$$P^\text{opt}(p) \propto P_\mathfrak{a}(p) \exp\left[w\nu_\mathfrak{n} \int_{o_\mathfrak{a}} F_\mathfrak{n}(o_\mathfrak{a}) \ln(M_\mathfrak{a}(o_\mathfrak{a}|p)) do_\mathfrak{a}\right], \quad w = \frac{\lambda_\mathfrak{n}}{\lambda_\mathfrak{a} + \lambda_\mathfrak{n}} \in [0,1]$$

$$F_\mathfrak{n}(o_\mathfrak{a}) = F_\mathfrak{n}(o_\mathfrak{a}|a_\mathfrak{a}, r_\mathfrak{a}), \quad M_\mathfrak{a}(o_\mathfrak{a}|p) = M_\mathfrak{a}(o_\mathfrak{a}|p, a_\mathfrak{a}, r_\mathfrak{a}).$$

*Kuhn-Tucker's multipliers [25], $\lambda_\mathfrak{a} \ge 0$, $\lambda_\mathfrak{n} \ge 0$, ensure that the constraints determining the set (23) are met. They fully replace the $\mathbf{S}_\mathfrak{f}$-parametrisation by bounds $b_\mathfrak{a}$, $b_\mathfrak{n}$.*

*If the agent has no prior wish on the fusing rule and identifies $S_\mathfrak{a} = \mathbf{S}_\mathfrak{f}$, [18] then the optimal fusing rule is deterministic and concentrates on $P^\text{opt}(p)$ (24).*

*Proof:* It is omitted as it in fact copies that of Proposition 1.                      □

### Discussion 4

✓ *Proposition 1 is the special case of Proposition 3 for $\nu_\mathfrak{n} = 1$. It confirms the claim introducing and motivating this section.*

✓ *If $F_\mathfrak{n}(o_\mathfrak{a}) = F_\mathfrak{n}(o_\mathfrak{a}|a_\mathfrak{a}, r_\mathfrak{a})$ is indeed the sample pd gained from $\nu_\mathfrak{n}$ realised data records then $P^\text{opt}(p)$ corrects $P_\mathfrak{a}(p)$ by Bayes' rule with the likelihood flattened by the trust weight $w \in [0,1]$.*

✓ *Proposition 3 specialises to the EF by using $w\nu_\mathfrak{n}$ instead of $w \in [0,1]$ in Proposition 2. The weight w remains to be the learnable trust weight. It may counteract a too high self-confidence of the neighbour expressed by a high offered value of $\nu_\mathfrak{n}$.*

## 6 SIMULATION EXAMPLES

This section illustrates the theory forming the core of the paper. The desirable real-life tests will be published independently.

---

[7] It uses the KLDs of posterior pds not the KLDs of joint pds.

## 6.1 Monte Carlo study with static environments

The example is intentionally simple to meet its illustrative purpose. Its description follows the basic blocks in Figure 1 and uses the next auxiliary vectors

$$v_{\mathfrak{a}} = [1,2,9,3,1,1,1,1,1], \qquad v_{\mathfrak{n}} = [1,1,1,1,1,3,9,2,1]. \tag{25}$$

**Simulated environments** generated sequences of observations $o_{t\mathfrak{a}}$ in the set

$$\boldsymbol{o}_{\mathfrak{a}} = \{1,2,\dots,9\} \ \text{ at time moments } \ t \in \boldsymbol{t} = \{1,2,\dots,10\}. \tag{26}$$

The observations were independent and influenced neither by actions, $a_{\mathfrak{a}}$, nor by regressors, $r_{\mathfrak{a}}, r_{\mathfrak{n}}$. The static environments were described by the next pds, see (25),

$$p_{o_{\mathfrak{a}}|a_{\mathfrak{a}},r_{\mathfrak{a}},r_{\mathfrak{n}},\alpha} = p_{o_{\mathfrak{a}}|\alpha} = \alpha \frac{v_{\mathfrak{n}o_{\mathfrak{a}}}}{\sum_{o \in \boldsymbol{o}_{\mathfrak{a}}} v_{\mathfrak{n}o}} + (1-\alpha) \frac{v_{\mathfrak{a}o_{\mathfrak{a}}}}{\sum_{o \in \boldsymbol{o}_{\mathfrak{a}}} v_{\mathfrak{a}o}} \tag{27}$$

$$\alpha \in \boldsymbol{\alpha} = \{0, 0.2, 0.4, 0.6, 0.8, 1\}.$$

**The agent,** $\mathfrak{a}$**,** was aware that observations are mutually independent and independent of its actions and of any regressor. It implied no need to generate actions and gave the parametric model (15)

$$\mathsf{M}_{\mathfrak{a}}(o_{\mathfrak{a}}|p, a_{\mathfrak{a}}, r_{\mathfrak{a}}) = \mathsf{M}_{\mathfrak{a}}(o_{\mathfrak{a}}|p) = p_{o_{\mathfrak{a}}} \geq 0, \quad \sum_{o \in \boldsymbol{o}_{\mathfrak{a}}} p_o = 1.$$

Dirichlet's conjugated pd (16) $\mathsf{P}_{\mathfrak{a}}(p|V_{\mathfrak{a}})$ was used. In each experiment, the statistic determining it was set $V_{\mathfrak{a}} \equiv v_{\mathfrak{a}}$ (25) and updated according to the algebraic version of Bayes' rule (17) by 10 observations $o_{t\mathfrak{a}} \in \boldsymbol{o}_{\mathfrak{a}}$, see (26).

**The neighbour,** $\mathfrak{n}$**,** provided the fixed forecasting pd, $\mathsf{F}_{\mathfrak{n}}(o_{\mathfrak{a}}) \propto v_{\mathfrak{n}o_{\mathfrak{a}}}$ (25), together with the effective number of data records, $v_{\mathfrak{n}} = 200$, see Section 5. This reflects a high self-confidence of the neighbour in its forecasting pd.

**The fuser,** $\mathfrak{f}$**,** used the forecasting pd, $\mathsf{F}_{\mathfrak{n}}$, and the effective number of data records, $v_{\mathfrak{n}}$, see Proposition 3. Thus, it increased the initial statistic to

$$V_{\mathfrak{a},o_{\mathfrak{a}}} \equiv V_{\mathfrak{a},o_{\mathfrak{a}}} + w v_{\mathfrak{n}} \mathsf{F}_{\mathfrak{n}}(o), \qquad o \in \boldsymbol{o}_{\mathfrak{a}} = \{1,\dots,9\}, \tag{28}$$

for a fixed trust weight

$$w \in \boldsymbol{w} \equiv \{0, 0.02, \dots, 0.98, 1\}. \tag{29}$$

**The evaluation** used Monte Carlo with $10^5$ runs for each $(\alpha, w) \in \boldsymbol{\alpha} \times \boldsymbol{w}$, (27), 29), determining the simulated environment (27) and the trust weight allocated to the neighbour's information. The runs differed in realised samples distributed according to (27). Bayes' rule processed the observed 10 samples. The KLD's of the pd of the simulated environment to their final point estimates, based on the realised observations, $(o_{t\mathfrak{a}})_{t=1}^{10}$, and the fuser information, $I_{\mathfrak{f}}$ (3), were evaluated. They were averaged over the Monte Carlo runs. Note that in this case the point estimate coincides with the forecasting pd.

**The results in Figure 2** have a direct interpretation. If a trust weight is properly chosen then the proposed fusion increases the estimation rate so vital for the short

data sequences. The observable deterioration of the estimation quality when the fused information is inadequate and $w$ is improperly chosen supports the interpretation of $w$ as the trust weight. It also confirms that a proper trust weight is learnable.

To accept these claims, it suffices to recall that the zero trust case into the neighbour's forecasting pd ($w = 0$, blue dashed line in Figure 2) gives the standard Bayesian estimation. Thus, the fusion results (red full line in Figure 2) above/below this base level indicate deterioration/improvement due to the fusion with a positive trust weight of the neighbour's information, $I_\mathfrak{n}$ (2).

Sub-figures in Figure 2 indicate: ▷ if the simulated environment, given by $\alpha \in \boldsymbol{\alpha}$ (27), differs much from the information (2) offered by the neighbour, it is better not to take the offer seriously ▷ if the reliable information is taken with the full trust then it helps significantly. The transition between these extremes is smooth. Its almost deterministic nature is due to: ▷ the deterministic contribution of the neighbour to the fused statistics (28) ▷ the volatility suppression by the number of Monte Carlo runs.

## 6.2 A case with dynamic environment

This example provides an additional insight into the proposed knowledge fusion. It deals with: ▷ the simulated *dynamic* environment ▷ the under-modelled parametric model ▷ the forecaster based on another under-modelled parametric model.

The example shows that even under these conditions the information brought by the forecaster notably improves the agent's forecasting ability.

**Simulated environments** generated sequences of observations $o_t$ in the set

$$o \equiv \{1,2\} \quad \text{at time moments} \quad t \in \boldsymbol{t} \equiv \{1,2,\dots,10\}. \tag{30}$$

The observations were generated by the second order Markov chain uninfluenced by agent's actions. The presented cases simulated the transition probabilities $p_{o_t|r_{t-1}}$

**Case 1**

$p_{o=1|r=[1,1]} \equiv 0.0$

$p_{o=1|r=[1,2]} \equiv 0.5$

$p_{o=1|r=[2,1]} \equiv 0.5$

$p_{o=1|r=[2,2]} \equiv 1.0$

**Case 2**

$p_{o=1|r=[1,1]} \equiv 0.2$

$p_{o=1|r=[1,2]} \equiv 0.4$    with    $r_{t-1} \equiv [o_{t-1}, o_{t-2}]$.

$p_{o=1|r=[2,1]} \equiv 0.4$

$p_{o=1|r=[2,2]} \equiv 0.9$

In all runs, initial regressor was $r_0 \equiv [o_0, o_{-1}] = [2,1]$. The random seed generator was reset to a common value when a new trust weight was inspected.

**The agent,** $\mathfrak{a}$**,** was aware that its observations, $o_\mathfrak{a} \equiv o \in \boldsymbol{o} = \{1,2\}$, are independent of its actions, $a_\mathfrak{a}$. The agent underestimated the environment dynamics. It used the first-order parametric model (15) $\mathsf{M}_\mathfrak{a}(o_\mathfrak{a}|p,a_\mathfrak{a},r_\mathfrak{a}) = p_{o_\mathfrak{a}|r_\mathfrak{a}}$ with $r_{(t-1)\mathfrak{a}} \equiv o_{(t-1)\mathfrak{a}}$. Dirichlet's conjugated pd (16) $\mathsf{P}_\mathfrak{a}(p|V_\mathfrak{a})$ was used. In each experiment, the prior statistics were set

$$
\begin{array}{ccc}
& \textbf{Case 1} & \textbf{Case 2} \\
V_\mathfrak{a} \equiv [V_{\mathfrak{a},o|r_\mathfrak{a}}]_{o,r_\mathfrak{a} \in \boldsymbol{o}} \equiv \begin{bmatrix} 10^{-6} & 1 \\ 1 & 10^{-6} \end{bmatrix} & V_\mathfrak{a} \equiv \begin{bmatrix} 2 & 1 \\ 1 & 3 \end{bmatrix}.
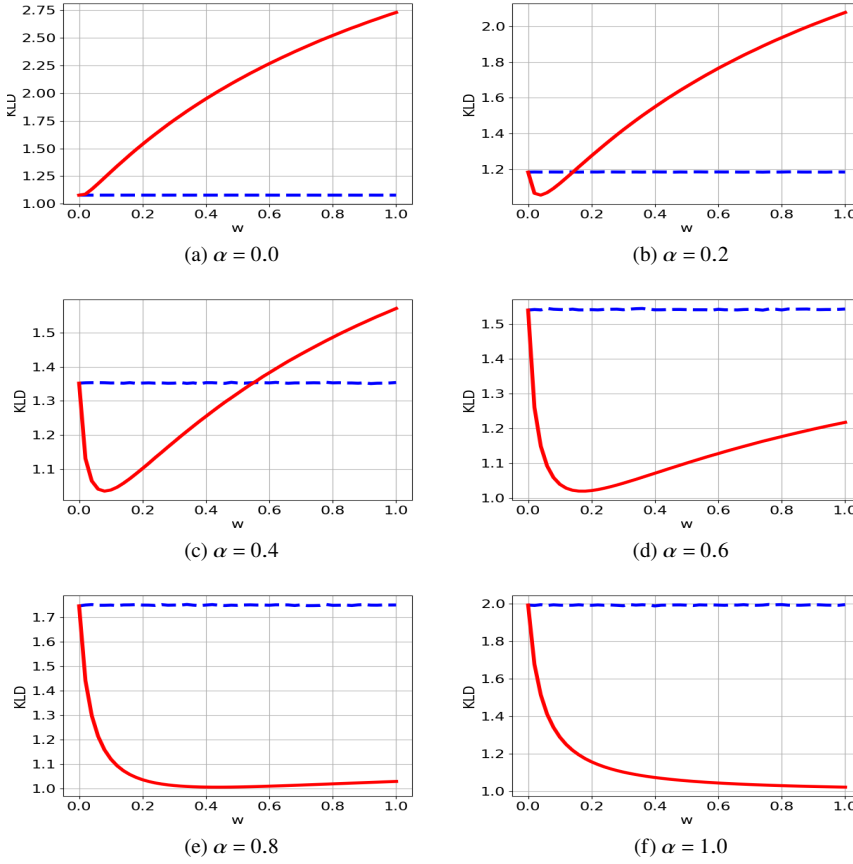\end{array}
$$

Fig. 2: The averaged KLD of simulated environment models to the estimated ones for different simulated environments given by values of parameter $\alpha$ (27). The red full line corresponds to the agent, $\mathfrak{a}$, that uses both (17) and (18) with various trust weights $w$. The blue dashed line corresponds to the agent, $\mathfrak{a}$, that just uses Bayes' rule (17).

The statistics were sequentially updated via the algebraic version of Bayes' rule (17) by 10 observations.

**The neighbour,** $\mathfrak{n}$**,** provided the forecasters gained by sequentially learning the approximate parametric model $\mathsf{M}_\mathfrak{n}(o_\mathfrak{a}|p_\mathfrak{n},a_\mathfrak{a},r_\mathfrak{n}) = p_{\mathfrak{n}o_\mathfrak{a}|r_\mathfrak{n}}$ of the observation $o_{t\mathfrak{a}}$ with the regressor $r_{(t-1)\mathfrak{n}} \equiv o_{(t-2)\mathfrak{a}}$, which gave $\mathsf{F}_\mathfrak{n}(o_\mathfrak{a}|r_\mathfrak{n}) \propto \mathsf{V}_{\mathfrak{n}o_\mathfrak{a}|r_\mathfrak{n}}$.

    The initial statistics were set

$$
\begin{array}{cc}
\textbf{Case 1} & \textbf{Case 2} \\
V_\mathfrak{n} \equiv [V_{\mathfrak{n},o|r_\mathfrak{n}}]_{o,r_\mathfrak{n}\in\boldsymbol{o}} \equiv \begin{bmatrix} 10^{-6} & 1 \\ 1 & 10^{-6} \end{bmatrix} & V_\mathfrak{n} \equiv \begin{bmatrix} 1 & 2 \\ 1 & 1 \end{bmatrix}
\end{array}
$$

but *unlike* $V_{\mathfrak{a}}$, which was *updated by* the observed *pairs* $o_{t\mathfrak{a}},o_{(t-1)\mathfrak{a}}$, $V_{\mathfrak{n}}$ was *updated* by the *pairs* $o_{t\mathfrak{a}},o_{(t-2)\mathfrak{a}}$. The effective number of processed data records in each time step was $\nu_{\mathfrak{n}} = 1$, cf. Section 5.

**The fuser,** $\mathfrak{f}$**,** used the forecasting pd, $\mathsf{F}_{\mathfrak{n}}$, and the effective number of data records, $\nu_{\mathfrak{n}} = 1$, Proposition 3 reduced to Proposition 2. Thus, it increased the updated statistic

$$V^{\text{opt}}_{o|r_{\mathfrak{a}}} = V_{\mathfrak{a},o|r_{\mathfrak{a}}} + w\mathsf{F}_{\mathfrak{n}}(o|r_{\mathfrak{n}}), \quad o \in \boldsymbol{o},$$

for a fixed trust weight $w \in \boldsymbol{w}$ (29).

**The evaluation** used Monte Carlo with $10^5$ runs for each $w \in \boldsymbol{w}$, (29). The runs differed in realised samples. The under-modelling made comparison in parameter spaces meaning-less. The quality of individual forecasters $\mathsf{F} \in \{\mathsf{F}_{\mathfrak{a}},\mathsf{F}_{\mathfrak{n}},\mathsf{F}_{\mathfrak{f}}\}$ was evaluated by the accuracy defined as

$$\text{Accuracy} = \frac{\text{number of } [\text{argmax}_{o \in \boldsymbol{o}}(\mathsf{F}_t(o)) = o_{t\mathfrak{a}}]_{t \in \boldsymbol{t}}}{10} \times 100 \quad [\%].$$

**The results in Figure 3** have a direct interpretation in **Case 1**. If a trust weight is properly chosen then the proposed fusion increases the agent's forecasting quality even when both the agent and its neighbour learn parametric models of wrong structures. Again, it confirms that a proper trust weight is learnable.

    **Case 2** represents the configuration in which the properly weighted forecaster again improves forecasting abilities of the agent but it does not guarantee that the gained quality will cross that reached by the neighbour. This is the cost for the considered unwillingness (inability) of the agent to employ another parametric model. The agent can see the price paid for the unwillingness and modify its parametric model if other circumstances allow this change.



Fig. 3: The averaged forecast accuracy. Blue dashed line reflects the agent's accuracy, green dots characterise the neighbour's accuracy and red full line corresponds with the fuser's outcomes. The left panel concerns **Case 1**, the right panel reflects **Case 2**.

## 7 REMARKS ON RELATED WORKS AND OPEN PROBLEMS

The presented scenario of information fusion complements the existing rich set of tools [31, 46]. To our best knowledge the incorporation of the information brought by the forecasting pd into the posterior pd of an unknown parameter has only been developed by us and by our colleagues. This explains why self-citations dominate our discussion. At the same time, we hope that this paper will contribute to a further development of our approach that has an extreme use range. We have used it, for instance, for knowledge elicitation [17] or distributed control [20]. The foreseen direct applications are discussed in Section 1. They concern mainly but not exclusively advisory and recommendation systems.

Mathematically, the inspected fusion is a case of combining pds [11]. The combination of pds operating on non-identical domains is the specificity of the proposed fusion. It primarily serves the targeted DM. The adopted Bayes' framework is important even when no specific DM objective is set but learning faces a lack of data. In the big-data era, it is a surprisingly frequent case. For instance, analysis of gene regulatory networks [13] or structure estimation of Bayesian networks [42] suffer from the data lack. Other learning problems like fraud detection [30] or building of sparse models [10] and many others are difficult due to the lack of *informative* data.

Methodologically, extensional or intensional combinations of partial probabilistic information exist. The insightful paper [35] favours the intensional, top-down approach and supports Bayesian networks [16]. Fuzzy methodology [48] represents clever extensional, bottom-up approaches. Our approach lies between these extremes. It steps out from the preferable intensional way in order to respect the limited agent's abilities. Still, it meets some challenges inherent to extensional technique, cf. [47].

On the other hand, comparing to our nearest predecessor [38], the proposed way derives the trust weight intensionally. It uses KLD "balls" when defining the sets (9), (23) of the suitable fusing rules. The choice of these balls is justified in [7, 19].

A survey of fusion techniques clearly shows two aspects relevant to our work: ▷ a lot of excellent work was done, e.g. [29] ▷ too much was done without clear guidelines, which makes a selection of an appropriate method for a specific problem error-prone [28]. A similar situation arises in artificial intelligence and surely in other areas as well. The deductive solutions as ours diminish this problem.

Technically, the proposed fusion can be simply extended to more neighbours, to more information sources. Importantly, the fusion can jointly use forecasting pds obtained by the objective-data-based estimation [36] and subjective experts' judgement [34]. It also processes a crisp data record in a way, which coincides with Bayes' rule whenever the observation source is qualified as reliable.

A closer inspection of the proposed information fusion way reveals a range of small open technical problems and surely-solvable problems like Bayesian learning of the trust weight. Open, conceptually hard, problems include the analysis of emergent behaviours of extensive networks sharing the information in the proposed way.

Even under the current research state, the proposed fusion can already be applied in the advisory and recommendation systems as well as within the internet of things [45] or cyber-physical-social systems [8]. They often need the fusion sketched in Figure 1.

## Declaration

*Conflict of interests:* The authors have no affiliation with any organization with a direct
or indirect financial interest in the subject matter discussed in the manuscript. This
manuscript has not been submitted to, nor is under review at, another journal or other
publishing venue.
*Availability of data and material:* Not applicable
*Code availability:* The code of examples is available at https://gitlab.com/hula-phd/bks.
*Authors' contributions:* Both authors tightly cooperated on the paper. MK dominated
in writing the text and FH in experiments.

## References

1. Antoniak, C.: Mixtures of Dirichlet processes with applications to Bayesian nonparametric problems. The Annals of Statistics **2**(6), 1152–1174 (1974)
2. Åström, K.: Introduction to Stochastic Control. Acad. Press, N.Y. (1970)
3. Bader, K., Lussier, B., Schon, W.: A fault tolerant architecture for data fusion: A real application of Kalman filters for mobile robot localization. Robotics and Autonomous Systems **88**, 11 – 23 (2017)
4. Bar-Shalom, Y., Li, X., Kirubarajan, T.: Estimation with Applications to Tracking and Navigation. Wiley (2003)
5. Barndorff-Nielsen, O.: Information and Exponential Families in Statistical Theory. Wiley, N.Y. (1978)
6. Berger, J.: Statistical Decision Theory and Bayesian Analysis. Springer (1985)
7. Bernardo, J.: Expected information as expected utility. The An. of Stat. **7**, 686–690 (1979)
8. Bogdan, P., Pedram, M.: Toward enabling automated cognition and decision-making in complex cyber-physical systems. In: 2018 IEEE ISCAS, pp. 1–4 (2018)
9. Foley, C., Quinn, A.: Fully probabilistic design for knowledge transfer in a pair of Kalman filters. IEEE Signal Proc. Letters **25**(4), 487–490 (2018)
10. Galeano, P., Pena, D.: Data science, big data and statistics. Test **28**, 289–325 (2019)
11. Genest, C., Zidek, J.: Combining probability distributions: A critique and annotated bibliography. Stat. Sci. **1**(1), 114–148 (1986)
12. Hall, D., Llinas, J.: An introduction to multisensor data fusion. Proc. of the IEEE **85**(1), 6–23 (1997)
13. Hlaváčková-Schindler, K., Naumova, V., Pereverzyev, S.: Granger causality for ill-posed problems: Ideas, methods, and application in life sciences. In: W. Wiedermann, A. von Eye (eds.) Statistics and Causality: Methods for Applied Empirical Research, pp. 249–276. Wiley (2016)
14. Hoshino, T., Igari, R.: Quasi-Bayesian Inference for Latent Variable Models with External Information: Application to generalized linear mixed models for biased data. Keio-IES Discussion Paper Series 2017-014, Institute for Economics Studies, Keio University (2017)
15. Jazwinski, A.: Stochastic Processes and Filtering Theory. Ac. Press (1970)
16. Jensen, F.: Bayesian Networks and Decision Graphs. Springer, N.Y. (2001)
17. Kárný, M., Bodini, A., Guy, T., Kracík, J., Nedoma, P., Ruggeri, F.: Fully probabilistic knowledge expression and incorporation. Statistics and Its Interface **7**(4), 503–515 (2014)
18. Kárný, M., Böhm, J., Guy, T., Jirsa, L., Nagy, I., Nedoma, P., Tesař, L.: Optimized Bayesian Dynamic Advising: Theory and Algorithms. Springer, London, UK (2006)
19. Kárný, M., Guy, T.: On support of imperfect Bayesian participants. In: T. Guy, et al (eds.) Decision Making with Imperfect Decision Makers, vol. 28, pp. 29–56. Springer, Int. Syst. Ref. Lib. (2012)
20. Kárný, M., Herzallah, R.: Scalable harmonization of complex networks with local adaptive controllers. IEEE Trans. on SMC: Systems **47**(3), 394–404 (2017)
21. Kasabov, N., Hu, Y.: Integrated optimisation method for personalised modelling and case studies for medical decision support. Int. J. Functional Informatics and Personalised Medicine **3**(3) (2010)

22. Kern-Isberner, G., Lukasiewicz, T.: Special issue on challenges for reasoning under uncertainty, inconsistency, vagueness, and preferences. Künstl. Intell. **31**, 5–8 (2017). DOI https://doi.org/10.1007/s13218-016-0479-z

23. Koopman, R.: On distributions admitting a sufficient statistic. Trans. of Am. Math. Society **39**, 399 (1936)

24. Kracík, J., Kárný, M.: Merging of data knowledge in Bayesian estimation. In: J. Filipe, et al (eds.) Proc. of the 2nd Int. Conf. on Informatics in Control, Automation and Robotics, pp. 229–232. Barcelona (2005)

25. Kuhn, H., Tucker, A.: Nonlinear programming. In: Proc. of 2nd Berkeley Symp., pp. 481–492. Univ. of California Press (1951)

26. Kulhavý, R., Zarrop, M.B.: On a general concept of forgetting. Int. J. of Control **58**(4), 905–924 (1993)

27. Kullback, S., Leibler, R.: On information and sufficiency. Ann Math Stat **22**, 79–87 (1951)

28. van Laere, J.: Challenges for IF performance evaluation in practice. In: 12th Intern. Conf. on Information Fusion, pp. 866 – 873. IEEE, Seattle, WA (2009)

29. Lee, H., Lee, B., Park, K., Elmasri, R.: Fusion techniques for reliable information: A survey. Intern. Journal of Digital Content Technology and its Applications **4**(2), 74–88 (2010)

30. Leevy, J., Khoshgoftaar, T., Bauder, R., Seliya, N.: A survey on addressing high-class imbalance in big data. Journal of Big Data **5**(42) (2018)

31. Meng, T., Jing, X., Yan, Z., Pedrycz, W.: A survey on machine learning for data fusion. Information Fusion (2019). DOI https://doi.org/10.1016/j.inffus.2019.12.001

32. Mine, H., Osaki, S.: Markovian Decision Processes. Elsevier (1970)

33. Nelsen, R.: An Introduction to Copulas. Springer, N.Y. (1999)

34. O'Hagan, A., et al: Uncertain Judgement: Eliciting Experts' Probabilities. J. Wiley (2006)

35. Pearl, J.: Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference. Morgan Kaufman (1988)

36. Peterka, V.: Bayesian system identification. In: P. Eykhoff (ed.) Trends and Progress in System Identification, pp. 239–304. Perg. Press (1981)

37. Quinn, A., Ettler, P., Jirsa, L., Nagy, I., Nedoma, P.: Probabilistic advisory systems for data-intensive applications. Int. J. of Adapt. Control & Signal Proc. **17**(2), 133–148 (2003)

38. Quinn, A., Kárný, M., Guy, T.: Optimal design of priors constrained by external predictors. Int. J. Approximate Reasoning **84**, 150–158 (2017)

39. Rao, M.: Measure Theory and Integration. J. Wiley (1987)

40. Sassani, B., Alahmadi, A., Sharifzadeh, H.: A cluster based collaborative filtering method for improving the performance of recommender systems in e-commerce. In: K. Arai, et al (eds.) Proceedings of the Future Technologies Conference (FTC) 2018, *Advances in Intelligent Systems and Computing*, vol. 881. Springer, Cham (2019)

41. Savage, L.: Foundations of Statistics. Wiley (1954)

42. Scanagatta, M., et al: A survey on Bayesian network structure learning from data. Progress in AI **8**, 425–439 (2019)

43. Shore, J., Johnson, R.: Axiomatic derivation of the principle of maximum entropy & the principle of minimum cross-entropy. IEEE Tran. on Inf. Th. **26**(1), 26–37 (1980)

44. Smith, A., Makov, U.: A quasi-Bayes sequential procedures for mixtures. J. of the Royal Statistical Society **40**(1), 106–112 (1978)

45. Tsai, C., Lai, C., Chiang, M., Yang, L.: Data mining for internet of things: A survey. IEEE Communications Surveys & Tutorials **16**(1), 77–95 (2014)

46. Wang, P., Yang, L., Li, J., Chen, J., Hu, S.: Data fusion in cyber-physical-social systems: State-of-the-art and perspectives. Information Fusion **51**, 42 – 57 (2019)

47. Xu, Z., He, Y., Wang, X.: An overview of probabilistic-based expressions for qualitative decision-making: techniques, comparisons and developments. International Journal of Machine Learning and Cybernetics **1513–1528**, 10 (2019)

48. Zadeh, L.: A fuzzy-algorithmic approach to the definition of complex or imprecise concepts. Systems Theory in the Social Sciences pp. 202–282 (1976)