# AN EFFICIENT ALGORITHM FOR ESTIMATING STATE SEQUENCES IN IMPRECISE HIDDEN MARKOV MODELS

#### JASPER DE BOCK AND GERT DE COOMAN

ABSTRACT. We present an efficient exact algorithm for estimating state sequences from outputs (or observations) in imprecise hidden Markov models (iHMM), where both the uncertainty linking one state to the next, and that linking a state to its output, are represented using coherent lower previsions. The notion of independence we associate with the credal network representing the iHMM is that of epistemic irrelevance. We consider as best estimates for state sequences the (Walley–Sen) maximal sequences for the posterior joint state model conditioned on the observed output sequence, associated with a gain function that is the indicator of the state sequence. This corresponds to (and generalises) finding the state sequence with the highest posterior probability in HMMs with precise transition and output probabilities (pHMMs). We argue that the computational complexity is at worst quadratic in the length of the Markov chain, cubic in the number of states, and essentially linear in the number of maximal state sequences. For binary iHMMs, we investigate experimentally how the number of maximal state sequences depends on the model parameters. We also present a simple toy application in optical character recognition, demonstrating that our algorithm can be used to robustify the inferences made by precise probability models.

# 1. Introduction

In Artificial Intelligence, probabilistic graphical models are becoming an increasingly powerful tool. Amongst these, hidden Markov models (HMMs) are definitely amongst the simplest, and perhaps also amongst the more popular ones.

An important application for HMMs involves finding the *sequence* of (hidden) states with the highest posterior probability after observing a sequence of outputs [12]. For HMMs with precise local transition and emission probabilities, there is a quite efficient dynamic programming algorithm, due to Viterbi [12, 14], for performing this task. For imprecise-probabilistic local models, such as coherent lower previsions, we know of no algorithm in the literature for which the computational complexity comes even close to that of Viterbi's.

In this paper, we take the first steps towards remedying this situation. We describe imprecise hidden Markov models as special cases of credal trees (a special case of credal networks) under epistemic irrelevance in Section 3. We show in particular how we can use the ideas underlying the MePiCTIr<sup>1</sup> algorithm [4], involving independent natural extension and marginal extension, to construct a most conservative joint model from imprecise local transition and emission models. We also derive a number of interesting and useful formulas from that construction.

The results in Section 3 assume basic knowledge of the theory of coherent lower previsions, a generalisation of classical probability that allows for incomplete specification of probabilities. We include a short introduction to this theory in Section 2.

In Section 4 we explain how a sequence of observations leads to (a collection of) socalled maximal state sequences. Finding all of them seems a daunting task at first: it has a search space that grows exponentially in the length of the Markov chain. However, in

Key words and phrases. Hidden Markov model, state sequence estimation, imprecise probabilities, maximality, coherent lower previsions.

<sup>&</sup>lt;sup>1</sup>MePiCTIr: Message Passing in Credal Trees under Irrelevance.

Section 5 we use the basic formulas found in Section 3 to derive an appropriate version of Bellman's [1] Principle of Optimality, which allows for an exponential reduction of the search space. By using a number of additional tricks, we are able in Section 6 to devise the EstiHMM² algorithm, which efficiently constructs all maximal state sequences. We prove in Section 7 that this algorithm is essentially linear in the number of maximal sequences, quadratic in the length of the chain, and cubic in the number of states. We perceive this complexity to be comparable to that of the Viterbi algorithm, especially after realising that the latter makes the simplifying step of resolving ties more or less arbitrarily in order to produce only a single optimal state sequence. This is something we will not allow our algorithm to do, for reasons that should become clear further on.

In Section 8, we consider the special case of binary iHMMs, and investigate experimentally how the number of maximal state sequences depends on the model parameters. We comment on the very interesting structures that emerge, and give them an heuristic explanation.

We show off the algorithm's efficiency in Section 9 by calculating the maximal sequences for a specific iHMM of length 100.

We conclude in Section 10 with a simple toy application in optical character recognition. It demonstrates the advantages of our algorithm and gives a clear indication that the EstiHMM algorithm is able to robustify the existing Viterbi algorithm in an intelligent manner.

In order to make our main argumentation as readable as possible, we have relegated all technical proofs to an appendix.

# 2. Freshening up on coherent lower previsions

We begin with some basic theory of coherent lower previsions. See Ref. [15] for an in-depth study, and Ref. [9] for a recent survey.

Coherent lower previsions are a special type of imprecise probability model. Roughly speaking, whereas classical probability theory assumes that a subject's uncertainty can be represented by a single probability mass function, the theory of imprecise probabilities effectively works with sets of possible probability mass functions, and thereby allows for imprecision as well as indecision to be modelled and represented. For people who are unfamiliar with the theory, looking at it as a way of robustifying the classical theory is perhaps the easiest way to understand and interpret it, and we will use this approach here.

Consider a set  $\mathscr{M}$  of probability mass functions, defined on a discrete set  $\mathscr{X}$  of possible states. With each mass function  $p \in \mathscr{M}$ , we can associate a *linear prevision* (or expectation operator)  $P_p$ , defined on the set  $\mathscr{G}(\mathscr{X})$  of all real-valued maps on  $\mathscr{X}$ . Any  $f \in \mathscr{G}(\mathscr{X})$  is also called a *gamble* on  $\mathscr{X}$ , and  $P_p(f) := \sum_{x \in \mathscr{X}} p(x) f(x)$  is the expected value of f, associated with the probability mass function p. We can now define the *lower prevision*  $\underline{P}_{\mathscr{M}}$  that corresponds with the set  $\mathscr{M}$  as the following *lower envelope* of linear previsions:

$$\underline{P}_{\mathscr{M}}(f) := \inf \{ P_p(f) \colon p \in \mathscr{M} \} \text{ for all gambles } f \text{ in } \mathscr{X}.$$
 (1)

Similarly, we define the *upper prevision*  $\overline{P}_{\mathcal{M}}$  as

$$\overline{P}_{\mathscr{M}}(f) := \sup \left\{ P_p(f) \colon p \in \mathscr{M} \right\} \\
= -\inf \left\{ -P_p(f) \colon p \in \mathscr{M} \right\} = -\inf \left\{ P_p(-f) \colon p \in \mathscr{M} \right\} = -\underline{P}_{\mathscr{M}}(-f) \quad (2)$$

for all gambles f on  $\mathscr{X}$ . We will mostly talk about lower previsions, since it follows from the *conjugacy relation* (2) that the two models are mathematically equivalent.

An *event A* is a subset of the set of possible values  $\mathcal{X}$ :  $A \subseteq \mathcal{X}$ . With such an event, we can associate an *indicator*  $\mathbb{I}_A$ , which is the gamble on  $\mathcal{X}$  that assumes the value 1 on A, and

<sup>&</sup>lt;sup>2</sup>Estimation in imprecise Hidden Markov Models

0 outside A. We call

$$\underline{P}_{\mathscr{M}}(A) := \underline{P}_{\mathscr{M}}(\mathbb{I}_A) = \inf \left\{ \sum_{x \in A} p(x) \colon p \in \mathscr{M} \right\}$$

the *lower probability* of the event A, and similarly  $\overline{P}_{\mathscr{M}}(A) := \overline{P}_{\mathscr{M}}(\mathbb{I}_A)$  its *upper probability*. It can be shown [15] that the functional  $\underline{P}_{\mathscr{M}}$  satisfies the following set of interesting mathematical properties, which define a *coherent lower prevision*:

C1.  $P_{\mathscr{M}}(f) \ge \min f$  for all  $f \in \mathscr{G}(\mathscr{X})$ ,

C2.  $\underline{P}_{\mathscr{M}}(\lambda f) = \lambda \underline{P}_{\mathscr{M}}(f)$  for all  $f \in \mathscr{G}(\mathscr{X})$  and all real  $\lambda \geq 0$ , [non-negative homogeneity] C3.  $\underline{P}_{\mathscr{M}}(f+g) \geq \underline{P}_{\mathscr{M}}(f) + \underline{P}_{\mathscr{M}}(g)$  for all  $f, g \in \mathscr{G}(\mathscr{X})$ . [superadditivity]

Every set of mass functions  $\mathcal{M}$  uniquely defines a coherent lower prevision  $\underline{P}_{\mathcal{M}}$ , but in general the converse does not hold. However, if we limit ourselves to sets of mass functions  $\mathcal{M}$  that are closed and convex—which makes them *credal sets*—they are in a one-to-one correspondence with coherent lower previsions [15]. This implies that we can use the theory of coherent lower previsions as a tool for reasoning with closed convex sets of probability mass functions. From now on, we will no longer explicitly refer to credal sets  $\mathcal{M}$ , but we will simply talk about coherent lower previsions  $\underline{P}$ . It is useful to keep in mind that there always is a unique credal set that corresponds with such a coherent lower prevision:  $\underline{P} = \underline{P}_{\mathcal{M}}$  for some unique credal set  $\mathcal{M}$ , given by  $\mathcal{M} = \left\{ p : (\forall f \in \mathcal{G}(\mathcal{X})) P_p(f) \geq \underline{P}(f) \right\}$ .

A special kind of imprecise model on  $\mathscr{X}$  is the *vacuous* lower prevision. It is a model that represents complete ignorance and therefore has the set of all possible mass functions on  $\mathscr{X}$  as its credal set  $\mathscr{M}$ . It can be shown easily that for every  $f \in \mathscr{G}(\mathscr{X})$ , the corresponding lower prevision is given by  $\underline{P}(f) = \min f$ .

Conditional lower and upper previsions, which are extensions of the classical conditional expectation functionals, can be defined in a similar, intuitively obvious way as lower envelopes associated with sets of conditional mass functions.

Consider a variable X in  $\mathscr{X}$  and a variable Y in  $\mathscr{Y}$ . A conditional lower prevision  $\underline{P}(\cdot|Y)$  on the set  $\mathscr{G}(\mathscr{X})$  of all gambles on  $\mathscr{X}$  is a two-place real-valued function. For any gamble f on  $\mathscr{X}$ ,  $\underline{P}(f|Y)$  is a gamble on  $\mathscr{Y}$ , whose value  $\underline{P}(g|y)$  in  $y \in \mathscr{Y}$  is the lower prevision of g, conditional on the event Y = y. If for any  $y \in \mathscr{Y}$ , the lower prevision  $\underline{P}(\cdot|y)$  is coherent—satisfies conditions C1–C3—then we call the conditional lower prevision  $\underline{P}(\cdot|Y)$  separately coherent. It will sometimes be useful to extend the domain of the conditional lower prevision  $\underline{P}(\cdot|y)$  from  $\mathscr{G}(\mathscr{X})$  to  $\mathscr{G}(\mathscr{X} \times \mathscr{Y})$  by letting  $\underline{P}(f|y) \coloneqq \underline{P}(f(\cdot,y)|y)$  for all gambles f on  $\mathscr{X} \times \mathscr{Y}$ .

If we have a number of conditional lower previsions involving a number of variables, then each of them must be separately coherent, but we also have to make sure that they satisfy a more stringent *joint coherence* requirement. Explaining this in detail would take us too far, but we refer to Ref. [15] for a detailed discussion, with motivation. For our present purposes, it suffices to say that joint coherence is very closely related to making sure that these conditional lower previsions are lower envelopes associated with conditional mass functions that satisfy Bayes's Rule.

For a given lower prevision  $\underline{P}$  on  $\mathscr{G}(\mathscr{X} \times \mathscr{Y})$ , a corresponding conditional lower prevision  $\underline{P}(\cdot|Y)$  that is jointly coherent with  $\underline{P}$  is not uniquely defined. It is however shown in Ref. [10] that it always lies between the so-called natural and regular extensions.

Using *natural extension*, the conditional coherent lower prevision  $\underline{P}(\cdot|Y)$  is defined by  $\underline{P}(f|y) := \max \left\{ \mu \in \mathbb{R} \colon \underline{P}(\mathbb{I}_{\{y\}}[f-\mu]) \geq 0 \right\}$  if  $\underline{P}(\{y\}) > 0$ , and it is vacuous and thus given by  $\underline{P}(f|y) := \min f$  if  $\underline{P}(\{y\}) = 0$ . This is the smallest (most conservative) way of conditioning a lower prevision. If  $\underline{P}(\{y\}) > 0$ , it corresponds to conditioning every probability mass function in the credal set of  $\underline{P}$  on the observation that Y = y and taking the lower envelope of all these conditioned mass functions.

Using *regular extension*, the conditional coherent lower prevision  $\underline{P}(\cdot|Y)$  is defined by  $\underline{P}(f|y) := \max \{ \mu \in \mathbb{R} : \underline{P}(\mathbb{I}_y[f-\mu]) \ge 0 \}$  if  $\overline{P}(\{y\}) > 0$ , and it is vacuous if  $\overline{P}(\{y\}) = 0$ .

This gives us the greatest (most informative) conditional lower prevision that is jointly coherent with the original unconditional lower prevision. It corresponds to taking all mass functions p in the credal set of  $\underline{P}$  for which  $p(y) \neq 0$ , conditioning them on the observation that Y = y and taking their lower envelope.

Natural and regular extension coincide if  $\underline{P}(\{y\}) > 0$  or  $\overline{P}(\{y\}) = 0$  but are different if  $\overline{P}(\{y\}) > \underline{P}(\{y\}) = 0$ . In the latter case, natural extension is vacuous, but regular extension usually remains more informative.

In this introduction, coherent lower previsions were interpreted as an alternative representation for closed and convex sets of probability mass functions. This approach is often adopted by sensitivity analysts and is rather intuitive for people who are used to working in classical probability theory. For the sake of completeness, we mention here that coherent lower previsions can also be given a behavioural interpretation, without using the notion of a probability mass function. The lower prevision  $\underline{P}(f)$  of a gamble  $f \in \mathcal{G}(\mathcal{X})$  can be interpreted as the supremum acceptable buying price that a subject is willing to pay in order to gain the (possibly negative) reward f(x) after the outcome  $x \in \mathcal{X}$  of the experiment has been determined. See Ref. [15] for more information regarding this interpretation.

#### 3. Basic notions

An imprecise hidden Markov model can be depicted using the following probabilistic graphical model:

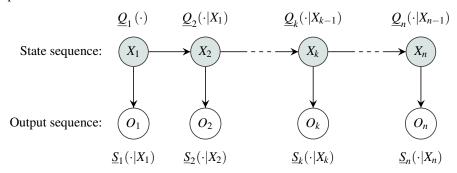


FIGURE 1. Tree representation of a hidden Markov model

Here n is some natural number. The *state variables*  $X_1, \ldots, X_n$  assume values in the respective finite sets  $\mathcal{X}_1, \ldots, \mathcal{X}_n$ , and the *output variables*  $O_1, \ldots, O_n$  assume values in the respective finite sets  $\mathcal{O}_1, \ldots, \mathcal{O}_n$ . We denote generic values of  $X_k$  by  $x_k$ ,  $\hat{x}_k$  or  $z_k$ , and generic values of  $O_k$  by  $O_k$ .

3.1. **Local uncertainty models.** We assume that we have the following local uncertainty models for these variables. For  $X_1$ , we have a marginal lower prevision  $\underline{Q}_1$ , defined on the set  $\mathscr{G}(\mathscr{X}_1)$  of all real-valued maps (or gambles) on  $\mathscr{X}_1$ . For the subsequent states  $X_k$ , with  $k \in \{2, \ldots, n\}$ , we have a conditional lower prevision  $\underline{Q}_k(\cdot|X_{k-1})$  defined on  $\mathscr{G}(\mathscr{X}_k)$ , called a  $transition\ model$ . In order to maintain uniformity of notation, we will also denote the marginal lower prevision  $\underline{Q}_1$  as a conditional lower prevision  $\underline{Q}_1(\cdot|X_0)$ , where  $X_0$  denotes a variable that may only assume a single value, and whose value is therefore certain. For any gamble  $f_k$  in  $\mathscr{G}(\mathscr{X}_k)$ ,  $\underline{Q}_k(f_k|X_{k-1})$  is interpreted as a gamble on  $\mathscr{X}_{k-1}$ , whose value  $\underline{Q}_k(f_k|z_{k-1})$  in any  $z_{k-1} \in \mathscr{X}_{k-1}$  is the lower prevision of the gamble  $f_k(X_k)$ , conditional on  $X_{k-1} = z_{k-1}$ .

In addition, for each output  $O_k$ , with  $k \in \{1, \ldots, n\}$ , we have a conditional lower prevision  $\underline{S}_k(\cdot|X_k)$  defined on  $\mathscr{G}(\mathscr{O}_k)$ , called an *emission model*. For any gamble  $g_k$  in  $\mathscr{G}(\mathscr{O}_k)$ ,  $\underline{S}_k(g_k|X_k)$  is interpreted as a gamble on  $\mathscr{X}_k$ , whose value  $\underline{S}_k(g_k|z_k)$  in any  $z_k \in \mathscr{X}_k$  is the lower prevision of the gamble  $g_k(O_k)$ , conditional on  $X_k = z_k$ .

We take all these local (marginal, transition and emission) uncertainty models to be *separately coherent*. Recall that this simply means that for any  $k \in \{1, ..., n\}$ , the lower prevision  $\underline{Q}_k(\cdot|z_{k-1})$  should be coherent (as an unconditional lower prevision) for every  $z_{k-1} \in \mathcal{X}_{k-1}$  and  $\underline{S}_k(\cdot|z_k)$  should be coherent for every  $z_k \in \mathcal{X}_k$ .

3.2. **Interpretation of the graphical structure.** We will assume that the graphical representation in Figure 1 represents the following irrelevance assessments: *conditional on its mother variable, the non-parent non-descendants of any variable in the tree are epistemically irrelevant to this variable and its descendants.* We say that a variable X is *epistemically irrelevant* to a variable Y if observing X does not affect our beliefs about Y. Mathematically stated in terms of lower previsions: P(f(Y)) = P(f(Y)|x) for all  $f \in \mathcal{G}(Y)$  and all  $x \in \mathcal{X}$ .

Before we go on, it will be useful to introduce some mathematical short-hand notation for describing joint variables in the tree of Figure 1. For any  $1 \le k \le \ell \le n$ , we denote the tuple  $(X_k, X_{k+1}, \dots, X_\ell)$  by  $X_{k:\ell}$ , and the tuple  $(O_k, O_{k+1}, \dots, O_\ell)$  by  $O_{k:\ell}$ .  $X_{k:\ell}$  is a (joint) variable that can assume all values in the set  $\mathscr{X}_{k:\ell} := \times_{r=k}^\ell \mathscr{X}_r$ , and  $O_{k:\ell}$  is a (joint) variable that can assume all values in the set  $\mathscr{O}_{k:\ell} := \times_{r=k}^\ell \mathscr{O}_r$ . Generic values of  $X_{k:\ell}$  are denoted by  $X_{k:\ell}$  or  $X_{k:\ell}$ , and generic values of  $X_{k:\ell}$  by  $X_{k:\ell}$  by  $X_{k:\ell}$ .

**Example 1.** Consider the variable  $X_k$  with mother variable  $X_{k-1}$  in Figure 1. The variables  $X_{1:k-2}$  and  $O_{1:k-1}$  are its non-parent non-descendants, and the variables  $X_{k+1:n}$  and  $O_{k:n}$  its descendants. Our interpretation of the graphical structure of Figure 1 implies that once we know (conditional on) the value  $x_{k_1}$  of  $X_{k-1}$ , additionally learning the values of any of the variables  $X_1, \ldots, X_{k-2}$  and  $X_{k-1}$  and  $X_{k-1}$  will not change our beliefs about  $X_{k:n}$  and  $X_{k:n}$ 

Epistemic irrelevance is weaker than the so-called *strong independence* condition that is usually associated with *credal networks* [2], which is the name usually given to probabilistic graphical models with coherent lower previsions as local uncertainty models. Recent work [4] has shown that using this weaker condition guarantees that an efficient algorithm exists for updating beliefs about a single target node of a credal *tree*, that is essentially linear in the number of nodes in the tree.

3.3. A **joint uncertainty model.** Using the local uncertainty models, we now want to construct a global model: a joint lower prevision  $\underline{P}$  on  $\mathscr{G}(\mathscr{X}_{1:n} \times \mathscr{O}_{1:n})$  for all the variables  $(X_{1:n}, O_{1:n})$  in the tree. This joint lower prevision should (i) be jointly coherent with all the local models; (ii) encode all epistemic irrelevance assessments encoded in the tree; and (iii) be as small, or conservative,<sup>3</sup> as possible. This is a special case of a more general problem for credal trees, discussed and solved in great detail in Ref. [4]. In this section, we summarise the solution for iHMMs and give an heuristic justification for it, but we refer to Ref. [4] for a proof that the joint model we present below is indeed the most conservative lower prevision that is coherent with all the local models and captures all epistemic irrelevance assessments encoded in the tree.

We proceed in a recursive manner, and consider any  $k \in \{1, ..., n\}$ . For any  $x_{k-1} \in \mathcal{X}_{k-1}$ , we consider the smallest coherent joint lower prevision  $\underline{P}_k(\cdot|x_{k-1})$  on  $\mathcal{G}(\mathcal{X}_{k:n} \times \mathcal{O}_{k:n})$  for the variables  $(X_{k:n}, O_{k:n})$  on the iHMM depicted in Figure 2, representing a subtree of the tree represented in Figure 1, with the lower prevision  $\underline{Q}_k(\cdot|x_{k-1})$  acting as the marginal model for the 'first' state variable  $X_k$ . Note that the global model  $\underline{P}$  we are looking for can be identified with the conditional lower prevision  $\underline{P}_1(\cdot|X_0)$ , for the reasons given in Section 3.1.

Our aim is to develop recursive expressions that enable us to construct  $\underline{P}_k(\cdot|X_{k-1})$  out of  $\underline{P}_{k+1}(\cdot|X_k)$ . Using these expressions over and over again will eventually yield the global model  $\underline{P} = \underline{P}_1(\cdot|X_0)$ .

In a first step, we combine the joint model  $\underline{P}_{k+1}(\cdot|X_k)$  for the variables  $(X_{k+1:n}, O_{k+1:n})$ , defined on  $\mathscr{G}(\mathscr{X}_{k+1:n} \times \mathscr{O}_{k+1:n})$ —see the thick dotted lines in Figure 2—, with the local

<sup>&</sup>lt;sup>3</sup>Recall that point-wise smaller lower previsions correspond to larger credal sets.

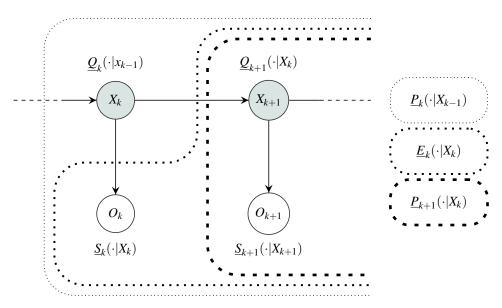


FIGURE 2. Subtree of the iHMM involving the variables  $(X_{k:n}, O_{k:n})$ 

model  $\underline{S}_k(\cdot|X_k)$  for the variable  $O_k$ , defined on  $\mathscr{G}(\mathscr{O}_k)$ . This will lead to a joint model  $\underline{E}_k(\cdot|X_k)$  for the variables  $(X_{k+1:n}, O_{k:n})$ , defined on  $\mathscr{G}(\mathscr{X}_{k+1:n} \times \mathscr{O}_{k:n})$ —see the semi-thick dotted lines in Figure 2. This is trivial for k = n, since we must have that  $\underline{E}_n(\cdot|X_n) = \underline{S}_n(\cdot|X_n)$ .

For  $k \neq n$ , the solution is less obvious. A joint model can be constructed in many different ways, so we will have to impose some conditions. A first condition is that  $\underline{E}_k(\cdot|X_k)$  should be a separately coherent conditional lower prevision that is jointly coherent with the 'marginal' models  $\underline{P}_{k+1}(\cdot|X_k)$  and  $\underline{S}_k(\cdot|X_k)$ . A second, rather obvious, condition is that  $\underline{E}_k(\cdot|X_k)$  should coincide with  $\underline{P}_{k+1}(\cdot|X_k)$  and  $\underline{S}_k(\cdot|X_k)$  on their respective domains. A third condition is that the model should capture the epistemic irrelevance assessments encoded in the tree. In particular these state that, conditional on  $X_k$ , the two variables  $(X_{k+1:n}, O_{k+1:n})$  and  $O_k$  should be *epistemically independent*, or in other words, epistemically irrelevant to one another.

Any model that meets all these conditions is called a (conditionally) *independent product* [5] of  $\underline{P}_{k+1}(\cdot|X_k)$  and  $\underline{S}_k(\cdot|X_k)$ . Generally speaking, such a (conditionally) independent product is not unique. We call the point-wise smallest, most conservative, of all possible (conditionally) independent products, which always exists, the (conditionally) *independent natural extension* [15, 5] of  $\underline{P}_{k+1}(\cdot|X_k)$  and  $\underline{S}_k(\cdot|X_k)$ , and we denote it as  $\underline{P}_{k+1}(\cdot|X_k) \otimes \underline{S}_k(\cdot|X_k)$ .

Summarising,  $\underline{E}_k(\cdot|X_k)$  is given by

$$\underline{E}_{k}(\cdot|X_{k}) := \begin{cases} \underline{S}_{n}(\cdot|X_{n}) & k = n\\ \underline{S}_{k}(\cdot|X_{k}) \otimes \underline{P}_{k+1}(\cdot|X_{k}) & k = n-1,\dots,1 \end{cases}$$
(3)

The (conditionally) independent natural extension and its properties were studied in great detail in Ref. [5]. For the purposes of this paper, it will suffice to recall from that study that—very much like independent products of precise probability models—such

independent natural extensions are factorising, which implies in particular that

$$\underline{E}_{k}(fg|z_{k}) = \underline{E}_{k}(g\underline{E}_{k}(f|z_{k})|z_{k})$$

$$= \underline{S}_{k}(g\underline{P}_{k+1}(f|z_{k})|z_{k})$$

$$= \begin{cases}
\underline{S}_{k}(g|z_{k})\underline{P}_{k+1}(f|z_{k}) & \text{if } \underline{P}_{k+1}(f|z_{k}) \ge 0 \\
\overline{S}_{k}(g|z_{k})\underline{P}_{k+1}(f|z_{k}) & \text{if } \underline{P}_{k+1}(f|z_{k}) \le 0
\end{cases}$$

$$= \underline{\overline{S}}_{k}(g|z_{k}) \odot \underline{P}_{k+1}(f|z_{k}), \tag{4}$$

for all  $z_k \in \mathcal{X}_k$ , all  $f \in \mathcal{G}(\mathcal{X}_{k+1:n} \times \mathcal{O}_{k+1:n})$  and all *non-negative*  $g \in \mathcal{G}(\mathcal{O}_k)$ —we call a gamble non-negative if all its values are. In this expression, the first equality is the actual factorisation property. The second equality holds because  $\underline{E}_k(\cdot|X_k)$  coincides with  $\underline{P}_{k+1}(\cdot|X_k)$  and  $\underline{S}_k(\cdot|X_k)$  on their respective domains. The third equality follows from the conjugacy relation (2) and coherence condition C2, and for the fourth we have used the shorthand notation  $\underline{m} \odot x := \underline{m} \max\{0, x\} + \overline{m} \min\{0, x\}$ . Further on, we will also use the analogous notation  $\underline{m} \overline{n} \odot x := \underline{m} \underline{n} \max\{0, x\} + \overline{m} \overline{n} \min\{0, x\}$ .

In a second and final step, we combine the joint model  $\underline{E}_k(\cdot|X_k)$  for the variables  $(X_{k+1:n}, O_{k:n})$ , defined on  $\mathscr{G}(\mathscr{X}_{k+1:n} \times \mathscr{O}_{k:n})$ , with the local model  $\underline{Q}_k(\cdot|x_{k-1})$  for the variable  $X_k$ , defined on  $\mathscr{G}(\mathscr{X}_k)$ , into the joint model  $\underline{P}_k(\cdot|X_{k-1})$  for the variables  $(X_{k:n}, O_{k:n})$ , defined on  $\mathscr{G}(\mathscr{X}_{k:n} \times \mathscr{O}_{k:n})$ . It has been shown elsewhere [15, 11] that the most conservative coherent way of doing this, is by means of *marginal extension*, also known as the law ot iterated (lower) expectations. This leads to  $\underline{P}_k(\cdot|x_{k-1}) := \underline{Q}_k(\underline{E}_k(\cdot|X_k)|x_{k-1})$ , or, if we now allow  $x_{k-1}$  to range over  $\mathscr{X}_{k-1}$ :

$$\underline{P}_k(\cdot|X_{k-1}) := Q_k(\underline{E}_k(\cdot|X_k)|X_{k-1}). \tag{5}$$

For practical purposes, it is useful to see that this is equivalent with

$$\underline{P}_k(f|X_{k-1}) = \underline{Q}_k \left( \sum_{z_k \in \mathscr{X}_k} \mathbb{I}_{\{z_k\}} \underline{E}_k(f(z_k, X_{k+1:n}, O_{k:n})|z_k) \Big| X_{k-1} \right)$$

for all  $f \in \mathcal{G}(\mathcal{X}_{k:n} \times \mathcal{O}_{k:n})$ . Recall that in this expression, the *indicator*  $\mathbb{I}_{\{z_k\}}$  is a gamble on  $\mathcal{X}_k$  that assumes the value 1 if  $X_k = z_k$  and 0 if  $X_k \neq z_k$ .

3.4. **Interesting lower and upper probabilities.** Without too much trouble,<sup>4</sup>, we can use Equations (3)–(5) to derive the following expressions for a number of interesting lower and upper probabilities:

$$\underline{P}_{k}(\{o_{k:n}\} \times \{z_{k:n}\} | z_{k-1}) = \prod_{i=k}^{n} \underline{S}_{i}(\{o_{i}\} | z_{i}) \underline{Q}_{i}(\{z_{i}\} | z_{i-1})$$
(6)

$$\overline{P}_{k}(\{o_{k:n}\} \times \{z_{k:n}\} | z_{k-1}) = \prod_{i=k}^{n} \overline{S}_{i}(\{o_{i}\} | z_{i}) \overline{Q}_{i}(\{z_{i}\} | z_{i-1})$$
(7)

for all  $z_{k-1} \in \mathscr{X}_{k-1}, z_{k:n} \in \mathscr{X}_{k:n}, o_{k:n} \in \mathscr{O}_{k:n}$  and  $k \in \{1, \dots, n\}$ , and

$$\underline{E}_{k}(\lbrace o_{k:n}\rbrace \times \lbrace z_{k+1:n}\rbrace | z_{k}) = \underline{S}_{k}(\lbrace o_{k}\rbrace | z_{k}) \prod_{i=k+1}^{n} \underline{S}_{i}(\lbrace o_{i}\rbrace | z_{i}) \underline{Q}_{i}(\lbrace z_{i}\rbrace | z_{i-1})$$
(8)

$$\overline{E}_{k}(\{o_{k:n}\} \times \{z_{k+1:n}\} | z_{k}) = \overline{S}_{k}(\{o_{k}\} | z_{k}) \prod_{i=k+1}^{n} \overline{S}_{i}(\{o_{i}\} | z_{i}) \overline{Q}_{i}(\{z_{i}\} | z_{i-1}).$$
(9)

for all  $z_k \in \mathcal{X}_k$ ,  $z_{k+1:n} \in \mathcal{X}_{k+1:n}$ ,  $o_{k:n} \in \mathcal{O}_{k:n}$  and  $k \in \{1, \dots, n\}$ . We will assume throughout that

$$\overline{P}(\{z_{1:n}\} \times \{o_{1:n}\}) > 0 \text{ for all } z_{1:n} \in \mathcal{X}_{1:n} \text{ and } o_{1:n} \in \mathcal{O}_{1:n}$$

<sup>&</sup>lt;sup>4</sup>As an example, we derive Equations (6) and (7) in Appendix A.

or equivalently, that all *local upper previsions are positive*, in the sense that [4]:

$$\overline{Q}_k(\{z_k\}|z_{k-1}) > 0$$
 and  $\overline{S}_k(\{o_k\}|z_k) > 0$   
for all  $z_{k-1} \in \mathcal{X}_{k-1}, z_k \in \mathcal{X}_k, o_k \in \mathcal{O}_k$  and  $k \in \{1, \dots, n\}$ . (10)

This assumption is very weak and not at all restrictive for practical purposes. The imprecise-probabilistic local models are usually constructed by adding some margin of error around a precise model, thereby making all upper transition probabilities positive by construction. We will however allow lower transition probabilities to be zero, which is something that does happen often in practical problems.

**Proposition 1.** The assumption (10) that all local upper previsions are positive implies that  $\overline{P}_k(\{o_{k:n}\}|z_{k-1}) > 0$  and  $\overline{E}_k(\{o_{k:n}\}|z_k) > 0$  for all  $k \in \{1, ..., n\}$ ,  $z_k \in \mathscr{X}_k$ ,  $z_{k-1} \in \mathscr{X}_{k-1}$  and  $o_{k:n} \in \mathscr{O}_{k:n}$ .

## 4. ESTIMATING STATES FROM OUTPUTS

In a hidden Markov model, the states are not directly observable, but the outputs are, and the general aim is to use the outputs to estimate the states. We concentrate on the following problem: Suppose we have observed the output sequence  $o_{1:n}$ , estimate the state sequence  $x_{1:n}$ . We will use an essentially Bayesian approach to do so, but need to allow for the fact that we are working with imprecise rather than precise probability models.

4.1. **Updating the iHMM.** The first step in our approach consists in updating (or conditioning) the joint model  $\underline{P} := \underline{P}_1(\cdot|X_0)$  on the observed outputs  $O_{1:n} = o_{1:n}$ . As mentioned in Section 2, there is no unique coherent way to perform this updating. However, for the particular problem we are solving in this paper, it so happens that it makes no difference which updating method is used, as long as it is coherent. For the time being, we choose to use the least conservative<sup>5</sup> (most informative) coherent updating method, which is *regular extension*. Later on in Section 4.2, we will show that any other coherent updating method yields the same results.

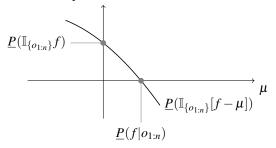
Since it follows from the positivity assumption (10) and Proposition 1 that  $\overline{P}(\{o_{1:n}\}) > 0$ , regular extension leads us to consider the updated lower prevision  $\underline{P}(\cdot|o_{1:n})$  on  $\mathscr{G}(\mathscr{X}_{1:n})$ , given by:

$$\underline{P}(f|o_{1:n}) := \max \left\{ \mu \in \mathbb{R} : \underline{P}(\mathbb{I}_{\{o_{1:n}\}}[f-\mu]) \ge 0 \right\} \text{ for all gambles } f \text{ on } \mathscr{X}_{1:n}. \tag{11}$$

Using the coherence of the joint lower prevision  $\underline{P}$ , it is not hard to prove that when  $\underline{P}(\{o_{1:n}\}) > 0$ ,  $\underline{P}(\mathbb{I}_{\{o_{1:n}\}}[f - \mu])$  is a strictly decreasing and continuous function of  $\mu$ , which therefore has a unique zero (see Lemma 7(i)&(iii) in Appendix A). As a consequence, we have for any  $f \in \mathcal{G}(\mathcal{X}_{1:n})$  that

$$\underline{P}(f|o_{1:n}) \le 0 \Leftrightarrow (\forall \mu > 0)\underline{P}(\mathbb{I}_{\{o_{1:n}\}}[f - \mu]) < 0 \Leftrightarrow \underline{P}(\mathbb{I}_{\{o_{1:n}\}}f) \le 0. \tag{12}$$

In fact, it is not hard to infer from the strictly decreasing and continuous character of  $\underline{P}(\mathbb{I}_{\{o_{1:n}\}}[f-\mu])$  that  $\underline{P}(f|o_{1:n})$  and  $\underline{P}(\mathbb{I}_{\{o_{1:n}\}}f)$  have the same sign. They are either both negative, both positive or both equal to zero; see also the illustration below.



<sup>&</sup>lt;sup>5</sup>The most conservative coherent way yields a vacuous model.

Equation (12) will be of crucial importance further on. However, in general, we want to allow  $\underline{P}(\{o_{1:n}\})$  to be zero (because this may happen if you allow lower transition probabilities to be zero), while requiring that  $\overline{P}(\{o_{1:n}\}) > 0$  (because this follows from the positivity assumption (10) and Proposition 1). This will, generally speaking, invalidate the second equivalence in Equation (12): it turns into an implication only. But, if we limit ourselves to the specific type of gambles on  $\mathscr{X}_{1:n}$  of the form  $f = \mathbb{I}_{\{\hat{x}_{1:n}\}} - \mathbb{I}_{\{x_{1:n}\}}$ , we can still prove the following important theorem.

**Theorem 2.** If all local upper previsions are positive, then  $\underline{P}(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}}])$  and  $\underline{P}(\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}}|o_{1:n})$  have the same sign for all fixed values of  $x_{1:n}, \hat{x}_{1:n} \in \mathscr{X}_{1:n}$  and  $o_{1:n} \in \mathscr{O}_{1:n}$ . They are both positive, both negative or both zero.

4.2. **Maximal state sequences.** The next step now consists in using the posterior model  $\underline{P}(\cdot|o_{1:n})$  to find best estimates for the state sequence  $x_{1:n}$ . On the Bayesian approach, this is usually done by solving a decision-making, or optimisation problem: we associate a gain function  $\mathbb{I}_{\{x_{1:n}\}}$  with every candidate state sequence  $x_{1:n}$ , and select as best estimates those state sequences  $\hat{x}_{1:n}$  that maximise the posterior expected gain, resulting in state sequences with maximal posterior probability.

Here we generalise this decision-making approach towards working with imprecise probability models. The criterion we use to decide which estimates are optimal for the given gain functions is that of (Walley–Sen) *maximality* [13, 15]. Maximality has a number of very desirable properties that make sure it works well in optimisation contexts [6, 8], and it is well-justified from a behavioural point of view, as well as in a robustness approach, as we shall see presently.

We can express a strict preference  $\succ$  between two state sequence estimates  $\hat{x}_{1:n}$  and  $x_{1:n}$  as follows:

$$\hat{x}_{1:n} \succ x_{1:n} \Leftrightarrow \underline{P}(\mathbb{I}_{\{\hat{x}_{1:n}\}} - \mathbb{I}_{\{x_{1:n}\}} | o_{1:n}) > 0.$$

On a behavioural interpretation, this expresses that a subject with lower prevision  $\underline{P}(\cdot|o_{1:n})$  is disposed to pay some strictly positive amount of utility to replace the (gain associated with the) estimate  $x_{1:n}$  with the (gain associated with the) estimate  $\hat{x}_{1:n}$ ; see Ref. [15, Section 3.9] for more details. Alternatively, from a robustness point of view, this expresses that for each conditional mass function  $p(\cdot|o_{1:n})$  in the credal set associated with the updated lower prevision  $\underline{P}(\cdot|o_{1:n})$ , the state sequence  $\hat{x}_{1:n}$  has a posterior probability  $p(\hat{x}_{1:n}|o_{1:n})$  that is *strictly higher* than the posterior probability  $p(x_{1:n}|o_{1:n})$  of the state sequence  $x_{1:n}$ .

The binary relation  $\succ$  thus defined is a strict partial order [an irreflexive and transitive binary relation] on the set of state sequences  $\mathcal{X}_{1:n}$ , and we consider an estimate  $\hat{x}_{1:n}$  to be *optimal* when it is *undominated*, or *maximal*, in this strict partial order:

$$\hat{x}_{1:n} \in \operatorname{opt}(\mathscr{X}_{1:n}|o_{1:n}) \Leftrightarrow (\forall x_{1:n} \in \mathscr{X}_{1:n}) x_{1:n} \not\succ \hat{x}_{1:n} 
\Leftrightarrow (\forall x_{1:n} \in \mathscr{X}_{1:n}) \underline{P}(\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}}|o_{1:n}) \leq 0 
\Leftrightarrow (\forall x_{1:n} \in \mathscr{X}_{1:n}) \underline{P}(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}}]) \leq 0,$$
(13)

where the very useful last equivalence follows from Theorem 2. In summary then, the aim of this paper is to develop an efficient algorithm for finding the set of maximal estimates opt  $(\mathcal{X}_{1:n}|o_{1:n})$ .

Our statement in Section 4.1, that any coherent updating method would yield the same results as regular extension, can now be justified. Since coherent updating is unique if  $\underline{P}(\{o_{1:n}\}) > 0$ , we only need to motivate our statement in the special case that  $\underline{P}(\{o_{1:n}\}) = 0$  and  $\overline{P}(\{o_{1:n}\}) > 0$ .

If we use regular extension to update our model, the optimal estimates are given by Eq. (13). For the special case  $\underline{P}(\{o_{1:n}\}) = 0$  however, we find for all  $x_{1:n} \in \mathcal{X}_{1:n}$  and  $\hat{x}_{1:n} \in \mathcal{X}_{1:n}$  that

$$\underline{P}(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}}]) \leq \underline{P}(\mathbb{I}_{\{o_{1:n}\}}) = \underline{P}(\{o_{1:n}\}) = 0,$$

where the first inequality follows from the monotonicity of coherent lower previsions (as a consequence of C1 and C2). Therefore, we find that if  $\underline{P}(\{o_{1:n}\}) = 0$ , all sequences are optimal, resulting in opt  $(\mathscr{X}_{1:n}|o_{1:n}) = \mathscr{X}_{1:n}$ .

If we use natural extension to update our joint model, the optimal state sequences are still given by Eq. (13), but the final equivalence would no longer hold because it uses Theorem 2, which assumes the use of regular extension to perform updating of the joint model. However, for the special case of  $\underline{P}(\{o_{1:n}\}) = 0$ , natural extension by definition leads to the updated model being equal to the vacuous one. Therefore, we find for all  $x_{1:n} \in \mathcal{X}_{1:n}$  and  $\hat{x}_{1:n} \in \mathcal{X}_{1:n}$  that

$$\underline{P}(\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}} | o_{1:n}) = \min(\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}}) \le 0.$$

This implies that for the special case of  $\underline{P}(\{o_{1:n}\}) = 0$  and  $\overline{P}(\{o_{1:n}\}) > 0$ —identical to what we found for regular extension—natural extension also results in all sequences being optimal, meaning that opt  $(\mathcal{X}_{1:n}|o_{1:n}) = \mathcal{X}_{1:n}$ .

We have thus shown that, in the special case when  $\underline{P}(\{o_{1:n}\}) = 0$  and  $\overline{P}(\{o_{1:n}\}) > 0$ , the set of optimal sequences is the same, regardless of whether we use natural or regular extension to update our joint model. Since every other coherent updating method lies in between those two methods, opt  $(\mathcal{X}_{1:n}|o_{1:n})$  does not depend on the updating method, as long as it is coherent. If  $\underline{P}(\{o_{1:n}\}) > 0$ , coherent updating is unique and thus equal to regular extension, thereby making this result trivial in that case. We can therefore conclude that the results in this paper do not depend on the particular updating method that is chosen, as long as it is coherent.

Instead of looking for the maximal state sequences, one could also use other decision criteria. A first approach that we will not consider here, could consist in trying to find the so-called  $\Gamma$ -maximin state sequences  $\bar{x}_{1:n}$ , which maximise the posterior lower probability:

$$\bar{x}_{1:n} \in \underset{x_1:n \in \mathcal{X}_{1:n}}{\operatorname{argmax}} \underline{P}(\{x_{1:n}\} | o_{1:n})$$

While it is well known that any such  $\Gamma$ -maximin sequence is in particular guaranteed to also be a maximal sequence, finding such  $\Gamma$ -maximin sequences seems to be a much more complicated affair. Of course, once we know all maximal solutions, we could determine which of them are the  $\Gamma$ -maximin solutions by comparing their posterior lower probabilities. As far as we can see, however, calculating these seems no trivial task from a computational point of view.

We expect similar computational difficulties with yet another approach, also not considered here, which consists in finding the so-called *E-admissable* sequences. They are those sequences that maximise the expected gain for at least one conditional mass function  $p(\cdot|o_{1:n})$  in the credal set associated with the updated lower prevision  $\underline{P}(\cdot|o_{1:n})$ . Similarly to the  $\Gamma$ -maximin solutions, the E-admissable ones are also known to be contained within the set of maximal ones that we will be constructing.

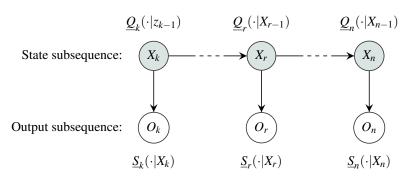
The main reason why our approach is so efficient compared to the other ones, is that we do not have to explicitly calculate the value of lower previsions, but only need to know their sign, thereby allowing us to work directly with the joint model, instead of the updated model.

4.3. **Maximal subsequences.** We shall see below that in order to find the set of maximal estimates, it is useful to consider more general sets of so-called maximal subsequences: for any  $k \in \{1, ..., n\}$  and  $z_{k-1} \in \mathcal{X}_{k-1}$ , we define opt  $(\mathcal{X}_{k:n}|z_{k-1}, o_{k:n})$ :

$$\hat{x}_{k:n} \in \text{opt}(\mathscr{X}_{k:n}|z_{k-1}, o_{k:n}) \Leftrightarrow (\forall x_{k:n} \in \mathscr{X}_{k:n}) \ \underline{P}_{k}(\mathbb{I}_{\{o_{k:n}\}}[\mathbb{I}_{\{x_{k:n}\}} - \mathbb{I}_{\{\hat{x}_{k:n}\}}]|z_{k-1}) \le 0. \quad (14)$$

The interpretation of these sets is immediate: consider the following part of the original iHMM, where we take  $Q_k(\cdot|z_{k-1})$  as the marginal model for the first state  $X_k$ :

<sup>&</sup>lt;sup>6</sup>Private communication from Cassio de Campos.



Then, as we have argued in Section 3.3, the corresponding joint lower prevision on  $\mathscr{G}(\mathscr{X}_{k:n} \times \mathscr{O}_{k:n})$  is precisely  $\underline{P}_k(\cdot|z_{k-1})$ , and if we have a sequence of outputs  $o_{k:n}$ , then opt  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$  is the set of state sequence estimates that are undominated by any other estimate in  $\mathscr{X}_{k:n}$ . It should be clear that the set opt  $(\mathscr{X}_{1:n}|o_{1:n})$  we are eventually looking for, can also be written as opt  $(\mathscr{X}_{1:n}|z_0,o_{1:n})$ .

4.4. **Useful recursion equations.** Fix any k in  $\{1,\ldots,n\}$ . If we look at Equation (14), we see that it will be useful to derive a manageable expression for the lower prevision  $\underline{P}_k(\mathbb{I}_{\{o_{k:n}\}}[\mathbb{I}_{\{x_{k:n}\}} - \mathbb{I}_{\{\hat{x}_{k:n}\}}]|z_{k-1})$ . This can be easily done (see Appendix A) using Equations (3)–(7) together with a few algebraic manipulations. We consider three different cases. If  $\hat{x}_k = x_k$  and  $k \in \{1,\ldots,n-1\}$  then, using the notation introduced in Section 3.3:

$$\underline{P}_{k}(\mathbb{I}_{\{o_{k:n}\}}[\mathbb{I}_{\{x_{k:n}\}} - \mathbb{I}_{\{\hat{x}_{k:n}\}}|z_{k-1}) \\
= \underline{\overline{Q}}_{k}(\{\hat{x}_{k}\}|z_{k-1})\underline{\overline{S}}_{k}(\{o_{k}\}|\hat{x}_{k}) \odot \underline{P}_{k+1}(\mathbb{I}_{\{o_{k+1:n}\}}[\mathbb{I}_{\{x_{k+1:n}\}} - \mathbb{I}_{\{\hat{x}_{k+1:n}\}}|\hat{x}_{k}). \quad (15)$$

If  $\hat{x}_n = x_n$  then

$$\underline{P}_{n}(\mathbb{I}_{\{o_{n}\}}[\mathbb{I}_{\{x_{n}\}} - \mathbb{I}_{\{\hat{x}_{n}\}}|z_{n-1}) = 0.$$
(16)

If  $\hat{x}_k \neq x_k$  and  $k \in \{1, ..., n\}$  then

$$\underline{P}_{k}(\mathbb{I}_{\{o_{k:n}\}}[\mathbb{I}_{\{x_{k:n}\}} - \mathbb{I}_{\{\hat{x}_{k:n}\}}|z_{k-1}) = \underline{Q}_{k}(\mathbb{I}_{\{x_{k}\}}\beta(x_{k:n}) - \mathbb{I}_{\{\hat{x}_{k}\}}\alpha(\hat{x}_{k:n})|z_{k-1}), \tag{17}$$

where we define, for any  $z_{k:n} \in \mathscr{X}_{k:n}$ :

$$\beta(z_{k:n}) := \underline{E}_k(\mathbb{I}_{\{o_{k:n}\}}\mathbb{I}_{\{z_{k+1:n}\}}|z_k) = \underline{S}_k(\{o_k\}|z_k) \prod_{i=k+1}^n \underline{S}_i(\{o_i\}|z_i)\underline{Q}_i(\{z_i\}|z_{i-1})$$
(18)

$$\alpha(z_{k:n}) := \overline{E}_k(\mathbb{I}_{\{o_{k:n}\}} \mathbb{I}_{\{z_{k+1:n}\}} | z_k) = \overline{S}_k(\{o_k\} | z_k) \prod_{i=k+1}^n \overline{S}_i(\{o_i\} | z_i) \overline{Q}_i(\{z_i\} | z_{i-1}).$$
 (19)

For any given sequence of states  $z_{k:n} \in \mathscr{X}_{k:n}$ , the  $\alpha(z_{k:n})$  and  $\beta(z_{k:n})$  can be found by simple backward recursion:

$$\alpha(z_{k:n}) := \alpha(z_{k+1:n})\overline{S}_k(\lbrace o_k \rbrace \vert z_k)\overline{Q}_{k+1}(\lbrace z_{k+1} \rbrace \vert z_k)$$
(20)

$$\beta(z_{k:n}) := \beta(z_{k+1:n}) \underline{S}_k(\{o_k\}|z_k) Q_{k+1}(\{z_{k+1}\}|z_k), \tag{21}$$

for  $k \in \{1, ..., n-1\}$ , and starting from:

$$\alpha(z_{n:n}) = \alpha(z_n) := \overline{S}_n(\lbrace o_n \rbrace \vert z_n) \text{ and } \beta(z_{n:n}) = \beta(z_n) := S_n(\lbrace o_n \rbrace \vert z_n).$$

# 5. The Principle of Optimality

Determining the state sequences in opt  $(\mathscr{X}_{1:n}|o_{1:n})$  directly using Equation (13) clearly has exponential complexity (in the length of the chain). We are now going to take a dynamic programming approach [1] to reducing this complexity by deriving a recursion equation for the sets of optimal (sub)sequences opt  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ .

**Theorem 3** (Principle of Optimality). For  $k \in \{1, ..., n-1\}$ , all  $z_{k-1} \in \mathcal{X}_{k-1}$  and all  $\hat{x}_{k:n} \in \mathcal{X}_{k:n}$ : if  $\underline{Q}_k(\{\hat{x}_k\}|z_{k-1}) > 0$  and  $\underline{S}_k(\{o_k\}|\hat{x}_k) > 0$ , then

$$\hat{x}_{k:n} \in \operatorname{opt}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n}) \Rightarrow \hat{x}_{k+1:n} \in \operatorname{opt}(\mathscr{X}_{k+1:n}|\hat{x}_k,o_{k+1:n}).$$

As an immediate consequence, we find that

$$\operatorname{opt}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n}) \subseteq \operatorname{cand}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n}), \tag{22}$$

with cand  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$  being the set of sequences in  $\mathscr{X}_{k:n}$  that can still be an element of opt  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$  according the the theorem above:

cand  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ 

$$:= \left(\bigcup_{z_k \in \operatorname{Pos}_k(z_{k-1})} z_k \oplus \operatorname{opt}\left(\mathscr{X}_{k+1:n}|z_k, o_{k+1:n}\right)\right) \cup \left(\bigcup_{z_k \notin \operatorname{Pos}_k(z_{k-1})} z_k \oplus \mathscr{X}_{k+1:n}\right). \tag{23}$$

Here  $\oplus$  denotes concatenation of state sequences and the set of states  $\operatorname{Pos}_k(z_{k-1}) \subseteq \mathscr{X}_k$  is defined as

$$z_k \in \operatorname{Pos}_k(z_{k-1}) \Leftrightarrow \underline{Q}_k(\{z_k\}|z_{k-1}) > 0 \text{ en } \underline{S}_k(\{o_k\}|z_k) > 0.$$
 (24)

Equation (23) simplifies to

$$\operatorname{cand}\left(\mathscr{X}_{k:n}|z_{k-1},o_{k:n}\right) = \bigcup_{z_k \in \mathscr{X}_k} z_k \oplus \operatorname{opt}\left(\mathscr{X}_{k+1:n}|z_k,o_{k+1:n}\right)$$
(25)

if all local lower previsions are positive, but this is not generally true in the more general case we are considering here, where only the upper previsions are required to be positive. We also introduce the following notation:

$$\operatorname{cand}_{x_{k:n}}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n}) := \{z_{k:n} \in \operatorname{cand}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n}) : z_{k:s} = x_{k:s}\}$$
for all  $k \in \{1,\ldots,n\}, s \in \{k,\ldots,n\}, z_{k-1} \in \mathscr{X}_{k-1}, x_{k:s} \in \mathscr{X}_{k:s} \text{ and } o_{k:n} \in \mathscr{O}_{k:n}.$ 

## 6. AN ALGORITHM FOR FINDING MAXIMAL STATE SEQUENCES

We now use Equation (22) to devise an algorithm for constructing the set opt  $(\mathcal{X}_{1:n}|o_{1:n})$  of maximal state sequences in a recursive manner.

6.1. **Initial set-up using backward recursion.** We begin by defining a few auxiliary notions. First of all, we consider the thresholds:

$$\theta_k(\hat{x}_k, x_k | z_{k-1}) := \min \left\{ a \ge 0 \colon \underline{Q}_k(\mathbb{I}_{\{x_k\}} - a \mathbb{I}_{\{\hat{x}_k\}} | z_{k-1}) \le 0 \right\} \tag{27}$$

for all  $k \in \{1, ..., n\}$ ,  $z_{k-1} \in \mathscr{X}_{k-1}$  and  $x_k, \hat{x}_k \in \mathscr{X}_k$ .

Next, we define

$$\alpha_k^{\max}(x_k) := \max_{\substack{z_{k:n} \in \mathscr{X}_{k:n} \\ z_1 = y_1}} \alpha(z_{k:n}) \text{ and } \beta_k^{\max}(x_k) := \max_{\substack{z_{k:n} \in \mathscr{X}_{k:n} \\ z_1 = y_1}} \beta(z_{k:n})$$
 (28)

for all  $k \in \{1, ..., n\}$  and  $x_k \in \mathcal{X}_k$ . Using Equations (20)–(21), these can be calculated efficiently using the following backward recursive (dynamic programming) procedure:

$$\alpha_{k}^{\max}(x_{k}) = \max_{z_{k+1} \in \mathcal{X}_{k+1}} \alpha_{k+1}^{\max}(z_{k+1}) \overline{S}_{k}(\{o_{k}\}|x_{k}) \overline{Q}_{k+1}(\{z_{k+1}\}|x_{k})$$

$$= \overline{S}_{k}(\{o_{k}\}|x_{k}) \max_{z_{k+1} \in \mathcal{X}_{k+1}} \alpha_{k+1}^{\max}(z_{k+1}) \overline{Q}_{k+1}(\{z_{k+1}\}|x_{k}), \tag{29}$$

and

$$\beta_{k}^{\max}(x_{k}) = \max_{z_{k+1} \in \mathscr{X}_{k+1}} \beta_{k+1}^{\max}(z_{k+1}) \underline{S}_{k}(\{o_{k}\}|x_{k}) \underline{Q}_{k+1}(\{z_{k+1}\}|x_{k})$$

$$= \underline{S}_{k}(\{o_{k}\}|x_{k}) \max_{z_{k+1} \in \mathscr{X}_{k+1}} \beta_{k+1}^{\max}(z_{k+1}) \underline{Q}_{k+1}(\{z_{k+1}\}|x_{k}), \tag{30}$$

for  $k \in \{1, \dots, n-1\}$ , starting from

$$\alpha_n^{\max}(x_n) = \alpha(x_n) = \overline{S}_n(\{o_n\}|x_n) \text{ and } \beta_n^{\max}(x_n) = \beta(x_n) = \underline{S}_n(\{o_n\}|x_n). \tag{31}$$

Finally, we let

$$\alpha_k^{\text{opt}}(\hat{x}_k|z_{k-1}) := \max_{\substack{x_k \in \mathcal{X}_k \\ x_k \neq \hat{x}_k}} \beta_k^{\max}(x_k) \theta_k(\hat{x}_k, x_k|z_{k-1}), \tag{32}$$

for all  $k \in \{1, ..., n\}$ ,  $z_{k-1} \in \mathcal{X}_{k-1}$  and  $\hat{x}_k \in \mathcal{X}_k$ .

6.2. **Reformulation of the optimality condition.** It turns out that the  $\alpha_k^{\text{opt}}(\hat{x}_k|z_{k-1})$ , calculated by Equation (32), are extremely useful. As proved in Appendix A, they allow us to significantly simplify Equation (14) as follows:

$$\operatorname{opt}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n}) = \left\{ \hat{x}_{k:n} \in \operatorname{cand}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n}) : \alpha(\hat{x}_{k:n}) \ge \alpha_k^{\operatorname{opt}}(\hat{x}_k|z_{k-1}) \right\}, \quad (33)$$

which, for k = n, reduces to

$$\operatorname{opt}(\mathscr{X}_n|z_{n-1},o_n) = \left\{ \hat{x}_n \in \mathscr{X}_n \colon \alpha(\hat{x}_n) \ge \alpha_n^{\operatorname{opt}}(\hat{x}_n|z_{n-1}) \right\}. \tag{34}$$

6.3. A recursive solution method. The aim of the algorithm is to determine the set opt  $(\mathcal{X}_{1:n}|o_{1:n})$  efficiently. We will do so recursively.

For k = n, opt  $(\mathcal{X}_n | z_{n-1}, o_n)$  can be determined in a straightforward manner for every  $z_{n-1} \in \mathcal{X}_{n-1}$  using Criterion (34).

**Example 2.** We consider a simple binary HMM with  $\mathscr{X} = \{0,1\}$ . For k = n, the maximal elements are simply states, which are trivially represented. We could for example find that opt  $(\mathscr{X}_n|0,o_n) = \{0,1\}$  for  $z_{n-1} = 0$ , and opt  $(\mathscr{X}_n|1,o_n) = \{0\}$  for  $z_{n-1} = 1$ .

Next, we let k run backward from n-1 to 1. For each k < n and all  $z_{k-1} \in \mathcal{Z}_{k-1}$ , we first build up the set cand  $(\mathcal{Z}_{k:n}|z_{k-1},o_{k:n})$ , using its definition in Equation (23) and the results of the previous recursion step. This set is then used to determine opt  $(\mathcal{Z}_k|z_{k-1},o_{k:n})$  with Criterion (33).

**Example 3.** We continue the discussion of Example 2. For k = n - 1 and  $z_{n-2} = 0$ , the set cand  $(\mathcal{X}_{n-1:n}|0,o_{n-1:n})$  is constructed using Equation (23). If, for instance  $\operatorname{Pos}_{n-1}(0) = \{0,1\}$ , this reduces to Equation (25) and we find that

$$\begin{aligned} \operatorname{cand}\left(\mathscr{X}_{n-1:n}|0,o_{n-1:n}\right) &= \bigcup_{z_{n-1}\in\{0,1\}} z_{n-1} \oplus \operatorname{opt}\left(\mathscr{X}_{n}|z_{n-1},o_{n}\right) \\ &= 0 \oplus \{0,1\} \cup 1 \oplus \{0\} = \{00,01\} \cup \{10\} = \{00,01,10\}. \end{aligned}$$

Applying Criterion (33) to every element of this set, we find the set opt  $(\mathscr{X}_{n-1:n}|0,o_{n-1:n})$ , which for instance could be equal to  $\{00,10\}$ . For  $z_{n-2}=0$ , an analoguous method can be used.

Continuing in this way, we eventually reach k = 1, which yields the desired set of maximal sequences opt  $(\mathcal{X}_{1:n}|o_{1:n}) = \text{opt}(\mathcal{X}_{1:n}|z_0,o_{1:n})$ .

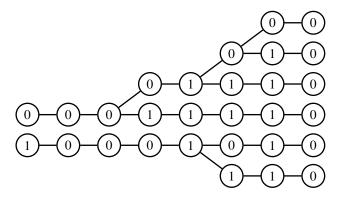
The possible bottleneck in this solution lies in the use of Criterion (33). While this criterion is already much more efficient than the original one, it can still lead to an exponential complexity if the set cand  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$  has a number of elements that is exponential in the length of the considered sequences. We therefore present a method that avoids checking the inequality in Criterion (33) for all elements of cand  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ .

The first trick consists in using an efficient data structure to store the sets of optimal sequences. For k = n, this is simply a list of the elements. For k < n, we could also just list the optimal sequences, but this would imply storing the same information multiple times, since parts of those sequences will be the same. We therefore choose to represent this list of optimal sequences as a collection of tree structures. The way these trees are constructed should be obvious from the following example.

# **Example 4.** Consider the following set of sequences:

 $\{00001000,00001010,00001110,00011110,10001010,10001110\}$ 

By representing this set in this way, useful information gets lost and memory space is waisted. For example, some of these sequences all start out the same way. It would be much more efficient to store such common subsequences only once.



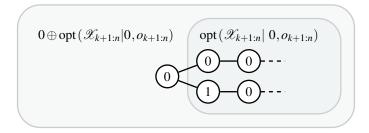
We therefore prefer to represent the above set as the collection of trees depicted above. •

The next step is now to exploit this data structure in order to apply Criterion (33) efficiently. We start by constructing the set cand  $(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$  and representing it in the same type of data structure.

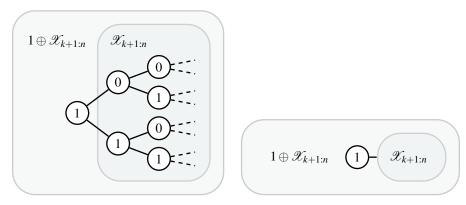
**Example 5.** We consider the set of sequences in Example 4 to be opt  $(\mathscr{X}_{k+1:n}|\ 0, o_{k+1:n})$ , where k = n - 8, since the length of the sequences is 8. Suppose we have already constructed this set in the previous recursion step. Furthermore, for the sake of this example, lets assume that  $0 \in \operatorname{Pos}_{k-1}(0)$  and  $1 \notin \operatorname{Pos}_{k-1}(0)$ . We will now use Equation (23) to construct the set cand  $(\mathscr{X}_{k:n}|0, o_{k:n})$ :

$$\operatorname{cand}\left(\mathscr{X}_{k:n}|0,o_{k:n}\right) = 0 \oplus \operatorname{opt}\left(\mathscr{X}_{k+1:n}|0,o_{k+1:n}\right) \cup 1 \oplus \mathscr{X}_{k+1:n}.$$

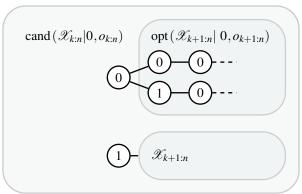
The set cand  $(\mathscr{X}_{k:n}|0,o_{k:n})$  consist of two subsets, which we will construct separately. The subset  $0 \oplus \operatorname{opt}(\mathscr{X}_{k+1:n}|0,o_{k+1:n})$  would normally take quite some effort to compose, since we have to concatenate 0 with each individual element of  $\operatorname{opt}(\mathscr{X}_{k+1:n}|0,o_{k+1:n})$ . However, using our representation, this comes down to adding one node and two links to the already existing data structure for  $\operatorname{opt}(\mathscr{X}_{k+1:n}|0,o_{k+1:n})$ :



Conceptually, we want to represent the set  $1 \oplus \mathscr{X}_{k+1:n}$  as a tree, which would look like the figure below on the left.



However, storing it this way in a computer is a bad idea, as this would mean constructing a complete binary tree, which is exponential in the depth of this tree. We therefore remember that the set of sequences can be represented as a tree, without actually constructing it, as is depicted above on the right.



The set cand  $(\mathscr{X}_{k:n}|0,o_{k:n})$  we are looking for is then trivially constructed by joining the two subsets  $0 \oplus \text{opt}(\mathscr{X}_{k+1:n}|0,o_{k+1:n})$  and  $1 \oplus \mathscr{X}_{k+1:n}$ , as depicted above.

It follows from Equation (33) that the data structure representing opt  $(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$  is contained in the data structure representing cand  $(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$ . All that is now left to do is find this subset in an efficient manner. We present a method that constructs a subset of cand  $(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$ , and will prove that this subset is indeed opt  $(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$ .

cand  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ , and will prove that this subset is indeed opt  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ . We first define  $\alpha_k^{\text{opt}}(z_{k:s}|z_{k-1})$  for every  $k \in \{1,\ldots,n\}$ ,  $s \in \{k,\ldots,n\}$ ,  $z_{k-1} \in \mathscr{X}_{k-1}$  and  $z_{k:s} \in \mathscr{X}_{k:s}$ . If s = k, we let  $\alpha_k^{\text{opt}}(z_{k:k}|z_{k-1}) := \alpha_k^{\text{opt}}(z_k|z_{k-1})$ , defined by Equation (32).  $\alpha_k^{\text{opt}}(z_{k:s}|z_{k-1})$  is then recursively defined by

$$\alpha_k^{\text{opt}}(z_{k:s}|z_{k-1}) = \frac{\alpha_k^{\text{opt}}(z_{k:s-1}|z_{k-1})}{\overline{S}_{s-1}(\{o_{s-1}\}|z_{s-1})\overline{O}_s(\{z_s\}|z_{s-1})} \text{ for every } s \in \{k+1,\dots,n\}.$$
 (35)

**Optimal tree construction.** The following method will select a subset out of a given set cand  $(\mathcal{L}_{k:n}|z_{k-1},o_{k:n})$  constructed using Equation (23).

First, for every  $x_k \in \mathcal{X}_k$ , check whether

$$\alpha_k^{\max}(x_k) \ge \alpha_k^{\text{opt}}(x_k|z_{k-1}). \tag{36}$$

From now on, we will use the generic notation  $\hat{x}_k$  for those  $x_k \in \mathcal{X}_k$  for which this condition is satisfied.

Next, choose an arbitrary  $\hat{x}_k$  and check, for every  $x_{k+1} \in \mathcal{X}_{k+1}$  that has a non-empty set cand  $\hat{x}_{k} \oplus x_{k+1}$  ( $\mathcal{X}_{k:n} | z_{k-1}, o_{k:n}$ ), if the following condition is satisfied:

$$\alpha_{k+1}^{\max}(x_{k+1}) \ge \alpha_k^{\text{opt}}(\hat{x}_k \oplus x_{k+1}|z_{k-1}).$$
 (37)

Notice that  $\alpha_k^{\text{opt}}(\hat{x}_k \oplus x_{k+1}|z_{k-1})$  can be easily calculated using Equation (35), because  $\alpha_k^{\text{opt}}(\hat{x}_k|z_{k-1})$  is already known from the previous recursion step. Denote those  $x_{k+1} \in \mathcal{X}_{k+1}$  for which the inequality (37) is true generically by  $\hat{x}_{k+1}$  and concatenate them with the state  $\hat{x}_k$ , creating a set of state sequences  $\hat{x}_{k:k+1}$ . Do this for every  $\hat{x}_k$  of the previous step and bundle the sets, obtaining a larger set of state sequences  $\hat{x}_{k:k+1}$ .

In a next step, consider an arbitrary  $\hat{x}_{k:k+1}$  and check, for every  $x_{k+2} \in \mathcal{X}_{k+2}$  that has a non-empty set cand  $\hat{x}_{k:k+1} \oplus x_{k+2}$  ( $\mathcal{X}_{k:n} | z_{k-1}, o_{k:n}$ ), if the following condition is satisfied:

$$\alpha_{k+2}^{\max}(x_{k+2}) \ge \alpha_k^{\text{opt}}(\hat{x}_{k:k+1} \oplus x_{k+2}|z_{k-1}).$$
 (38)

As before,  $\alpha_k^{\text{opt}}(\hat{x}_{k:k+1} \oplus x_{k+2}|z_{k-1})$  can be calculated easily using Equation (35), since  $\alpha_k^{\text{opt}}(\hat{x}_{k+1}|z_{k-1})$  has already been calculated in the previous step. Denote those  $x_{k+2} \in \mathcal{X}_{k+2}$  for which the inequality (38) holds generically by  $\hat{x}_{k+2}$  and concatenate them with  $\hat{x}_{k:k+1}$ , creating a set of state sequences  $\hat{x}_{k:k+2}$ . Do this for every  $\hat{x}_{k:k+1}$  from the previous step and bundle the sets to obtain a larger set of state sequences  $\hat{x}_{k:k+2}$ .

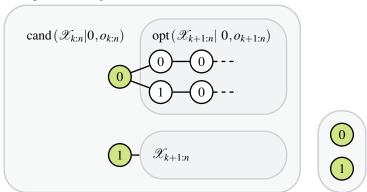
It should be clear that we can go on this way, to eventually end up with a set of sequences  $\hat{x}_{k:n-1}$ . Now consider an arbitrary  $\hat{x}_{k:n-1}$  and check, for every  $x_n \in \mathscr{X}_n$  that has a non-empty set  $\operatorname{cand}_{\hat{x}_{k:n-1} \oplus x_n} (\mathscr{X}_{k:n} | z_{k-1}, o_{k:n})$ , if the following condition holds:

$$\alpha_n^{\max}(x_n) \ge \alpha_k^{\text{opt}}(\hat{x}_{k:n-1} \oplus x_n | z_{k-1}). \tag{39}$$

Denote those  $x_n \in \mathscr{X}_n$  for which this is the case as  $\hat{x}_n$ , and concatenate them with  $\hat{x}_{k:n-1}$ , creating a set of state sequences  $\hat{x}_{k:n}$ . Do this for every  $\hat{x}_{k:n-1}$  from the previous step and bundle the sets to finally obtain a set of state sequences  $\hat{x}_{k:n}$ , which is a subset of the set cand  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$  we started out from.

**Theorem 4.** The subset of cand  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$  that is obtained by using the optimal tree construction is equal to opt  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ .

**Example 6.** We continue with Example 5. Following the optimal tree construction, we start by checking for every  $x_k \in \{0,1\}$  whether  $\alpha_k^{\max}(x_k) \ge \alpha_k^{\mathrm{opt}}(x_k|0)$ . Suppose this is the case. We will symbolise this by giving the corresponding nodes in our representation a green colour, as in the leftmost part of the figure below. It then follows by Theorem 4 that every sequence in opt  $(\mathcal{X}_{k:n}|0,o_{k:n})$  will either start with 0 or 1, since the set of  $\hat{x}_k$  is  $\{0,1\}$ . In this example, this is of course trivial, but if the set of  $\hat{x}_k$  would have been  $\{0\}$ , we would have obtained the non-trivial result that every sequence in opt  $(\mathcal{X}_{k:n}|0,o_{k:n})$  starts with 0. We can represent this partial information about the set opt  $(\mathcal{X}_{k:n}|0,o_{k:n})$  in a trivial way, as in the rightmost part of the figure below.



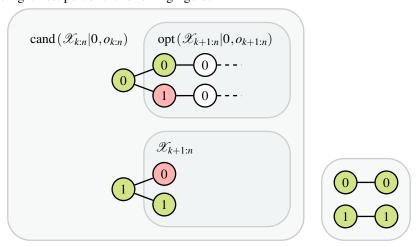
In the next step, we need to check some criteria for every  $\hat{x}_k$  we have found in the previous step. We begin with  $\hat{x}_k = 0$  and start by looking at  $x_{k+1} = 0$ . The set  $\operatorname{cand}_{\hat{x}_k \oplus x_{k+1}}(\mathcal{X}_{k:n}|0,o_{k:n})$  is then  $\operatorname{cand}_{00}(\mathcal{X}_{k:n}|0,o_{k:n})$ , which is simply the subset of sequences in  $\operatorname{cand}(\mathcal{X}_{k:n}|0,o_{k:n})$  that start with 00. In our tree representation of  $\operatorname{cand}(\mathcal{X}_{k:n}|0,o_{k:n})$ , checking whether this set is non-empty comes down to checking if the node  $\hat{x}_k = 0$  has a daughter with value 0. Since

this is indeed the case, we need to check whether  $\alpha_{k+1}^{\max}(0) \ge \alpha_k^{\text{opt}}(\hat{x}_k \oplus x_{k+1}|0) = \alpha_k^{\text{opt}}(00|0)$ . Suppose this criterion is met, then we have found our first subsequence  $\hat{x}_{k:k+1}$ , namely 00. We symbolise this in the figure below by giving the child  $x_{k+1} = 0$  of the node  $\hat{x}_k = 0$  a green colour.

The node  $\hat{x}_k = 0$  also has a daughter  $x_{k+1} = 1$ . If  $\alpha_{k+1}^{\max}(1) < \alpha_k^{\text{opt}}(01|0)$ , this daughter gets coloured red and 01 is not part of the set of sequences  $\hat{x}_{k:k+1}$  we are constructing in this step. By Theorem 4, this also means that none of the elements of  $\text{opt}(\mathcal{L}_{k:n}|0,o_{k:n})$  will start with the subsequence 01.

For  $\hat{x}_k=1$ , we know that the tree representing the sequences in cand  $(\mathscr{X}_{k:n}|0,o_{k:n})$  that start with 1 is a complete tree, which we have not explicitly constructed. This does not create a problem, since we only need that tree to check whether  $\operatorname{cand}_{1\oplus x_{k+1}}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$  is a non-empty set, which is a condition that is trivially met for all  $x_{k+1}\in\mathscr{X}_{k+1}$  because of the completeness of the set  $\operatorname{cand}_1(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ . We are therefore left to check Criterion (37) for  $\hat{x}_k=1$  and every  $x_{k+1}\in\{0,1\}$ . For  $x_{k+1}=0$ , we might for instance find that  $\alpha_{k+1}^{\max}(0)<\alpha_k^{\operatorname{opt}}(10|0)$  and for  $x_{k+1}=1$  we might find that  $\alpha_{k+1}^{\max}(1)\geq\alpha_k^{\operatorname{opt}}(11|0)$ . The results of these checks are summarised in the leftmost part of the figure below. The

The results of these checks are summarised in the leftmost part of the figure below. The corresponding sequences  $\hat{x}_{k:k+1}$ , which by Theorem 4 are the possible starting sequences for the elements of opt  $(\mathcal{X}_{k:n}|0,o_{k:n})$ , can be easily stored and depicted in our tree representation; see the rightmost part of the following figure.



If we keep performing the steps of optimal tree construction in this way, Theorem 4 states that the data structure that is built up while checking all these criteria represents the set opt  $(\mathscr{X}_{k:n}|0,o_{k:n})$ . This set might look like this:

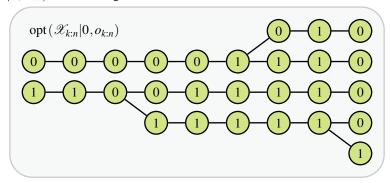


Figure 3 should clarify how this set was constructed. Notice that we have indeed never explicitly constructed the set  $\mathcal{X}_{k+1:n}$  in the tree representation, since every time we reached a red node, the descendants of this node were not constructed.

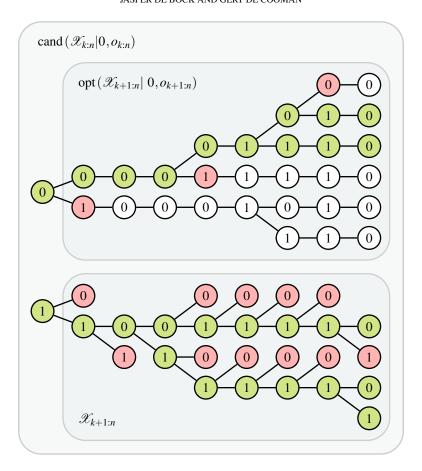


FIGURE 3. Clarification of the construction of cand  $(\mathcal{X}_{k:n}|0,o_{k:n})$ 

6.4. **Additional comments.** All that is needed in order to produce the  $\alpha$ - and  $\beta$ -functions are assessments for the lower and upper transition and emission probabilities:

$$\underline{Q}_k(\{z_k\}|z_{k-1}), \overline{Q}_k(\{z_k\}|z_{k-1}), \underline{S}_k(\{o_k\}|z_k) \text{ and } \overline{S}_k(\{o_k\}|z_k)$$

for all  $k \in \{1, ..., n\}$ ,  $z_{k-1} \in \mathcal{X}_{k-1}$ ,  $z_k \in \mathcal{X}_k$  and  $o_k \in \mathcal{O}_k$ . The most conservative coherent models  $\underline{Q}_k(\cdot|X_{k-1})$  that correspond to such assessments are 2-monotone [3, 7]. Due to their comonotone additivity [7], this implies that:

$$\underline{Q}_k(\mathbb{I}_{\{x_k\}}-a\mathbb{I}_{\{\hat{x}_k\}}|z_{k-1})=\underline{Q}_k(\{x_k\}|z_{k-1})-a\overline{Q}_k(\{\hat{x}_k\}|z_{k-1})$$

for all  $a \ge 0$ , and therefore Equation (27) leads to

$$\theta_k(\hat{x}_k, x_k | z_{k-1}) = \frac{\underline{Q}_k(\{x_k\} | z_{k-1})}{\overline{Q}_k(\{\hat{x}_k\} | z_{k-1})}.$$
(40)

The right-hand side is the smallest possible value of the threshold  $\theta_k(\hat{x}_k, x_k|z_{k-1})$  corresponding to the assessments  $\underline{Q}_k(\{x_k\}|z_{k-1})$  and  $\overline{Q}_k(\{\hat{x}_k\}|z_{k-1})$ , leading to the most conservative inferences, and therefore the largest possible sets of maximal sequences, that correspond to these assessments.

# 7. DISCUSSION OF THE ALGORITHM'S COMPLEXITY

7.1. **Preparatory calculations.** We begin with the preparatory calculations of the quantities in Equations (27)–(32). For the thresholds  $\theta_k(\hat{x}_k, x_k | z_{k-1})$  in Equation (27), the computational complexity is clearly cubic in the number of states, and linear in the number of nodes. Calculating the  $\alpha_k^{\max}(x_k)$  and  $\beta_k^{\max}(x_k)$  in Equations (29) and (30) is linear in the number

of nodes, and quadratic in the number of states. The complexity of finding the  $\alpha_k^{\text{opt}}(\hat{x}_k|z_{k-1})$  in Equation (32) is linear in the number of nodes, and cubic in the number of states.

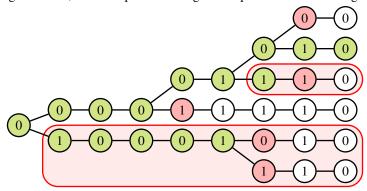
7.2. Complexity of the optimal tree construction. The computational complexity of the optimal tree construction is less trivial. Let us start by noting that this construction essentially consists in repeating the same small step over and over again, namely adding a state  $\hat{x}_s$  to an already constructed  $\hat{x}_{k:s-1}$ .

To perform such a step for a sequence  $\hat{x}_{k:s-1}$ , we first have to check for all  $x_s \in \mathcal{X}_s$  whether  $\operatorname{cand}_{\hat{x}_{k:s-1} \oplus x_s}(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$  is non-empty. This can be done in constant time, since our representation reduces this step to checking whether the node  $x_s$  is a daughter of  $\hat{x}_{s-1}$  in the data structure of  $\operatorname{cand}_{\hat{x}_{k:s-1}}(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$ . Next, for those  $x_s \in \mathcal{X}_s$  for which this is indeed the case, we need to check if  $\alpha_s^{\max}(x_s) \geq \alpha_k^{\operatorname{opt}}(\hat{x}_{k:s-1} \oplus x_s|z_{k-1})$ . Checking those two criteria for every  $x_s \in \mathcal{X}_s$  will from now on be called *performing a search step*, and its complexity is linear in the number of states. Those  $x_s \in \mathcal{X}_s$  that meet both criteria will be noted as  $\hat{x}_s$  and concatenated with  $\hat{x}_{k:s-1}$ .

We will now prove that performing such a search step will always yield at least one  $\hat{x}_s$  that can be concatenated with  $\hat{x}_{k:s-1}$ .

**Theorem 5.** Consider an arbitrary sequence  $\hat{x}_{k:s-1}$  that is created while performing the optimal tree construction, with  $k \in \{1,...,n\}$  and  $s \in \{k,...,n\}$ . Then there is always at least one  $x_s \in \mathcal{X}_s$  for which both  $\operatorname{cand}_{\hat{x}_{k:s-1} \oplus x_s}(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$  is non-empty and the inequality  $\alpha_s^{\max}(x_s) \geq \alpha_t^{\operatorname{opt}}(\hat{x}_{k:s-1} \oplus x_s|z_{k-1})$  holds.

**Example 7.** In our visual representations, this means that every green node will alway have at least one green child, which implies that all green sequences will have length n - k + 1.



The situation depicted above is therefore impossible.

Next, notice that every optimal sequence  $\hat{x}_{k:n}$  yielded by the optimal tree construction is built up by adding extra states  $\hat{x}_s$  to an already constructed sequence  $\hat{x}_{k:s-1}$ , repeating this for s going from k to n. Adding such a state means performing one search step, but Theorem 5 implies that performing a search step also means adding at least one state. Therefore, constructing one maximal sequence  $\hat{x}_{k:n}$  will never take more search steps than the length of this sequence. Since performing one search step is linear in the number of states, constructing one maximal sequence is linear in the length of the sequence and the number of states. Determining the set opt  $(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$  of all maximal sequences will thus be linear in the number of sequences, in the length of the sequences and in the number of states.

7.3. The recursive construction of the solutions. To construct opt  $(\mathcal{X}_{1:n}|o_{1:n})$  recursively, we let k run from n to 1. For a fixed k, we construct the set opt  $(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$  for every  $z_{k-1} \in \mathcal{X}_{k-1}$ , by means of the optimal tree construction. We have already shown that

<sup>&</sup>lt;sup>7</sup>If s = k, we identify  $\hat{x}_{k:s-1} = \hat{x}_{k:k-1}$  with a sequence of length zero.

constructing such a set is linear in the number of sequences, the length of the sequences and the number of states. This means that performing this recursive construction is quadratic in the length of the sequences, quadratic in number of states and roughly speaking<sup>8</sup> linear in the number of maximal sequences.

- 7.4. **General complexity.** The complete algorithm consist of the preparatory calculations and the recursive construction of the solutions. We conclude that it is quadratic in the number of nodes, cubic in the number of states, and roughly speaking linear in the number of maximal sequences.
- 7.5. **Comparison with Viterbi's algorithm.** For precise HMMs, the state sequence estimation problem can be solved very efficiently by the Viterbi algorithm [12, 14], whose complexity is linear in the number of nodes, and quadratic in the number of states. However, this algorithm only emits a single optimal (most probable) state sequence, even in cases where there are multiple (equally probable) optimal solutions: this of course simplifies the problem. If we would content ourselves with giving only a single maximal solution, the ensuing version of our algorithm would have a complexity that is similar to Viterbi's.

So, to allow for a fair comparison between Viterbi's algorithm and ours, we would need to alter Viterbi's algorithm in such a way that it no longer resolves ties arbitrarily, and emits all (equally probable) optimal state sequences. This new version will remain linear in the number of nodes, and quadratic in the number of states, but will also have added complexity. This can easily be seen by noting that emitting the optimal sequences will be linear in the number of them and thus possibly exponential, if all possible solutions would for example be equally probable.

For the complexity for the most time-consuming part of our algorithm (the recursive construction of the solutions), the only difference is this: Viterbi's approach is linear and ours is quadratic in the number of nodes. Where does this difference come from? In iHMMs we have mutually incomparable solutions, whereas in pHMMs the optimal solutions are indifferent, or equally probable. This makes sure that the algorithm for pHMMs requires no forward loops, as is the case in the EstiHMM algorithm, when we perform the optimal tree construction. We believe that this added complexity is a reasonable price to pay for the robustness that working with imprecise-probabilistic models offers.

## 8. Some experiments

While a linear complexity in the number of maximal sequences is probably the best we can hope for, we also see that we will only be able to find all maximal sequences efficiently provided their number is reasonably small. Should it, say, tend to increase exponentially with the length of the chain, then no algorithm, however cleverly designed, could overcome this hurdle. Because this number of maximal sequences is so important, we study its behaviour in more detail. In order to do so, we take a closer look at how this number of maximal sequences depends on the transition probabilities of the model, and how it evolves when we let the imprecision of the local models grow. We shall see that this number displays very interesting behaviour that can be explained, and even predicted to some extent. To allow for easy visualisation, we limit this discussion to binary iHMMs, where both the state and output variables can assume only two possible values, say 0 and 1.

8.1. **Describing a binary stationary iHMM.** We first consider a binary stationary HMM. The (precise) transition probabilities for going from one state to the next are completely determined by numbers in the unit interval: the probability p to go from state 0 to state 0, and the probability q to go from state 1 to state 0. To further pin down the HMM we also need to specify the (marginal) probability m for the first state to be 0, and the two emission

<sup>&</sup>lt;sup>8</sup>For every k, constructing the set opt  $(\mathscr{Z}_{k:n}|z_{k-1},o_{k:n})$  has linear complexity in the number of maximal elements at that stage.

probabilities: the probability r of emitting output 0 from state 0 and the probability s of emitting output 0 from state 1.

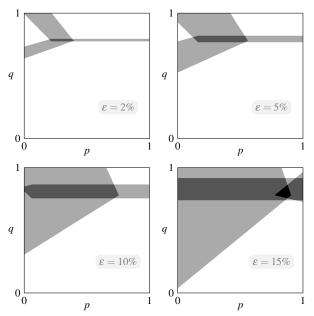
In this binary case, all coherent imprecise-probabilistic models can be found by contamination: taking convex mixtures of precise models, with mixture coefficient  $1-\varepsilon$ , and the vacuous model, with mixture coefficient  $\varepsilon$ , leading to a so-called linear-vacuous model [15]. To simplify the analysis, we let the emission model remain precise, and use the same mixture coefficient  $\varepsilon$  for the marginal and the transition models. As  $\varepsilon$  ranges from zero to one, we then evolve from a precise HMM towards an iHMM with vacuous marginal and transition models (and precise emission models).

8.2. **Explaining the basic ideas using a chain of length two.** We now examine the behaviour of an iHMM of length two, with the following (precise) probabilities fixed:<sup>9</sup>

$$m = 0.1$$
,  $r = 0.8$  and  $s = 0.3$ .

Fixing an output sequence and a value for  $\varepsilon$ , we can use our algorithm to calculate the corresponding numbers of maximal state sequences as p and q range over the unit interval. The results can be represented conveniently in the form of a heat plot. The plots below correspond to the output sequence  $o_{1:2} = 01$ .

The number of maximal state sequences clearly depends on the transition probabilities p and q. In the rather large parts of 'probability space' that are coloured white, we get a single maximal sequence-as we would for HMMs—, but there are contiguous regions where we see a higher number appear. In the present example (binary chain of length two), the highest possible number of maximal sequences is of course four. In the dark grey area, there are three maximal sequences, and two in the light grey regions. The plots show what happens when we let  $\varepsilon$  increase: the



grey areas expand and the number of maximal sequences increases. For  $\varepsilon = 15\%$ , we even find a small area (coloured black) where all four possible state sequences are maximal: locally, due to the relatively high imprecision of our local models, we cannot give any useful robust estimates of the state sequence producing the output sequence  $o_{1:2} = 01$ .

For small  $\varepsilon$ , the areas with more than one maximal state sequence are quite small and seem to resemble strips that narrow down to lines as  $\varepsilon$  tends to zero. This suggests that we should be able to explain at least qualitatively where these areas come from by looking at compatible precise models: the regions where an iHMM produces different maximal (mutually incomparable) sequences, are widened versions of loci of indifference for precise HMMs.

By a *locus of indifference*, we mean the set of (p,q) that correspond to two given state sequences  $x_{1:2}$  and  $\hat{x}_{1:2}$  having equal posterior probability:

$$p(x_{1:2}|o_{1:2}) = p(\hat{x}_{1:2}|o_{1:2}),$$

<sup>&</sup>lt;sup>9</sup>This choice is of course arbitrary. Different values would yield comparable results.

or, provided that  $p(o_{1:2}) > 0$ ,

$$p(x_{1:2}, o_{1:2}) = p(\hat{x}_{1:2}, o_{1:2}).$$

In our example where  $o_{1:2} = 01$ , we find the following expressions for each of the four possible state sequences:

$$p(00,01) = mr(1-r)p$$

$$p(01,01) = mr(1-s)(1-p)$$

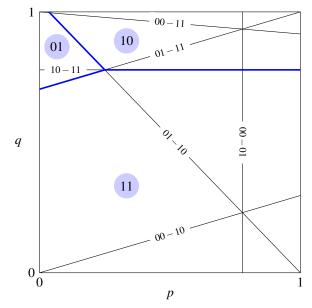
$$p(10,01) = (1-m)s(1-r)q$$

$$p(11,01) = (1-m)s(1-s)(1-q)$$

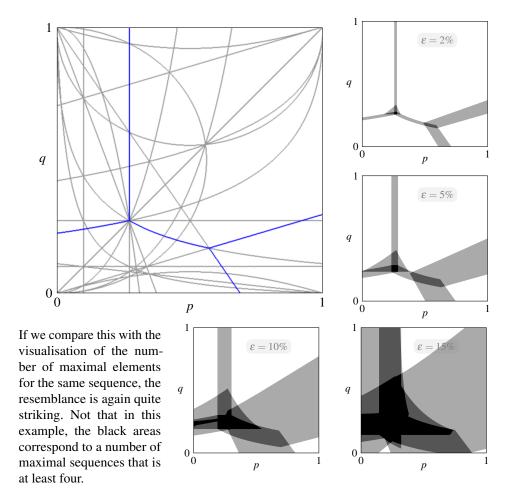
By equating any two of these expressions, we express that the corresponding two state sequences have an equal posterior probability. Since the resulting equations are a function of p and q only, each of these six possible combinations defines a locus of indifference. All of them are depicted as lines in the following figure.

Parts of these loci, depicted in blue (darker and bolder in monochrome versions of this paper) demarcate the three regions where the state sequences 01, 10 and 11 are optimal (have the highest posterior probability).

What happens when the transition models become imprecise? Roughly speaking, nearby values of the original *p* and *q* enter the picture, effectively turning the loci (lines) of indifference into bands of incomparability: the emergence of regions with two and more maximal sequences can be seen to originate from the loci of indifference; compare the figure for these loci with the heat plots given above.



8.3. **Extending the argument to a chain of length three.** For a chain of length three, we can determine the loci of indifference for precise models in a completely analogous manner. If we use the same marginal model and emission model as in the previous example, the resulting lines of indifference for the output sequence 000 look as follows.



# 9. Showing off the algorithm's power

In order to demonstrate that our algorithm is indeed quite efficient, we let it determine the maximal sequences for a random output sequence of length 100.

We consider the same binary stationary HMM as we presented above, but with the following precise marginal and emission probabilities:

$$m = 0.1$$
,  $r = 0.98$  and  $s = 0.01$ .

In practical applications, the probability for an output variable to have the same value as the corresponding hidden state variable is usually quite high, which explains why we have chosen r and s to be close to 1 and to 0, respectively. In contrast with the previous experiments, we do not let the transition probabilities vary, but fix them to the following values:

$$p = 0.6$$
 and  $q = 0.5$ .

The iHMM we use to determine the maximal sequences is then generated by mixing these precise local models with a vacuous one, using the same mixture coefficient  $\varepsilon$  for the marginal, transition and emission models. In Figure 4, we display the five maximal sequences corresponding to the highlighted output sequence, and  $\varepsilon = 2\%$ . Since the emission probabilities were chosen to be quite accurate, it is no surprise that the output sequence itself is one of the maximal sequences. In addition, we have indicated in bold face the state values that differ from the outputs in the output sequence. We see that the model represents more indecision about the values of the state variables as we move further away from the end of the sequence. This is a result of a phenomenon called *dilation*, which—as has been

noted in another paper [4]—tends to occur when inferences in a credal tree proceed from the leaves towards the root.

As for the efficiency of our algorithm: it took about 0.2 seconds to calculate these 5 maximal sequences. <sup>10</sup> The reason why this could be done so fast is that the algorithm is linear in the number of solutions, which in this case is only 5. If we let  $\varepsilon$  grow to for example 5%, the number of maximal sequences for the same output sequence is 764 and these can be determined in about 32 seconds. This demonstrates that the complexity is indeed linear in the number of solutions and that the algorithm can efficiently calculate the maximal sequences even for long output sequences.

#### 10. AN APPLICATION IN OPTICAL CHARACTER RECOGNITION

As a first and simple toy application, we use the EstiHMM algorithm to try and detect mistakes in words. A written word is regarded as a hidden sequence  $x_{1:n}$ , generating an output sequence  $o_{1:n}$  by artificially corrupting the word. In this way, we simulate observation processes that are not perfectly reliable, such as the output of an Optical Character Recognition (OCR) device. This leads to observed output sequences that may contain errors, which we will try and detect. We compare our results with those of the Viterbi algorithm and show that our algorithm offers a more robust solution.

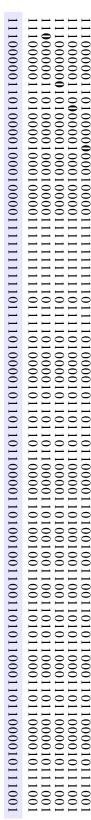
10.1. **Generating the HMM.** A local uncertainty model must be identified for each original and observed letter: a marginal model  $\underline{Q}_1$  for the first letter  $X_1$  of the original word, a transition model  $\underline{Q}_k(\cdot|X_{k-1})$  for the subsequent letters  $X_k$ , with  $k \in \{2, \ldots, n\}$ , and an emission model  $\underline{S}_k(\cdot|X_k)$  for the observed letters  $O_k$ , with  $k \in \{1, \ldots, n\}$ . For the sake of simplicity, we assume stationarity, making the transition and emission models independent of k.

For the identification of the local models of the iHMM, we use the imprecise Dirichlet model (IDM, [16]). For example, for the marginal model  $\underline{Q}_1$ , applying the IDM leads to the following simple identification:

$$\underline{Q}_{1}\left(\left\{x\right\}\right) = \frac{n_{x}}{s + \sum_{z \in \mathcal{X}} n_{z}} \text{ and } \overline{Q}_{1}\left(\left\{x\right\}\right) = \frac{s + n_{x}}{s + \sum_{z \in \mathcal{X}} n_{z}},$$

where  $n_z$  counts the words in the sample text for which the first letter  $X_1 = z$  and s is a (positive real) hyperparameter that expresses the degree of caution in the inferences. In this example, we let s = 2. For the transition and emission models, we can proceed similarly, by counting the transitions of one character to another, respectively in the original word or during the observation process. In this way we obtain lower and upper transition and emission probabilities for singletons, which, as pointed out in Section 6.4, suffice to run the algorithm. Note that if s were chosen to be zero, the local models would become precise and the EstiHMM algorithm would reduce to the Viterbi algorithm (or a version of it that does not resolve ties arbitrarily, see Section 7.5).

For the identification of the local models in the precise HMM, we use a similar but now precise Dirichlet model approach, with a Perks's prior that has the same prior strength s = 2. As an example,



<sup>&</sup>lt;sup>10</sup>Running a Python programme on a 2012 MacBookPro.

for the precise marginal model  $Q_1$ , this leads to the following simple identification:

$$Q_1(\lbrace x\rbrace) = \frac{s/|\mathcal{X}| + n_x}{s + \sum_{z \in \mathcal{X}} n_z},$$

where  $|\mathcal{X}|$  is the number of states.

10.2. **Results.** Let us first discuss a specific example of the difference between the actual results we obtained using the Viterbi and the EstiHMM algorithms, in order to illustrate an important advantage of the latter. OCR software has mistakenly read the Italian word QUANTO as OUANTO. Using a precise model, the Viterbi algorithm does not correct this mistake, as it suggests that the original correct word is DUANTO. The EstiHMM algorithm on the other hand, using an imprecise model, returns CUANTO, DUANTO, FUANTO and QUANTO as maximal (undominated) solutions, including the correct one. Of course we would still have to pick the correct solution out of this set of suggestions—for example by using a dictionary or a human opinion—, but by using the EstiHMM algorithm, we have managed to reduce the search space from all possible five letter words to the much smaller set of four words given above. Notice that the solution of the Viterbi algorithm is included in the maximal solutions EstiHMM returns. One can easily prove that this will always be the case.

To simulate an OCR device, we have artificially corrupted the first 200 words of the first *canto* of Dante's *Divina Commedia*, resulting in 137 correctly read words and 63 words containing errors. We try and correct these errors using both the EstiHMM and the Viterbi algorithm, and compare both approaches. The results are summarised in the following table.

	total number	correct after OCR	wrong after OCR
total number	200 (100%)	137 (68.5%)	63 (31.5%)
Viterbi			
correct solution	157 (78.5%)	132	25
wrong solution	43 (21.5%)	5	38
EstiHMM			
correct solution included	172 (86%)	137	35
correct solution not included	28 (14%)	0	28

For the Viterbi algorithm, the main conclusion is that applying it to the output of the OCR device results in a decreased number of incorrect words. The number of correct words rises from 68.5% to 78.5%. However, the Viterbi algorithm also introduces new errors for 5 correctly read words.

The EstiHMM algorithm manages to suggest the original correct word as one of her solutions in 86% of the cases. Assuming we are able to detect this correct word, the percentage of correct words rises from 68.5% to 86% by applying the EstiHMM algorithm, thereby outperforming the Viterbi algorithm by almost 10%. Secondly, we also notice that the EstiHMM algorithm has never introduced new errors in words that were already correct.

Of course, since the EstiHMM algorithm allows for multiple solutions, instead of a single one, it is no surprise that we manage to increase the amount of times we suggest the correct solution. This would happen even if we added random extra solutions to the solution of the Viterbi algorithm. Giving extra solutions can only be seen as an improvement if this is done smartly. To investigate this, we distinguish between the cases where the EstiHMM algorithm returns a single solution, and those where it returns multiple solutions; and look at how the Viterbi and EstiHMM algorithms compare in those two cases.

The EstiHMM algorithm returned a single solution for 155 of the 200 words. As we have already mentioned above, this single solution will always coincide with the one given by the Viterbi algorithm. The results for the EstiHMM (and Viterbi) algorithms are summarised in the following table.

EstiHMM (single solutions)	total number	correct after OCR	wrong after OCR
total number	155 (100%)	129 (83.2%)	26 (16.8%)
single correct solution	134 (86.5%)	129	5
single wrong solution	21 (13.5%)	0	21

The percentage of words correctly read by the OCR software is now 83.2% instead of the global 68.5%. When the result of the EstiHMM algorithm is a single solution, this serves as an indication that the word we are trying to correct has a fairly high probability of already being correct. We also see that the eventual percentage of correct words is 86.5%, which is only a slight improvement over the 83.2% that were already correct before applying the algorithms.

Next, we look at the remaining 45 words, for which the EstiHMM algorithm returns more than one maximal element. In this case, we do see a significant difference between the results of the Viterbi and the EstiHMM algorithm, since the Viterbi algorithm never returned more than one solution. <sup>11</sup> The results for both algorithms are listed in the following table.

	total number	correct after OCR	wrong after OCR
total number	45 (100%)	8 (17.8%)	37 (82.2%)
<b>EstiHMM (multiple solutions)</b>			
correct solution included	38 (84.4%)	8	30
correct solution not included	7 (15.6%)	0	7
Viterbi			
correct solution	23 (51.1%)	3	20
wrong solution	22 (48.9%)	5	17

A first and very important conclusion to be drawn from this table, is that EstiHMM's being indecisive serves as a rather strong indication that the word we are applying the algorithm to does indeed contain errors: when the EstiHMM algorithm returns multiple solutions, the original word has been incorrectly read by the OCR software in 82.2% of cases.

A second conclusion, related to the first, is that EstiHMM's being indecisive also serves as an indication that the result returned by the Viterbi algorithm is less reliable: the percentage of correct words after applying the Viterbi algorithm has dropped to 51.1%, in contrast with the global percentage of 78.5%. The EstiHMM algorithm, however, still gives the correct word as one of its solutions in 84.4% of cases, which is almost as high as its global percentage of 86%. If the set given by the EstiHMM algorithm contains the correct solution, the Viterbi algorithm manages to pick this correct solution out of the set in 60.5% of cases. We see that the EstiHMM algorithm seems to notice that we are dealing with more difficult words and therefore gives us multiple solutions, between which it cannot decide.

We conclude from this experiment that EstiHMM can be usefully applied to make the results of the Viterbi algorithm more robust, and to gain an appreciation of where it is likely to go wrong. If the EstiHMM algorithm returns multiple solutions, this serves as an indication for robustness issues that would occur if we solved the same problem with the Viterbi algorithm. In that case, EstiHMM returns multiple solutions, between which it cannot decide, whereas the Viterbi algorithm will pick one out of this set in a fairly arbitrary way—depending on the choice of the prior—, thereby increasing the amount of errors made. The advantage of our method is that it detects such robustness issues, leaving us with the option of solving them in different ways. A first method would be to pick the correct word out of the set of possible solutions in some non-arbitrary way. For the current application this could be done using a dictionary or a human expert. Another method for dealing with robustness issues would be to conclude that we need more data in order to build a better model, less sensitive to the choice of the prior. After applying the EstiHMM algorithm again,

<sup>&</sup>lt;sup>11</sup>In theory, the Viterbi algorithm can return multiple indifferent solutions, but in practice it almost never does.

using the new model, we could check whether the robustness issues have been satisfactorily dealt with.

#### 11. CONCLUSIONS

Interpreting the graphical structure of an imprecise hidden Markov model as a credal network under epistemic irrelevance leads to an efficient algorithm for finding the maximal (undominated) state sequences for a given output sequence. Preliminary simulations show that, even for transition models with non-negligible imprecision, the number of maximal elements seems to be reasonably low in fairly large regions of parameter space, with high numbers of maximal elements concentrated in fairly small regions. It remains to be seen whether this observation can be corroborated by a deeper theoretical analysis.

A first and simple toy application clearly shows that the EstiHMM algorithm is able to robustify the results of the Viterbi algorithm. Not only does it reduce the amount of wrong conclusions by giving extra possible solutions, but it does so in an intelligent manner. It adds extra solutions in the specific cases where the Viterbi algorithm has robustness issues, thereby also serving as an indicator of the reliability of the result given by the Viterbi algorithm. An interesting further avenue of research would be to compare the EstiHMM algorithm with other methods that also try to robustify the Viterbi algorithm. Although most of these methods start from a precise model and introduce safety rather than imprecision by for example trying to find the *k* most probable solutions, their practical applications are similar. A comparison of their results with ours could therefore prove to be interesting. We leave this as a topic of future research.

It is not clear to us, at this point, whether ideas similar to the ones we discussed above could be used to derive similarly efficient algorithms for imprecise hidden Markov models whose graphical structure is interpreted as a credal network under strong independence [2]. This could be interesting and relevant, as the more stringent independence condition leads to joint models that are less imprecise, and therefore produce fewer maximal state sequences (although they will be contained in our solutions).

# ACKNOWLEDGEMENTS

Jasper De Bock is a Ph.D. Fellow of the Research Foundation - Flanders (FWO) at Ghent University, and has developed the algorithm described here in the context of his Master's thesis, in close cooperation with Gert de Cooman, who acted as his thesis supervisor. The present article describes the main results of this Master's thesis.

Research by De Cooman has been supported by SBO project 060043 of the IWT-Vlaanderen. This paper has benefitted from discussions with Marco Zaffalon, Alessandro Antonucci, Alessio Benavoli, Cassio de Campos, Erik Quaeghebeur and Filip Hermans. We are grateful to Marco Zaffalon for providing travel funds allowing us to visit IDSIA and discuss practical applications.

## REFERENCES

- [1] Richard Bellman. Dynamic Programming. Princeton University Press, Princeton, 1957.
- [2] Fabio G. Cozman. Credal networks. Artificial Intelligence, 120:199-233, 2000.
- [3] L. M. de Campos, J. F. Huete, and S. Moral. Probability intervals: a tool for uncertain reasoning. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 2:167–196, 1994.
- [4] Gert de Cooman, Filip Hermans, Alessandro Antonucci, and Marco Zaffalon. Epistemic irrelevance in credal nets: the case of imprecise Markov trees. *International Journal of Approximate Reasoning*, 51:1029–1052, 2010.
- [5] Gert de Cooman, Enrique Miranda, and Marco Zaffalon. Independent natural extension. Artificial Intelligence, 2010. Accepted for publication.
- [6] Gert de Cooman and Matthias C. M. Troffaes. Dynamic programming for deterministic discrete-time systems with uncertain gain. *International Journal of Approximate Reasoning*, 39:257–278, 2005.
- [7] Gert de Cooman, Matthias C. M. Troffaes, and Enrique Miranda. n-Monotone exact functionals. Journal of Mathematical Analysis and Applications, 347:143–156, 2008.

- [8] Nathan Huntley and Matthias C. M. Troffaes. Normal form backward induction for decision trees with coherent lower previsions. *Annals of Operations Research*, 2010. Submitted for publication.
- [9] Enrique Miranda. A survey of the theory of coherent lower previsions. *International Journal of Approximate Reasoning*, 48(2):628–658, January 2008.
- [10] Enrique Miranda. Updating coherent lower previsions on finite spaces. Fuzzy Sets and Systems, 160(9):1286–1307, January 2009.
- [11] Enrique Miranda and Gert de Cooman. Marginal extension in the theory of coherent lower previsions. International Journal of Approximate Reasoning, 46(1):188–225, September 2007.
- [12] Lawrence R. Rabiner. A tutorial on HMM and selected applications in speech recognition. Proceedings of the IEEE, 77(2):257–286, February 1989.
- [13] Matthias C. M. Troffaes. Decision making under uncertainty using imprecise probabilities. *International Journal of Approximate Reasoning*, 45(1):17–29, January 2007.
- [14] Andrew J. Viterbi. Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. IEEE Transactions on Information Theory, 13(2):260–269, 1967.
- [15] Peter Walley. Statistical Reasoning with Imprecise Probabilities. Chapman and Hall, London, 1991.
- [16] Peter Walley. Inferences from multinomial data: learning about a bag of marbles. *Journal of the Royal Statistical Society, Series B*, 58:3–57, 1996. With discussion.

### APPENDIX A. PROOFS OF MAIN RESULTS

In this appendix, we justify the formulas (6), (7), (15), (16), (17), (33) and (34); and we give proofs for Proposition 1 and Theorems 2–5. We will frequently use terms such as positive, negative, decreasing and increasing. We therefore start by clarifying what we mean by them. For  $x \in \mathbb{R}$ , we say that x is *positive* if x > 0, *negative* if x < 0, *non-negative* if  $x \ge 0$  and *non-positive* if  $x \le 0$ . We call a real-valued function f defined on  $\mathbb{R}$ :

- (i) increasing if  $(\forall x, y \in \mathbb{R})(x > y \Rightarrow f(x) > f(y))$ ;
- (ii) *decreasing* if  $(\forall x, y \in \mathbb{R})(x > y \Rightarrow f(x) < f(y))$ ;
- (iii) *non-decreasing* if  $(\forall x, y \in \mathbb{R})(x > y \Rightarrow f(x) \ge f(y))$ ;
- (iv) *non-increasing* if  $(\forall x, y \in \mathbb{R})(x > y \Rightarrow f(x) \leq f(y))$ .

*Proof of Equation* (6). For all  $k \in \{1, ..., n\}$ ,  $z_{k-1} \in \mathcal{X}_{k-1}$ ,  $z_{k:n} \in \mathcal{X}_{k:n}$  and  $o_{k:n} \in \mathcal{O}_{k:n}$  we infer from Equation (5) that

$$\begin{split} \underline{P}_k(\mathbb{I}_{\{z_{k:n}\}}\mathbb{I}_{\{o_{k:n}\}}|z_{k-1}) &= \underline{Q}_k(\underline{E}_k(\mathbb{I}_{\{z_{k:n}\}}\mathbb{I}_{\{o_{k:n}\}}|X_k)|z_{k-1}) \\ &= \underline{Q}_k\bigg(\sum_{x_k \in \mathscr{X}_k} \mathbb{I}_{\{x_k\}}\underline{E}_k(\mathbb{I}_{\{z_k\}}(x_k)\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k:n}\}}|x_k)\Big|z_{k-1}\bigg) \\ &= \underline{Q}_k(\mathbb{I}_{\{z_k\}}\underline{E}_k(\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k:n}\}}|z_k)|z_{k-1}). \end{split}$$

Since  $\underline{E}_k(\mathbb{I}_{\{z_{k+1},r_k\}}\mathbb{I}_{\{o_{k,r_k}\}}|z_k) \geq 0$  by C1, we see that C2 transforms the above into

$$= \underline{Q}_{k}(\mathbb{I}_{\{z_{k}\}}|z_{k-1})\underline{E}_{k}(\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k:n}\}}|z_{k}),$$

which can be reformulated as

$$\begin{split} &= \underline{Q}_k(\mathbb{I}_{\{z_k\}}|z_{k-1})\underline{S}_k(\mathbb{I}_{\{o_k\}}|z_k)\underline{P}_{k+1}(\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k+1:n}\}}|z_k) \\ &= \underline{Q}_k(\{z_k\}|z_{k-1})\underline{S}_k(\{o_k\}|z_k)\underline{P}_{k+1}(\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k+1:n}\}}|z_k), \end{split}$$

if we take into account Equation (4), since  $\underline{P}_{k+1}(\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k+1:n}\}}|z_k)\geq 0$  by C1. Repeating these steps again and again eventually yields Equation (6):

$$\underline{P}_k(\mathbb{I}_{\{z_{k:n}\}}\mathbb{I}_{\{o_{k:n}\}}|z_{k-1}) = \prod_{i=k}^n \underline{Q}_i(\{z_i\}|z_{i-1})\underline{S}_i(\{o_i\}|z_i).$$

In the last step, for k = n, we have used the equality  $\underline{E}_n(\{o_n\}|z_n) = \underline{S}_n(\{o_n\}|z_n)$ , which follows from Equation (3).

*Proof of Equation* (7). For all  $k \in \{1, ..., n\}$ ,  $z_{k-1} \in \mathcal{X}_{k-1}$ ,  $z_{k:n} \in \mathcal{X}_{k:n}$  and  $o_{k:n} \in \mathcal{O}_{k:n}$  we infer from conjugacy and Equation (5) that

$$\begin{split} \overline{P}_{k}(\mathbb{I}_{\{z_{k:n}\}}\mathbb{I}_{\{o_{k:n}\}}|z_{k-1}) &= -\underline{P}_{k}(-\mathbb{I}_{\{z_{k:n}\}}\mathbb{I}_{\{o_{k:n}\}}|z_{k-1}) \\ &= -\underline{Q}_{k}(\underline{E}_{k}(-\mathbb{I}_{\{z_{k:n}\}}\mathbb{I}_{\{o_{k:n}\}}|X_{k})|z_{k-1}) \\ &= -\underline{Q}_{k}\left(\sum_{x_{k}\in\mathscr{X}_{k}}\mathbb{I}_{\{x_{k}\}}\underline{E}_{k}(-\mathbb{I}_{\{z_{k}\}}(x_{k})\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k:n}\}}|x_{k})\Big|z_{k-1}\right) \\ &= -\underline{Q}_{k}(\mathbb{I}_{\{z_{k}\}}\underline{E}_{k}(-\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k:n}\}}|z_{k})|z_{k-1}) \\ &= -\underline{Q}_{k}(-\mathbb{I}_{\{z_{k}\}}(-\underline{E}_{k}(-\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k:n}\}}|z_{k})))|z_{k-1}). \end{split}$$

Since  $-\underline{E}_k(-\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k:n}\}}|z_k) = \overline{E}_k(\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k:n}\}}|z_k) \geq 0$  by conjugacy and Lemma 6, we see that C2 and Equation (2) transform the above into

$$\begin{split} &= - \left( -\underline{E}_k(-\mathbb{I}_{\{z_{k+1:n}\}} \mathbb{I}_{\{o_{k:n}\}} | z_k) \right) \underline{Q}_k(-\mathbb{I}_{\{z_k\}} | z_{k-1}) \\ &= -\overline{Q}_k(\mathbb{I}_{\{z_k\}} | z_{k-1}) \underline{E}_k(-\mathbb{I}_{\{z_{k+1:n}\}} \mathbb{I}_{\{o_{k:n}\}} | z_k), \end{split}$$

which can be reformulated as

$$\begin{split} &= -\overline{Q}_{k}(\mathbb{I}_{\{z_{k}\}}|z_{k-1})\overline{S}_{k}(\mathbb{I}_{\{o_{k}\}}|z_{k})\underline{P}_{k+1}(-\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k+1:n}\}}|z_{k}) \\ &= \overline{Q}_{k}(\mathbb{I}_{\{z_{k}\}}|z_{k-1})\overline{S}_{k}(\mathbb{I}_{\{o_{k}\}}|z_{k})\overline{P}_{k+1}(\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k+1:n}\}}|z_{k}) \\ &= \overline{Q}_{k}(\{z_{k}\}|z_{k-1})\overline{S}_{k}(\{o_{k}\}|z_{k})\overline{P}_{k+1}(\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k+1:n}\}}|z_{k}), \end{split}$$

using conjugacy and Equation (4), since  $\underline{P}_{k+1}(-\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k+1:n}\}}|z_k) \leq 0$ . This last inequality is true because we know that  $\underline{P}_{k+1}(-\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k+1:n}\}}|z_k) = -\overline{P}_{k+1}(\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k+1:n}\}}|z_k)$  by conjugacy and that  $\overline{P}_{k+1}(\mathbb{I}_{\{z_{k+1:n}\}}\mathbb{I}_{\{o_{k+1:n}\}}|z_k) \ge 0$  by Lemma 6. Repeating the steps above again and again, eventually yields Equation (7):

$$\overline{P}_k(\mathbb{I}_{\{z_{k:n}\}}\mathbb{I}_{\{o_{k:n}\}}|z_{k-1}) = \prod_{i=-k}^n \overline{Q}_i(\{z_i\}|z_{i-1})\overline{S}_i(\{o_i\}|z_i).$$

In the last step, for k = n, we have used the equality  $\overline{E}_n(\{o_n\}|z_n) = \overline{S}_n(\{o_n\}|z_n)$ , which follows from Equation (3) and conjugacy.

**Lemma 6.** Consider a coherent lower prevision  $\underline{P}$  on  $\mathscr{G}(\mathscr{X})$ . Then  $\min f \leq \underline{P}(f) \leq \overline{P}(f) \leq \overline{$  $\max f$  for all  $f \in \mathcal{G}(\mathcal{X})$  and  $\underline{P}(f) = \overline{P}(\mu) = \mu$  for all  $\mu \in \mathbb{R}$ .

*Proof.* We prove the inequalities in min  $f \leq \underline{P}(f) \leq \overline{P}(f) \leq \max f$  one by one. The first one is the same as C1. It follows by C3 that  $\underline{P}(f-f) \ge \underline{P}(f) + \underline{P}(-f)$  and, since we know by C2 that  $\underline{P}(0) = 0\underline{P}(0) = 0$ , this implies that  $\underline{P}(f) \le -\underline{P}(-f) = \overline{P}(f)$ , using conjugacy for the last equality. For the gamble -f, C1 yields that  $\min -f \leq \underline{P}(-f)$  which implies that  $\max f = -\min -f \geq -\underline{P}(-f) = \overline{P}(f)$ .

To conclude,  $\underline{P}(f) = \overline{P}(\mu) = \mu$  follows by applying these inequalities for  $f = \mu$ .

Proof of Proposition 1. Observe that

$$\overline{P}_k(\mathbb{I}_{\{o_{k:n}\}}|x_{k-1}) = \overline{P}_k\bigg(\mathbb{I}_{\{o_{k:n}\}} \sum_{z_{k:n} \in \mathscr{X}_{k:n}} \mathbb{I}_{\{z_{k:n}\}} \Big| x_{k-1}\bigg) \ge \overline{P}_k\bigg(\mathbb{I}_{\{o_{k:n}\}} \mathbb{I}_{\{z_{k:n}^*\}} \Big| x_{k-1}\bigg) > 0,$$

where  $z_{k:n}^*$  is any element of  $\mathscr{X}_{k:n}$ . The equality follows from  $\sum_{z_{k:n}} \mathscr{Z}_{k:n} \mathbb{I}_{\{z_{k:n}\}} = 1$ , the first inequality from Lemma 8(ii), and the second one from the positivity assumption (10) and Equation (7).

In the same way, we can easily prove that

$$\overline{E}_k(\{o_{k:n}\}|x_k) = \overline{E}_k\left(\mathbb{I}_{\{o_{k:n}\}} \sum_{z_{k+1:n} \in \mathscr{X}_{k+1:n}} \mathbb{I}_{\{z_{k+1:n}\}} \Big| x_k\right) \ge \overline{E}_k\left(\mathbb{I}_{\{o_{k:n}\}} \mathbb{I}_{\{z_{k+1:n}^*\}} \Big| x_k\right) > 0.$$

This time, we have used the positivity assumption (10) and Equation (9) for the last inequal-

*Proof of Theorem 2.* Consider the function  $\rho$  defined by  $\rho(\mu) \coloneqq \underline{P}(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}}]$  $\mu$ ]) for all real  $\mu$ . It follows from Equation (11) that  $\underline{P}(\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}} | o_{1:n})$  is  $\rho$ 's rightmost zero, and we also know that  $\rho(0) = \underline{P}(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}}])$ .  $\rho$  is non-increasing and continuous by Lemma 7(i), and has at least one zero by Lemma 7(ii). Hence, if  $\rho(0) > 0$ , then  $\rho$  has at least one positive zero and  $\underline{P}(\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}}|o_{1:n}) > 0$ . If  $\rho(0) < 0$ , then  $\rho$  has only negative zeroes and  $\underline{P}(\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}} | o_{1:n}) < 0$ . Hence, proving the theorem comes down to proving that  $\rho(0) = 0$  implies that  $\rho(\varepsilon) < 0$  for all  $\varepsilon > 0$ , since this in turn implies that  $\underline{P}(\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}} | o_{1:n}) = 0$ . We now prove this implication. We consider two different

The case  $x_1 = \hat{x}_1$ . For any real  $\varepsilon > 0$ :

$$\rho(\varepsilon) = \underline{P}(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}} - \varepsilon]) 
= \underline{Q}_{1}(\underline{E}_{1}(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}} - \varepsilon]|X_{1})) 
= \underline{Q}_{1}\left(\mathbb{I}_{\{x_{1}\}}\underline{E}_{1}(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{2:n}\}} - \mathbb{I}_{\{\hat{x}_{2:n}\}} - \varepsilon]|x_{1}) + \sum_{z_{1} \neq x_{1}} \mathbb{I}_{\{z_{1}\}}\underline{E}_{1}(-\varepsilon\mathbb{I}_{\{o_{1:n}\}}|z_{1})\right).$$
(41)

The coefficients  $\underline{E}_1(-\varepsilon \mathbb{I}_{\{o_{1:n}\}}|z_1)$  can be written as  $-\varepsilon \overline{E}_1(\{o_{1:n}\}|z_1)$  by conjugacy and C2, which makes them negative, decreasing functions of  $\varepsilon$ , since  $\overline{E}_1(\{o_{1:n}\}|z_1) > 0$  by the positivity assumption (10) and Proposition 1.

For the coefficient  $\underline{E}_1(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{2:n}\}} - \mathbb{I}_{\{\hat{x}_{2:n}\}} - \varepsilon]|x_1)$ , we consider two possible cases. If  $\underline{E}_1(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{2:n}\}} - \mathbb{I}_{\{\hat{x}_{2:n}\}}]|x_1) > 0$ , we know that  $\underline{E}_1(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{2:n}\}} - \mathbb{I}_{\{\hat{x}_{2:n}\}} - \varepsilon]|x_1)$  is a decreasing function of  $\varepsilon$  by Lemma 7(vi). Therefore, the argument of  $\underline{\mathcal{Q}}_1$  in Equation (41) decreases pointwise in  $\varepsilon$ , which by Lemma 8(i) implies that  $\rho(\varepsilon)$  is a decreasing function of  $\varepsilon$  and therefore  $\rho(\varepsilon) < \rho(0) = 0$ .

If, on the other hand,  $\underline{E}_1(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{2:n}\}} - \mathbb{I}_{\{\hat{x}_{2:n}\}}]|x_1) \le 0$ , we know by Lemma 8(ii) that  $\underline{E}_1(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{2:n}\}} - \mathbb{I}_{\{\hat{x}_{2:n}\}} - \varepsilon]|x_1) \leq 0$ , implying that

$$\begin{split} & \rho(\varepsilon) \leq \underline{Q}_1 \bigg( \sum_{z_1 \neq x_1} \mathbb{I}_{\{z_1\}} \underline{E}_1 (-\varepsilon \mathbb{I}_{\{o_{1:n}\}} | z_1) \bigg) \\ & \leq \underline{Q}_1 \left( \mathbb{I}_{\{z_{1*}\}} \underline{E}_1 (-\varepsilon \mathbb{I}_{\{o_{1:n}\}} | z_{1*}) \right) = -\varepsilon \overline{E}_1 (\{o_{1:n}\} | z_{1*}) \overline{Q}_1 \{z_{1*}\} < 0. \end{split}$$

In this expression,  $z_{1*}$  is an arbitrary  $z_1 \neq x_1$ . The first two inequalities are due to Lemma 8(ii). Conjugacy and C2 yield the equality and the last inequality is a consequence of the positivity assumption (10) and Proposition 1. Also in this case, therefore, we find that  $\rho(\varepsilon) < 0$ .

The case  $x_1 \neq \hat{x}_1$ . For any real  $\varepsilon > 0$ :

$$\rho(\varepsilon) = \underline{P}(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}} - \varepsilon]) 
= \underline{Q}_{1}(\underline{E}_{1}(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{1:n}\}} - \mathbb{I}_{\{\hat{x}_{1:n}\}} - \varepsilon]|X_{1})) 
= \underline{Q}_{1}\left(\mathbb{I}_{\{x_{1}\}}\underline{E}_{1}(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{2:n}\}} - \varepsilon]|x_{1}) + \mathbb{I}_{\{\hat{x}_{1}\}}\underline{E}_{1}(\mathbb{I}_{\{o_{1:n}\}}[-\mathbb{I}_{\{\hat{x}_{2:n}\}} - \varepsilon]|\hat{x}_{1}) 
+ \sum_{z_{1} \neq x_{1},\hat{x}_{1}} \mathbb{I}_{\{z_{1}\}}\underline{E}_{1}(-\varepsilon\mathbb{I}_{\{o_{1:n}\}}|z_{1})\right)$$
(42)

In the proof for the case  $x_1 = \hat{x}_1$ , we have already shown that the coefficients  $\underline{E}_1(-\varepsilon \mathbb{I}_{\{o_{1:n}\}}|z_1)$ are negative, decreasing functions of  $\varepsilon$ . Together with Lemma 8(ii), this implies that  $\underline{E}_1(\mathbb{I}_{\{o_{1:n}\}}[-\mathbb{I}_{\{\hat{x}_{2:n}\}}-\varepsilon]|\hat{x}_1) \leq \underline{E}_1(-\varepsilon\mathbb{I}_{\{o_{1:n}\}}|\hat{x}_1) < 0, \text{ which in turn by Lemma 7(vii) implies}$ that  $\underline{E}_1(\mathbb{I}_{\{o_{1:n}\}}[-\mathbb{I}_{\{\hat{x}_{2:n}\}}-\varepsilon]|\hat{x}_1)$  is a decreasing function of  $\varepsilon$ . All that is left to consider is the coefficient  $\underline{E}_1(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{2:n}\}} - \varepsilon]|x_1)$ . There are two possibilities.

If  $\underline{E}_1(\mathbb{I}_{\{o_{1:n}\}}\mathbb{I}_{\{x_{2:n}\}}|x_1) > 0$ , then Lemma 7(vi) implies that  $\underline{E}_1(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{2:n}\}} - \varepsilon]|x_1)$  is a decreasing function of  $\varepsilon$ . Therefore, the argument of  $\underline{Q}_1$  in Equation (42) decreases

pointwise in  $\varepsilon$ , which by Lemma 8(i) implies that  $\rho(\varepsilon)$  is a decreasing function of  $\varepsilon$  and therefore  $\rho(\varepsilon) < \rho(0) = 0$ .

If, on the other hand,  $\underline{E}_1(\mathbb{I}_{\{o_{1:n}\}}\mathbb{I}_{\{x_{2:n}\}}|x_1) = 0$ , then we know that  $\underline{E}_1(\mathbb{I}_{\{o_{1:n}\}}[\mathbb{I}_{\{x_{2:n}\}} - \varepsilon]|x_1) \le 0$  by Lemma 8(ii), implying that

$$\begin{split} \rho(\varepsilon) &\leq \underline{Q}_1(\mathbb{I}_{\{\hat{x}_1\}}\underline{E}_1(\mathbb{I}_{\{o_{1:n}\}}[-\mathbb{I}_{\{\hat{x}_{2:n}\}} - \varepsilon]|\hat{x}_1)) \\ &\leq Q_1(\mathbb{I}_{\{\hat{x}_1\}}\underline{E}_1(-\varepsilon\mathbb{I}_{\{o_{1:n}\}}|\hat{x}_1)) = -\varepsilon\overline{E}_1(\{o_{1:n}\}|\hat{x}_1)\overline{Q}_1(\{\hat{x}_1\}) < 0. \end{split}$$

The first two inequalities follow from Lemma 8(ii). Conjugacy and C2 yield the equality, and the last inequality is a consequence of the positivity assumption (10) and Proposition 1. Also in this case, then, we find that  $\rho(\varepsilon) < 0$ .

**Lemma 7.** Let  $\underline{P}$  be a coherent lower prevision on  $\mathscr{G}(\mathscr{X})$ . For any  $f \in \mathscr{G}(\mathscr{X})$  and  $y \in \mathscr{Y}$ , consider the real-valued map  $\rho$  defined on  $\mathbb{R}$  by  $\rho(\mu) := \underline{P}(\mathbb{I}_{\{y\}}[f-\mu])$  for all real  $\mu$ . Then the following statements hold:

- (i)  $\rho$  is non-increasing, concave and continuous.
- (ii)  $\rho$  has at least one zero.
- (iii) If  $\underline{P}(\{y\}) > 0$ , then  $\rho$  is decreasing and has a unique zero.
- (iv) If  $\overline{P}(\{y\}) = 0$ , then  $\rho$  is identically zero.
- (v) If  $\underline{P}(\{y\}) = 0$  and  $\overline{P}(\{y\}) > 0$ , then  $\rho$  is zero on  $(-\infty, \underline{P}(f|y)]$ , and negative and decreasing on  $(\underline{P}(f|y), +\infty)$ .
- (vi) If  $\rho(a) > 0$  for some a, then  $\rho$  is decreasing and has a unique zero.
- (vii) If  $\rho$  is negative on an interval (a,b), then it is also decreasing on (a,b).

*Proof.* We start by proving (i). It follows directly from Lemma 8(ii) that  $\rho$  is non-increasing in  $\mu$ . Now consider  $\mu_1$  and  $\mu_2$  in  $\mathbb{R}$  and  $0 \le \lambda \le 1$ .  $\rho$  is concave because

$$\begin{split} \rho(\lambda \mu_1 + (1 - \lambda)\mu_2) &= \underline{P}(\mathbb{I}_{\{y\}}[f - (\lambda \mu_1 + (1 - \lambda)\mu_2)]) \\ &= \underline{P}(\lambda \mathbb{I}_{\{y\}}[f - \mu_1] + (1 - \lambda)\mathbb{I}_{\{y\}}[f - \mu_2]) \\ &\geq \underline{P}(\lambda \mathbb{I}_{\{y\}}[f - \mu_1]) + \underline{P}((1 - \lambda)\mathbb{I}_{\{y\}}[f - \mu_2]) \\ &= \lambda \underline{P}(\mathbb{I}_{\{y\}}[f - \mu_1]) + (1 - \lambda)\underline{P}(\mathbb{I}_{\{y\}}[f - \mu_2]) \\ &= \lambda \rho(\mu_1) + (1 - \lambda)\rho(\mu_2), \end{split}$$

where the inequality follows from C3 and the subsequent step is due to C2. To prove that  $\rho(\mu)$  is continuous, consider any  $\mu_1$  and  $\mu_2$  in  $\mathbb{R}$ , then we see that

$$\begin{split} \rho(\mu_2) &= \underline{P}(\mathbb{I}_{\{y\}}[f - \mu_2]) = \underline{P}(\mathbb{I}_{\{y\}}[f - \mu_1 + (\mu_1 - \mu_2)]) \\ &= \underline{P}(\mathbb{I}_{\{y\}}[f - \mu_1] + \mathbb{I}_{\{y\}}(\mu_1 - \mu_2)) \geq \underline{P}(\mathbb{I}_{\{y\}}[f - \mu_1]) + \underline{P}(\mathbb{I}_{\{y\}}(\mu_1 - \mu_2)) \\ &= \rho(\mu_1) - \underline{\overline{P}}(\{y\}) \odot (\mu_2 - \mu_1), \end{split}$$

where the inequality follows from C3, and the last equality is due to conjugacy and C2. Hence  $|\rho(\mu_1) - \rho(\mu_2)| \le |\mu_2 - \mu_1|\overline{P}(\{y\})$ , which proves that  $\rho$  is Lipschitz continuous, and therefore also continuous.

To prove (ii), notice that  $\rho(\min f) = \underline{P}(\mathbb{I}_{\{y\}}[f - \min f]) \geq \underline{P}(\mathbb{I}_{\{y\}}[\min f - \min f]) = 0$  and  $\rho(\max f) = \underline{P}(\mathbb{I}_{\{y\}}[f - \max f])E \leq \underline{P}(\mathbb{I}_{\{y\}}[\max f - \max f]) = 0$ . The inequalities are a consequence of Lemma 8(ii), and the last equalities follow from Lemma 6. Since  $\rho(\mu)$  is continuous, this implies the existence of a zero between  $\min f$  and  $\max f$ .

Property (iii) can be proved by considering  $\mu_1$  and  $\mu_2$  in  $\mathbb{R}$  with  $\mu_2 > \mu_1$ . If  $\underline{P}(\{y\}) > 0$ , we see that  $\rho$  is decreasing, since

$$\begin{split} \rho(\mu_1) &= \underline{P}(\mathbb{I}_{\{y\}}[f - \mu_1]) = \underline{P}(\mathbb{I}_{\{y\}}[f - \mu_2 + (\mu_2 - \mu_1)]) \\ &= \underline{P}(\mathbb{I}_{\{y\}}[f - \mu_2] + \mathbb{I}_{\{y\}}(\mu_2 - \mu_1)) \ge \underline{P}(\mathbb{I}_{\{y\}}[f - \mu_2]) + \underline{P}(\mathbb{I}_{\{y\}}(\mu_2 - \mu_1)) \\ &= \rho(\mu_2) + (\mu_2 - \mu_1)\underline{P}(\{y\}) > \rho(\mu_2), \end{split}$$

where the first inequality follows from C3 and the last equality from C2. We know by (ii) that  $\rho$  has at least one zero, which must be unique because  $\rho$  is decreasing.

To prove (iv), first note that  $\overline{P}(\{y\}) = 0$  also implies  $\underline{P}(\{y\}) = 0$ , because of Lemma 6. Now fix  $\mu$  in  $\mathbb{R}$  and choose a and b in  $\mathbb{R}$  such that

$$a < \min\{0, \min\{f - \mu\}\} \le \max\{0, \max\{f - \mu\}\} < b.$$

Then at the same time  $\rho(\mu) = \underline{P}(\mathbb{I}_{\{y\}}[f-\mu]) \geq \underline{P}(\mathbb{I}_{\{y\}}a) = a\overline{P}(\{y\}) = 0$  and  $\rho(\mu) = \underline{P}(\mathbb{I}_{\{y\}}[f-\mu]) \leq \underline{P}(\mathbb{I}_{\{y\}}b) = b\underline{P}(\{y\}) = 0$ , using Lemma 8(ii), C2 and conjugacy. We conclude that  $\rho(\mu) = 0$  for any  $\mu$  in  $\mathbb{R}$ .

The proof of (v) starts by noticing that  $\rho(\mu) \geq 0$  for  $\mu \in (-\infty, \underline{P}(f|y)]$  and  $\rho(\mu) < 0$  for  $\mu \in (\underline{P}(f|y), +\infty)$ , due to the definition of  $\underline{P}(f|y)$  (see Equation (11)), and the fact that  $\rho$  is non-increasing by (i). In the proof of (iv), we have already shown that  $\rho$  is non-positive if  $\underline{P}(\{y\}) = 0$ , which allows us to conclude that  $\rho(\mu) = 0$  for  $\mu \in (-\infty, \underline{P}(f|y)]$ . We are left to prove that  $\rho$  is decreasing on the interval  $(\underline{P}(f|y), +\infty)$ . We will do so by contradiction. Suppose that  $\rho$  is not decreasing on that interval, then there are  $\mu_1$  and  $\mu_2$  in this interval, such that  $\mu_2 > \mu_1$  and  $0 > \rho(\mu_2) \geq \rho(\mu_1)$ . Since  $\rho$  is zero on  $(-\infty, \underline{P}(f|y))$ , we can also choose  $\mu_0 < \mu_1$  such that  $\rho(\mu_0) = 0$ . The existence of such  $\rho_0$ ,  $\rho(\mu_0) = 0$  and  $\rho(\mu_0) = 0$ . The existence of such  $\rho(\mu_0) = 0$  and  $\rho(\mu_0) = 0$  are contradicts the concavity of  $\rho$ , established by (i).

To prove (vi), observe that  $\overline{P}(\{y\}) \ge \underline{P}(\{y\}) \ge 0$  by Lemma 6. This implies that the three cases considered in (iii), (iv) and (v) are exhaustive and mutually exclusive. If there is an a for which  $\rho(a) > 0$ , we can only have the case considered in (iii), which implies that  $\rho$  is decreasing and has a unique zero.

It now only remains to prove (vii). By repeating the argument in the proof of (vi), we see that  $\rho$  is negative on an interval (a,b), only the cases considered in (iii) and (v) can obtain. For (iii),  $\rho$  is decreasing on its entire domain. For (v),  $\rho$  is definitely decreasing on (a,b).

**Lemma 8.** Consider a coherent lower prevision  $\underline{P}$  on  $\mathscr{G}(\mathscr{X})$  and two gambles  $f,g \in \mathscr{G}(\mathscr{X})$ .

- (i) If f(x) > g(x) for all  $x \in \mathcal{X}$ , then  $\underline{P}(f) > \underline{P}(g)$ .
- (ii) If  $f(x) \ge g(x)$  for all  $x \in \mathcal{X}$ , then  $\underline{P}(f) \ge \underline{P}(g)$ .

*Proof.* We start with (i). If f-g is pointwise positive, then  $\min(f-g)>0$  and therefore  $\underline{P}(f-g)\geq \min(f-g)>0$ , using C1 for the first inequality. It follows from C3 that  $\underline{P}(f)=\underline{P}((f-g)+g)\geq \underline{P}(f-g)+\underline{P}(g)$ , and therefore that  $\underline{P}(f)-\underline{P}(g)\geq \underline{P}(f-g)>0$ , whence indeed  $\underline{P}(f)>\underline{P}(g)$ .

The proof for (ii) is analogous, but now we only have that  $\min(f-g) \ge 0$ , implying that  $\underline{P}(f) - \underline{P}(g) \ge \underline{P}(f-g) \ge \min(f-g) \ge 0$ .

Proof of Equation (15). Let  $\Delta[x_{k:n}, \hat{x}_{k:n}] := \mathbb{I}_{\{o_{k:n}\}}[\mathbb{I}_{\{x_{k:n}\}} - \mathbb{I}_{\{\hat{x}_{k:n}\}}]$ . Since  $k \in \{1, \dots, n-1\}$  and  $\hat{x}_k = x_k$ , this implies that

$$\begin{split} \Delta[x_{k:n}, \hat{x}_{k:n}] &= \mathbb{I}_{\{o_{k:n}\}} [\mathbb{I}_{\{x_{k:n}\}} - \mathbb{I}_{\{\hat{x}_{k:n}\}}] \\ &= \mathbb{I}_{\{o_{k}\}} \mathbb{I}_{\{x_{k}\}} \mathbb{I}_{\{o_{k+1:n}\}} [\mathbb{I}_{\{x_{k+1:n}\}} - \mathbb{I}_{\{\hat{x}_{k+1:n}\}}] \\ &= \mathbb{I}_{\{o_{k}\}} \mathbb{I}_{\{x_{k}\}} \Delta[x_{k+1:n}, \hat{x}_{k+1:n}], \end{split}$$

which in turn implies that

$$\begin{split} \underline{P}_k(\Delta[x_{k:n}, \hat{x}_{k:n}] | z_{k-1}) &= \underline{Q}_k(\underline{E}_k(\mathbb{I}_{\{o_k\}} \mathbb{I}_{\{x_k\}} \Delta[x_{k+1:n}, \hat{x}_{k+1:n}] | X_k) | z_{k-1}) \\ &= \underline{Q}_k(\mathbb{I}_{\{x_k\}} \underline{E}_k(\mathbb{I}_{\{o_k\}} \Delta[x_{k+1:n}, \hat{x}_{k+1:n}] | x_k) | z_{k-1}) \\ &= \underline{\overline{Q}}_k(\{x_k\} | z_{k-1}) \odot \underline{E}_k(\mathbb{I}_{\{o_k\}} \Delta[x_{k+1:n}, \hat{x}_{k+1:n}] | x_k) \\ &= \underline{\overline{Q}}_k(\{x_k\} | z_{k-1}) \underline{\overline{S}}_k(\{o_k\} | x_k) \odot \underline{P}_{k+1}(\Delta[x_{k+1:n}, \hat{x}_{k+1:n}] | x_k), \end{split}$$

proving Equation (15). The first equality follows from Equation (5). The second equality holds because  $\mathbb{I}_{\{x_k\}}(z_k) = 0$  for all  $z_k \neq x_k$ , implying that  $\underline{E}_k(\mathbb{I}_{\{o_k\}}\mathbb{I}_{\{x_k\}}\Delta[x_{k+1:n},\hat{x}_{k+1:n}]|X_k) =$ 

 $\mathbb{I}_{\{x_k\}}\underline{E}_k(\mathbb{I}_{\{o_k\}}\Delta[x_{k+1:n},\hat{x}_{k+1:n}]|x_k)$ . The third equality is follows from conjugacy and C2, and the last one follows from Equation (4).

*Proof of Equation* (16). Since  $\hat{x}_n = x_n$ , Lemma 6 yields:

$$\underline{P}_n(\mathbb{I}_{\{o_n\}}[\mathbb{I}_{\{x_n\}} - \mathbb{I}_{\{\hat{x}_n\}}|z_{n-1}) = \underline{P}_n(\mathbb{I}_{\{o_n\}}[\mathbb{I}_{\{x_n\}} - \mathbb{I}_{\{x_n\}}|z_{n-1}) = \underline{P}_n(0|z_{n-1}) = 0. \qquad \Box$$

*Proof of Equation* (17). If  $k \in \{1, ..., n\}$  and  $\hat{x}_k \neq x_k$ , then

$$\begin{split} \underline{P}_{k}(\mathbb{I}_{\{o_{k:n}\}}[\mathbb{I}_{\{x_{k:n}\}} - \mathbb{I}_{\{\hat{x}_{k:n}\}}|z_{k-1}) \\ &= \underline{Q}_{k}(\underline{E}_{k}(\mathbb{I}_{\{o_{k:n}\}}[\mathbb{I}_{\{x_{k:n}\}} - \mathbb{I}_{\{\hat{x}_{k:n}\}}]|X_{k})|z_{k-1}) \\ &= \underline{Q}_{k}(\mathbb{I}_{\{x_{k}\}}\underline{E}_{k}(\mathbb{I}_{\{o_{k:n}\}}\mathbb{I}_{\{x_{k+1:n}\}}|x_{k}) + \mathbb{I}_{\{\hat{x}_{k}\}}\underline{E}_{k}(-\mathbb{I}_{\{o_{k:n}\}}\mathbb{I}_{\{\hat{x}_{k+1:n}\}}|\hat{x}_{k})|z_{k-1}) \\ &= \underline{Q}_{k}(\mathbb{I}_{\{x_{k}\}}\underline{E}_{k}(\mathbb{I}_{\{o_{k:n}\}}\mathbb{I}_{\{x_{k+1:n}\}}|x_{k}) - \mathbb{I}_{\{\hat{x}_{k}\}}\overline{E}_{k}(\mathbb{I}_{\{o_{k:n}\}}\mathbb{I}_{\{\hat{x}_{k+1:n}\}}|\hat{x}_{k})|z_{k-1}) \\ &= \underline{Q}_{k}(\mathbb{I}_{\{x_{k}\}}\beta(x_{k:n}) - \mathbb{I}_{\{\hat{x}_{k}\}}\alpha(\hat{x}_{k:n})|z_{k-1}), \end{split}$$

proving Equation (17). The reasons why all these equalities hold, are analogous to the ones given in the proof of Equation (15).  $\Box$ 

*Proof of Theorem 3.* Fix  $k \in \{1, \dots, n-1\}$ ,  $z_{k-1} \in \mathscr{X}_{k-1}$  and  $\hat{x}_{k:n} \in \mathscr{X}_{k:n}$ . We assume that  $\hat{x}_{k+1:n} \notin \operatorname{opt}(\mathscr{X}_{k+1:n}|\hat{x}_k, o_{k+1:n})$  and then show that  $\hat{x}_{k:n} \notin \operatorname{opt}(\mathscr{X}_{k:n}|z_{k-1}, o_{k:n})$ . It follows from the assumption that  $\underline{P}_{k+1}(\mathbb{I}_{\{o_{k+1:n}\}}[\mathbb{I}_{\{x_{k+1:n}\}} - \mathbb{I}_{\{\hat{x}_{k+1:n}\}}|\hat{x}_k) > 0$  for some  $x_{k+1:n} \in \mathscr{X}_{k+1}$ . Now prefix this state sequence  $x_{k+1:n}$  with the state  $\hat{x}_k$  to form the state sequence  $x_{k:n}$ , implying that  $x_k = \hat{x}_k$ . We then infer from Equation (15) that

$$\begin{split} & \underline{P}_{k}(\mathbb{I}_{\{o_{k:n}\}}[\mathbb{I}_{\{x_{k:n}\}} - \mathbb{I}_{\{\hat{x}_{k:n}\}}|z_{k-1}) \\ & = \underline{Q}_{k}(\{\hat{x}_{k}\}|z_{k-1})\underline{S}_{k}(\{o_{k}\}|\hat{x}_{k})\underline{P}_{k+1}(\mathbb{I}_{\{o_{k+1:n}\}}[\mathbb{I}_{\{x_{k+1:n}\}} - \mathbb{I}_{\{\hat{x}_{k+1:n}\}}|\hat{x}_{k}) > 0, \\ & \text{which tells us that indeed } \hat{x}_{k:n} \notin \text{opt}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n}). \end{split}$$

*Proof of Equations* (33) *and* (34). First, we consider k = n. For every  $z_{n-1} \in \mathcal{X}_{n-1}$ , we determine opt  $(\mathcal{X}_n|z_{n-1},o_n)$  as the set of those elements  $\hat{x}_n$  of  $\mathcal{X}_n$  for which

$$(\forall x_n \in \mathscr{X}_n \setminus \{\hat{x}_n\}) Q_n(\mathbb{I}_{\{x_n\}} \beta_n^{\max}(x_n) - \mathbb{I}_{\{\hat{x}_n\}} \alpha(\hat{x}_n) | z_{n-1}) \le 0,$$

as this condition is equivalent to the optimality condition (14) for k = n, taking into account Equations (16), (17) and (31). We now show that this condition is also equivalent to

$$(\forall x_n \in \mathcal{X}_n \setminus \{\hat{x}_n\}) \alpha(\hat{x}_n) \ge \beta_n^{\max}(x_n) \theta_n(\hat{x}_n, x_n | z_{n-1}), \tag{43}$$

To see this, we consider two different cases. For those  $x_n$  for which  $\beta_n^{\max}(x_n) = 0$ , the inequalities  $\underline{Q}_n(\mathbb{I}_{\{x_n\}}\beta_n^{\max}(x_n) - \mathbb{I}_{\{\hat{x}_n\}}\alpha(\hat{x}_n)|z_{n-1}) \leq 0$  and  $\alpha(\hat{x}_n) \geq \beta_n^{\max}(x_n)\theta_n(\hat{x}_n,x_n|z_{n-1})$  are both trivially satisfied since  $\alpha(\hat{x}_n) = \overline{S}_n(\{o_n\}|\hat{x}_n) > 0$  by the positivity assumption (10). If  $\beta_n^{\max}(x_n) > 0$ , both inequalities are equivalent because of C2 and Equation (27):

$$\begin{split} \underline{Q}_n(\mathbb{I}_{\{x_n\}}\beta_n^{\max}(x_n) - \mathbb{I}_{\{\hat{x}_n\}}\alpha(\hat{x}_n)|z_{n-1}) &\leq 0 \Leftrightarrow \underline{Q}_n\bigg(\mathbb{I}_{\{x_n\}} - \mathbb{I}_{\{\hat{x}_n\}}\frac{\alpha(\hat{x}_n)}{\beta_n^{\max}(x_n)}\bigg|z_{n-1}\bigg) \leq 0 \\ &\Leftrightarrow \frac{\alpha(\hat{x}_n)}{\beta_n^{\max}(x_n)} \geq \theta_n(\hat{x}_n, x_n|z_{n-1}) \\ &\Leftrightarrow \alpha(\hat{x}_n) \geq \beta_n^{\max}(x_n)\theta_n(\hat{x}_n, x_n|z_{n-1}). \end{split}$$

Using Equation (32), Equation (43) can now be reformulated as  $\alpha(\hat{x}_n) \ge \alpha_n^{\text{opt}}(\hat{x}_n|z_{n-1})$ , which completes the proof of the equivalence.

Next, we consider any  $k \in \{1, \ldots, n-1\}$ . Fix  $z_{k-1} \in \mathscr{X}_{k-1}$ , then we must determine opt  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ . We know from the Principle of Optimality (22) that we can limit the candidate optimal sequences  $\hat{x}_{k:n}$  to the set cand  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ . Consider any such  $\hat{x}_{k:n}$ , then we must check for any  $x_{k:n} \in \mathscr{X}_{k:n}$  whether  $\underline{P}_k(\mathbb{I}_{\{o_{k:n}\}}[\mathbb{I}_{\{x_{k:n}\}} - \mathbb{I}_{\{\hat{x}_{k:n}\}}]|z_{k-1}) \leq 0$ ; see Equation (14).

If  $x_{k:n}$  is such that  $x_k = \hat{x}_k$ , this inequality is automatically satisfied. Indeed, if  $\hat{x}_k \notin \operatorname{Pos}_k(z_{k-1})$ , then we infer from Equation (24) that  $\underline{Q}_k(\{\hat{x}_k\}|z_{k-1}) = 0$  or  $\underline{S}_k(\{o_k\}|\hat{x}_k) = 0$ , and then Equation (15) tells us that  $\underline{P}_k(\mathbb{I}_{\{o_{k:n}\}}[\mathbb{I}_{\{x_{k:n}\}} - \mathbb{I}_{\{\hat{x}_{k:n}\}}]|z_{k-1}) = 0$ . If  $\hat{x}_k \in \operatorname{Pos}_k(z_{k-1})$ , we know from Equation (23) that  $\hat{x}_{k+1:n} \in \operatorname{opt}(\mathscr{X}_{k+1:n}|\hat{x}_k,o_{k+1:n})$ , which implies that  $\underline{P}_{k+1}(\mathbb{I}_{\{o_{k+1:n}\}}[\mathbb{I}_{\{x_{k+1:n}\}} - \mathbb{I}_{\{\hat{x}_{k+1:n}\}}]|\hat{x}_k) \leq 0$ . Hence  $\underline{P}_k(\mathbb{I}_{\{o_{k:n}\}}[\mathbb{I}_{\{x_{k:n}\}} - \mathbb{I}_{\{\hat{x}_{k:n}\}}]|z_{k-1}) \leq 0$ , again by Equation (15).

This means we can limit ourselves to checking the inequality for those  $x_{k:n}$  for which  $x_k \neq \hat{x}_k$ . So fix any  $x_k \neq \hat{x}_k$ , then we must check whether

$$(\forall x_{k+1:n} \in \mathscr{X}_{k+1:n})\underline{\mathcal{Q}}_{k}(\mathbb{I}_{\{x_k\}}\beta(x_{k:n}) - \mathbb{I}_{\{\hat{x}_k\}}\alpha(\hat{x}_{k:n})|z_{k-1}) \leq 0;$$

see Equation (17). By Equation (28) and Lemma 8, this is equivalent to

$$Q_{k}(\mathbb{I}_{\{x_{k}\}}\beta_{k}^{\max}(x_{k}) - \mathbb{I}_{\{\hat{x}_{k}\}}\alpha(\hat{x}_{k:n})|z_{k-1}) \leq 0,$$

which can in turn be seen to be equivalent to  $\alpha(\hat{x}_{k:n}) \geq \beta_k^{\max}(x_k)\theta_k(\hat{x}_k,x_k|z_{k-1})$ , using a course of reasoning completely analogous to the one used above for the case k=n. Since this inequality must hold for every  $x_k \neq \hat{x}_k$ , we infer from Equation (32) that we must have that  $\alpha(\hat{x}_{k:n}) \geq \alpha_k^{\text{opt}}(\hat{x}_k|z_{k-1})$ . So we must check this condition for all the candidate sequences  $\hat{x}_{k:n}$  in cand  $(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$ , which proves Equation (33).

*Proof of Theorem 4.* This proof consists of two parts. We will first prove that every sequence  $\hat{x}_{k:n}$  obtained by the optimal tree construction is an element of opt  $(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$ . Secondly we will prove that a sequence  $z_{k:n}$  that is not part of the set of sequences obtained by the optimal tree construction cannot be an element of opt  $(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$ .

Let us start by proving that every sequence  $\hat{x}_{k:n}$  obtained by the optimal tree construction is an element of opt  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ . It follows from the last step of the optimal tree construction that every  $\hat{x}_{k:n}$  of the constructed set is an element of  $\mathrm{cand}_{\hat{x}_{k:n}}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ , and therefore by Equation (26) also of  $\mathrm{cand}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ . This last step also implies that  $\alpha_n^{\max}(\hat{x}_n) \geq \alpha_k^{\mathrm{opt}}(\hat{x}_{k:n}|z_{k-1})$ , which can be seen to be equivalent with  $\alpha(\hat{x}_{k:n}) \geq \alpha_k^{\mathrm{opt}}(\hat{x}_k|z_{k-1})$ , by Equation (31) and the repeated use of Equations (35) and (20). It then follows from Equation (33) that  $\hat{x}_{k:n}$  is an element of opt  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ .

To conclude, we show that a sequence  $z_{k:n}$  that is not part of the set of sequences obtained by the optimal tree construction cannot be an element of opt  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ . If a sequence  $z_{k:n}$  is not part of the set of sequences obtained by the optimal tree construction, this either implies that it is not an element of cand  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ , or that there is some  $s \in \{k,\ldots,n\}$  for which  $\alpha_s^{\max}(z_s) < \alpha_k^{\text{opt}}(z_{k:s}|z_{k-1})$ . In the first case, it follows directly from Equation (33) that  $z_{k:n}$  cannot be an element of opt  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ . In the second case, we see that  $\alpha_s^{\max}(z_s) < \alpha_k^{\text{opt}}(z_{k:s}|z_{k-1})$  implies that  $\alpha(z_{k:n}) < \alpha_k^{\text{opt}}(z_k|z_{k-1})$ , which can be seen to be equivalent with  $\alpha(z_{k:n}) < \alpha_k^{\text{opt}}(z_k|z_{k-1})$  by the repeated use of Equations (35) and (20). It then follows from Equation (33) that  $z_{k:n}$  cannot be an element of opt  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ .

*Proof of Theorem 5.* If s = k, this can be proved by contradiction. If for all  $x_s \in \mathcal{X}_s$  both conditions would not be fulfilled, the optimal tree construction would stop and the set opt  $(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$  would be empty. This is a contradiction since every finite partially ordered set has at least one maximal element.

Now consider any  $s \in \{k+1,\ldots,n\}$ . Equation (28) implies that there is at least one sequence  $x_{s:n}^* \in \mathscr{X}_{s:n}$  for which  $\alpha_{s-1}(\hat{x}_{s-1} \oplus x_{s:n}^*) = \alpha_{s-1}^{\max}(\hat{x}_{s-1})$ . We prove that the first state  $x_s^*$  of this sequence meets both criteria of the theorem.

We know that  $\hat{x}_{k:s-1}$  is found using the optimal tree construction, which implies that  $\operatorname{cand}_{\hat{x}_{k:s-1}}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$  is a non-empty set and  $\alpha_{s-1}^{\max}(\hat{x}_{s-1}) \geq \alpha_k^{\operatorname{opt}}(\hat{x}_{k:s-1}|z_{k-1})$ . It follows from this inequality that  $\alpha_{s-1}(\hat{x}_{s-1} \oplus x_{s:n}^*) \geq \alpha_k^{\operatorname{opt}}(\hat{x}_{k:s-1}|z_{k-1})$ , which can be seen to be equivalent with  $\alpha_s(x_{s:n}^*) \geq \alpha_k^{\operatorname{opt}}(\hat{x}_{k:s-1} \oplus x_s^*|z_{k-1})$  by Equations (20) and (35). Since we know that  $\alpha_s(x_{s:n}^*) = \alpha_s^{\max}(x_s^*)$  by Equation (29), we find that  $\alpha_s^{\max}(x_s^*) \geq \alpha_k^{\operatorname{opt}}(\hat{x}_{k:s-1} \oplus x_s^*|z_{k-1})$ , meaning that  $x_s^*$  satisfies the first criterium.

To prove that the state  $x_s^*$  also satisfies the second criterium, which means that the set  $\operatorname{cand}_{\hat{x}_{k:s-1} \oplus x_s^*}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$  is non-empty, it suffices by Equation (26) to prove that  $\hat{x}_{k:s-1} \oplus x_{s:n}^*$  is an element of  $\operatorname{cand}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ .

Since  $\operatorname{cand}_{\hat{x}_{k:s-1}}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$  is non-empty, there is at least one  $z_{s:n} \in \mathscr{X}_{s:n}$  for which  $\hat{x}_{k:s-1} \oplus z_{s:n}$  is an element of  $\operatorname{cand}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ . Furthermore, we have chosen  $x_{s:n}^*$  such that  $\alpha_{s-1}(\hat{x}_{s-1} \oplus x_{s:n}^*) = \alpha_{s-1}^{\max}(\hat{x}_{s-1})$ . Lemma 9 now implies that  $\hat{x}_{k:s-1} \oplus x_{s:n}^*$  is an element of  $\operatorname{cand}(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ .

**Lemma 9.** Fix  $k \in \{1, ..., n-1\}$ ,  $s \in \{k+1, ..., n\}$ ,  $z_{k-1} \in \mathcal{X}_{k-1}$  and  $\hat{x}_{k:s-1} \in \mathcal{X}_{k:s-1}$ . Choose an arbitrary  $x_{s:n}^* \in \mathcal{X}_{s:n}$  for which  $\alpha_{s-1}(\hat{x}_{s-1} \oplus x_{s:n}^*) = \alpha_{s-1}^{\max}(\hat{x}_{s-1})$ . If there is some  $z_{s:n} \in \mathcal{X}_{s:n}$  for which  $\hat{x}_{k:s-1} \oplus z_{s:n}$  belongs to cand  $(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$ , then  $\hat{x}_{k:s-1} \oplus x_{s:n}^*$  belongs to cand  $(\mathcal{X}_{k:n}|z_{k-1},o_{k:n})$ .

*Proof.* To simplify the notations in this proof, it is convenient to use  $\hat{x}_{k-1}$  as an alternative notation for  $z_{k-1}$ . So from now on  $\hat{x}_{k-1} = z_{k-1}$ .

It follows by Lemma 10 that  $x_{s:n}^* \in \operatorname{opt}(\mathscr{X}_{s:n}|\hat{x}_{s-1},o_{s:n})$ . Together with Equation (23), this implies that  $\hat{x}_{s-1} \oplus x_{s:n}^* \in \operatorname{cand}(\mathscr{X}_{s-1:n}|\hat{x}_{s-2},o_{s-1:n})$ . If s=k+1, this concludes the proof. If  $s \in \{k+2,\ldots,n\}$ , consider all  $q \in \{k,\ldots,s-1\}$  and check af there is some q for which  $\hat{x}_q \notin \operatorname{Pos}_q(\hat{x}_{q-1})$  (see definition (24)). If such a q exists, denote the lowest  $q \in \{k,\ldots,s-2\}$  for which this is the case as  $q^*$ . By Equation (23) we know that  $\hat{x}_{q^*:s-1} \oplus x_{s:n}^*$  and  $\hat{x}_{q^*:s-1} \oplus z_{s:n}$  are both elements of cand  $(\mathscr{X}_{q^*:n}|\hat{x}_{q^*-1},o_{q^*:n})$ , since all sequences in the set  $\hat{x}_{q^*} \oplus \mathscr{X}_{q^*+1:n}$  belong to cand  $(\mathscr{X}_{q^*:n}|\hat{x}_{q^*-1},o_{q^*:n})$ .

If no  $q \in \{k, \dots, s-2\}$  exists for which  $\hat{x}_q \notin \operatorname{Pos}_q(\hat{x}_{q-1})$ , we choose  $q^* \coloneqq s-1$ . It then follows by the repeated use of Equations (22) and (23) that  $\hat{x}_{s-1} \oplus z_{s:n}$  belongs to  $\operatorname{cand}(\mathscr{X}_{s-1:n}|\hat{x}_{s-2},o_{s-1:n})$  and we already know that  $\hat{x}_{s-1} \oplus x_{s:n}^* \in \operatorname{cand}(\mathscr{X}_{s-1:n}|\hat{x}_{s-2},o_{s-1:n})$ .

We now have a  $q^* \in \{k, \ldots, s-1\}$  for which both  $\hat{x}_{q^*:s-1} \oplus x_{s:n}^*$  and  $\hat{x}_{q^*:s-1} \oplus z_{s:n}$  belong to cand  $(\mathscr{X}_{q^*:n}|\hat{x}_{q^*-1},o_{q^*:n})$  and for which it holds that  $\hat{x}_q \in \operatorname{Pos}_q(\hat{x}_{q-1})$  for all  $q \in \{k, \ldots, q^*-1\}$ . If  $q^* = k$ , this concludes the proof.

If  $q^* \in \{k+1, \ldots, s-1\}$ , notice that cand  $(\mathscr{X}_{k:n}|z_{k-1}, o_{k:n})$  is built up by repeatedly using Equations (33) and (23). We also know that  $\hat{x}_{q^*:s-1} \oplus z_{s:n} \in \operatorname{cand}(\mathscr{X}_{q^*:n}|\hat{x}_{q^*-1}, o_{q^*:n})$  and it is given that  $\hat{x}_{k:s-1} \oplus z_{s:n}$  belongs to cand  $(\mathscr{X}_{k:n}|z_{k-1}, o_{k:n})$ , which implies that

$$\alpha_q(\hat{x}_{q:s-1} \oplus z_{s:n}) \ge \alpha_q^{\text{opt}}(\hat{x}_q|\hat{x}_{q-1}) \text{ for all } q \in \{k, \dots, q^*-1\}.$$

Furthermore,  $\alpha_{s-1}(\hat{x}_{s-1} \oplus x_{s:n}^*) = \alpha_{s-1}^{\max}(\hat{x}_{s-1})$ , so  $\alpha_{s-1}(\hat{x}_{s-1} \oplus x_{s:n}^*) \ge \alpha_{s-1}(\hat{x}_{s-1} \oplus z_{s:n})$  by Equation (28). Equation (20) then tells us that

$$\alpha_t(\hat{x}_{t:s-1} \oplus x^*_{s:n}) \ge \alpha_t(\hat{x}_{t:s-1} \oplus z_{s:n}) \text{ for all } t \in \{k, \dots, s-1\},$$

so we know that

$$\alpha_q(\hat{x}_{q:s-1} \oplus x_{s:n}^*) \ge \alpha_q^{\text{opt}}(\hat{x}_q | \hat{x}_{q-1}) \text{ for all } q \in \{k, \dots, q^* - 1\}.$$

This implies (since cand  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$  is built up by repeatedly using Equations (33) and (23) and because  $\hat{x}_{q^*:s-1} \oplus x_{s:n}^*$  is an element of cand  $(\mathscr{X}_{q^*:n}|\hat{x}_{q^*-1},o_{q^*:n})$ ) that the sequence  $\hat{x}_{k:s-1} \oplus x_{s:n}^*$  belongs to cand  $(\mathscr{X}_{k:n}|z_{k-1},o_{k:n})$ , which concludes the proof.

**Lemma 10.** Consider any  $s \in \{1, ..., n\}$ ,  $\hat{x}_{s-1} \in \mathcal{X}_{s-1}$  and  $x^*_{s:n} \in \mathcal{X}_{s:n}$ . If  $\alpha_{s-1}(\hat{x}_{s-1} \oplus x^*_{s:n}) = \alpha^{\max}_{s-1}(\hat{x}_{s-1})$ , then  $x^*_{s:n} \in \text{opt}(\mathcal{X}_{s:n}|\hat{x}_{s-1}, o_{s:n})$ .

*Proof.* If  $\alpha_{s-1}(\hat{x}_{s-1} \oplus x_{s,n}^*) = \alpha_{s-1}^{\max}(\hat{x}_{s-1})$ , then we know by Equation (28) that

$$\alpha_{s-1}(\hat{x}_{s-1} \oplus x_{s:n}^*) \ge \alpha_{s-1}(\hat{x}_{s-1} \oplus z_{s:n})$$
 for all  $z_{s:n} \in \mathscr{X}_{s:n}$ ,

and therefore by Equations (19) and (7) that

$$\overline{S}_{s-1}(\{o_{s-1}\}|\hat{x}_{s-1})\overline{P}_s(\mathbb{I}_{\{x_{s:n}^*\}}\mathbb{I}_{\{o_{s:n}\}}|\hat{x}_{s-1}) \ge \overline{S}_{s-1}(\{o_{s-1}\}|\hat{x}_{s-1})\overline{P}_s(\mathbb{I}_{\{z_{s:n}\}}\mathbb{I}_{\{o_{s:n}\}}|\hat{x}_{s-1}).$$

Together with the positivity assumption (10), this implies that

$$\overline{P}_{s}(\mathbb{I}_{\{x_{s:n}^{*}\}}\mathbb{I}_{\{o_{s:n}\}}|\hat{x}_{s-1}) \ge \overline{P}_{s}(\mathbb{I}_{\{z_{s:n}\}}\mathbb{I}_{\{o_{s:n}\}}|\hat{x}_{s-1}) \text{ for all } z_{s:n} \in \mathscr{X}_{s:n}.$$
(44)

We now by C3 that

$$\underline{P}_{s}(-\mathbb{I}_{\{x_{s:n}^{*}\}}\mathbb{I}_{\{o_{s:n}\}}|\hat{x}_{s-1}) \geq \underline{P}_{s}(\mathbb{I}_{\{o_{s:n}\}}(\mathbb{I}_{\{z_{s:n}\}} - \mathbb{I}_{\{x_{s:n}^{*}\}})|\hat{x}_{s-1}) + \underline{P}_{s}(-\mathbb{I}_{\{z_{s:n}\}}\mathbb{I}_{\{o_{s:n}\}}|\hat{x}_{s-1})$$
 which by conjugacy implies that

$$\underline{P}_s(\mathbb{I}_{\{o_{s:n}\}}(\mathbb{I}_{\{z_{s:n}\}}-\mathbb{I}_{\{x_{s:n}^*\}})|\hat{x}_{s-1}) \leq \overline{P}_s(\mathbb{I}_{\{z_{s:n}\}}\mathbb{I}_{\{o_{s:n}\}}|\hat{x}_{s-1}) - \overline{P}_s(\mathbb{I}_{\{x_{s:n}^*\}}\mathbb{I}_{\{o_{s:n}\}}|x_{s-1}).$$
 Using Equation (44), we see that  $\underline{P}_s(\mathbb{I}_{\{o_{s:n}\}}(\mathbb{I}_{\{z_{s:n}\}}-\mathbb{I}_{\{x_{s:n}^*\}})|\hat{x}_s) \leq 0$  for all  $z_{s:n} \in \mathscr{X}_{s:n}$ , which concludes the proof, since  $x_{s:n}^* \in \text{opt}(\mathscr{X}_{s:n}|\hat{x}_{s-1},o_{s:n})$  by Equation (14).  $\square$