
A Novel Bayes' Theorem for Upper Probabilities

Michele Caprio¹

Yusuf Sale^{2,3}

Eyke Hüllermeier^{2,3}

Insup Lee¹

¹PRECISE Center, Department of Computer and Information Science, University of Pennsylvania, USA

²Institute of Informatics, University of Munich (LMU), Germany

³Munich Center for Machine Learning, Germany

Abstract

In their seminal 1990 paper, Wasserman and Kadane establish an upper bound for the Bayes' posterior probability of a measurable set A , when the prior lies in a class of probability measures \mathcal{P} and the likelihood is precise. They also give a sufficient condition for such upper bound to hold with equality. In this paper, we introduce a generalization of their result by additionally addressing uncertainty related to the likelihood. We give an upper bound for the posterior probability when both the prior and the likelihood belong to a set of probabilities. Furthermore, we give a sufficient condition for this upper bound to become an equality. This result is interesting on its own, and has the potential of being applied to various fields of engineering (e.g. model predictive control), machine learning, and artificial intelligence.

1 INTRODUCTION

Bayes' rule (BR) is arguably the best-known mechanism to update subjective beliefs. It prescribes the agent to elicit a prior distribution that encapsulates their initial opinion, and to come up with a likelihood that describes the data generating process. Combining prior and likelihood via BR produces the posterior distribution, which captures the agent's revised opinion in light of the collected data.

But what happens if the agent is not able to specify a single prior distribution? This may occur if they face *ambiguity* [Ellsberg, 1961, Gilboa and Marinacci, 2013]. In [Walley, 1991, Section 1.1.4] and Caprio and Gong [2023], the authors point out that missing information and bounded rationality may prevent the agent from assessing probabilities precisely in practice, even if doing so is possible in principle. This may be due to the lack of information on how likely events of interest are, lack of computational time or

ability, or because it is extremely difficult to analyze a complex body of evidence. Similarly, the agent may face difficulties in gauging the data generating process, so specifying a single likelihood may become a challenging task.

The notion of ambiguity is strictly related with that of epistemic uncertainty in machine learning (ML) and artificial intelligence (AI). Let us illustrate this clearly by borrowing concepts from [Caprio et al., 2023, Section 3.2]. Epistemic uncertainty (EU) corresponds to reducible uncertainty caused by a lack of knowledge about the best model. Notice how, in the precise case – that is, when the agent specifies a single distribution – EU is absent. In many applications, a single probability measure is only able to capture the idea of irreducible uncertainty, since it represents a case in which the agent knows exactly the true data generating process, and the prior probability that perfectly encapsulates their initial beliefs. This is a well-studied property of sets of probabilities [Hüllermeier and Waegeman, 2021, Page 458]. Due to the increasing relevance of reliable and trustworthy ML and AI applications, effective uncertainty representation and quantification have become vital research areas [Kendall and Gal, 2017, Smith and Gal, 2018, Depeweg et al., 2018, Kapoor et al., 2022, Sale et al., 2023, Wimmer et al., 2023]. Thus, theoretic underpinnings of imprecise probability theories emerge as a valuable methodology for improving the representation and quantification of uncertainties. Adopting concepts like (convex) sets of probabilities and upper and lower probabilities foster a more sophisticated and fine-tuned articulation of uncertainty.

We remark that EU should not be confused with the concept of *epistemic probability* [de Finetti, 1974, 1975, Walley, 1991]. In the subjective probability literature, epistemic probability can be captured by a single distribution. Its best definition can be found in [Walley, 1991, Sections 1.3.2 and 2.11.2]. There, the author specifies how epistemic probabilities model logical or psychological degrees of partial belief of the agent. We remark, though, how de Finetti and Walley work with finitely additive probabilities, while

in this paper we use countably additive probabilities.

The field of imprecise probabilities [Augustin, Thomas and Coolen, Frank P.A. and De Cooman, Gert and Troffaes, Matthias C.M., 2014, Walley, 1991], and in particular the classic paper by Wasserman and Kadane [1990], and successive works like Cozman [2000], Epstein and Wang [1994], Giacomini and Kitagawa [2021], Klibanoff and Hanany [2007], study the problem of an agent updating their beliefs using BR in the presence of ambiguity. Our paper belongs to this body of work. Our main result, Theorem 1, generalizes the theorem in [Wasserman and Kadane, 1990, Section 2] to the case where the agent faces ambiguity on both what prior and what likelihood to choose to model the phenomenon of interest. We find the upper posterior, that is, the “upper bound” to the set of posteriors, using only the upper prior and the upper likelihood. Thanks to the conjugacy property of upper probabilities, introduced in the next section, we derive the lower posterior, that is the “lower bound” to the set of posteriors. We also give a necessary condition for the bound to hold with equality. In addition, we hint at possible applications, in particular in the field of model predictive control, a method of process control that is used to control a process while satisfying a set of constraints Rakovi and Levine [2018].

The paper is divided as follows. Section 2.1 introduces the concepts that are needed to understand our result. In Section 2.2 we present the main theorem, and conclude our work Section 3. We prove our results in Section 4.

2 BAYES’ THEOREM FOR UPPER PROBABILITIES

2.1 PRELIMINARIES

In this section, we introduce the background notions from the IP literature [Augustin, Thomas and Coolen, Frank P.A. and De Cooman, Gert and Troffaes, Matthias C.M., 2014, Troffaes and de Cooman, 2014, Walley, 1991] that are needed to understand our main results.

Call $\Delta(\Omega, \mathcal{F})$ the space of (countably additive) probability measures on a generic measurable space (Ω, \mathcal{F}) . Pick a generic set $\mathcal{P} \subset \Delta(\Omega, \mathcal{F})$. We denote by \bar{P} the *upper probability* associated with \mathcal{P} , that is, $\bar{P}(A) = \sup_{P \in \mathcal{P}} P(A)$, for all $A \in \mathcal{F}$. Its conjugate is called *lower probability*, $\underline{P}(A) = 1 - \bar{P}(A^c) = \inf_{P \in \mathcal{P}} P(A)$, for all $A \in \mathcal{F}$. Because of the conjugacy property, in the remainder of this document we focus on upper probabilities only.

We say that upper probability \bar{P} is *concave* or *2-alternating* if $\bar{P}(A \cup B) \leq \bar{P}(A) + \bar{P}(B) - \bar{P}(A \cap B)$, for all $A, B \in \mathcal{F}$. Upper probability \bar{P} is *compatible* [Gong and Meng, 2021]

with the set

$$\begin{aligned} \text{core}(\bar{P}) &:= \{P \in \Delta(\Omega, \mathcal{F}) : P(A) \leq \bar{P}(A), \forall A \in \mathcal{F}\} \\ &= \{P \in \Delta(\Omega, \mathcal{F}) : \bar{P}(A) \geq P(A) \geq \underline{P}(A), \\ &\quad \forall A \in \mathcal{F}\} \end{aligned} \quad (1)$$

where (1) is a characterization [Cerrei-Vioglio et al., 2015, Page 3389]. The core is the set of all (countably additive) probability measures on Ω that are set-wise dominated by \bar{P} . It is convex: it is immediate to see that if P and Q are dominated by \bar{P} , then γP and $(1 - \gamma)Q$ are dominated by $\gamma \bar{P}$ and $(1 - \gamma)\bar{P}$, respectively, for all $\gamma \in [0, 1]$. In turn, $\gamma P + (1 - \gamma)Q$ is dominated by \bar{P} , thus giving the desired convex property. In addition, throughout the present work, we assume that $\text{core}(\bar{P})$ is nonempty and weak*-closed.¹ Then, as a result of [Walley, 1991, Section 3.6.1], it is also weak*-compact.

Remark 1. In [Walley, 1991, Section 3.6.1], the author shows that the finitely additive core is weak*-compact. The latter is defined as the set of all finitely additive probabilities that are set-wise dominated by \bar{P} . It is a superset of the countably additive core in (1). To see this, notice that, in general, there might well be a probability measure that is set-wise dominated by \bar{P} , but that is merely finitely additive. In fact, there might even be no countably additive probabilities that are set-wise dominated by \bar{P} . For this reason, we have to assume that the countably additive core in (1) is nonempty. If we further require that the countably additive core in (1) is weak*-closed, then, this implies that it is a weak*-closed subset of a weak*-compact space. By [Rudin, 1976, Theorem 2.35], closed subsets of compact sets are compact. In turn, we have that if the countably additive core is weak*-closed, it is also weak*-compact.

We now present a class of probabilities that (i) is well-studied and used in robust statistics Huber and Ronchetti [2009], and (ii) is the core of a concave upper probability. Other classes with similar properties are presented in [Wasserman and Kadane, 1990, Examples 3-7].

Example 1 (ε -contaminated class). Consider the space $\Delta(\Omega, \mathcal{F})$ of probability measures on a generic measurable space (Ω, \mathcal{F}) , and assume Ω is compact. Pick any $P \in \Delta(\Omega, \mathcal{F})$ and any $\varepsilon \in [0, 1]$. Define

$$\begin{aligned} \mathcal{Q}^{\text{co}} &:= \{Q \in \Delta(\Omega, \mathcal{F}) : Q(A) = (1 - \varepsilon)P(A) + \varepsilon R(A), \\ &\quad \forall A \in \mathcal{F}, R \in \Delta(\Omega, \mathcal{F})\}. \end{aligned} \quad (2)$$

Superscript “co” stands for convex and core. \mathcal{Q}^{co} is the ε -contaminated class induced by P ; it was studied in [Wasserman and Kadane, 1990, Example 3] and references

¹Recall that in the weak* topology, a net $(P_\alpha)_{\alpha \in I}$ converges to P if and only if $P_\alpha(A) \rightarrow P(A)$, for all $A \in \mathcal{F}$. See also results presented in [Walley, 1991, Appendix D3]

therein. We have that $\overline{Q}(A) = (1 - \varepsilon)P(A) + \varepsilon$, for all $A \in \mathcal{F} \setminus \{\emptyset\}$ and $\underline{Q}(A) = (1 - \varepsilon)P(A)$, for all $A \in \mathcal{F} \setminus \{\Omega\}$. In addition, $\mathcal{Q}^{\text{co}} = \text{core}(\overline{Q})$, and \overline{Q} is concave.

The ε -contaminated class is also instrumental for a future application of Theorem 1 to model predictive control. We will discuss this at length at the end of section 2.2. There, we also explain why it is important to account for the ambiguity in the likelihood model in real-world safety-critical scenarios.

2.2 A NOVEL BAYES' THEOREM FOR UPPER PROBABILITIES

Let (Θ, \mathcal{F}) be the measurable parameter space of interest and $\Delta(\Theta, \mathcal{F})$ the space of (countably additive) probability measures on (Θ, \mathcal{F}) . Let \mathcal{Y} be the set of all *bounded, non-negative, \mathcal{F} -measurable* functionals on Θ . Call \mathcal{D} the sample space endowed with sigma-algebra \mathcal{A} . That is, for any random variable Y of interest and all $\theta \in \Theta$, $Y(\theta) = y \in \mathcal{D}$. Let the agent elicit a set of probabilities $\mathcal{L}_\theta := \{P_\theta \in \Delta(\mathcal{D}, \mathcal{A}) : \theta \in \Theta\}$ on \mathcal{D} , parameterized by $\theta \in \Theta$. This captures the ambiguity faced by the agent in determining the true data generating process [Ellsberg, 1961, Gilboa and Marinacci, 2013]. We write $P_\theta \equiv P(\cdot | \theta)$ for notational convenience. Assume that each $P_\theta \in \mathcal{L}_\theta$ has density $L(\theta) = p(y | \theta)$ with respect to some sigma-finite dominating measure ν on $(\mathcal{D}, \mathcal{A})$; this represents the likelihood function for θ having observed data $y \in \mathcal{D}$.

Assumption 1. Every density L corresponding to an element P_θ of \mathcal{L}_θ belongs to \mathcal{Y} ; that is, every density is bounded and non-negative.

Assumption 1 is needed mainly for mathematical purposes; as we shall see later in this section, it can be relaxed. Let the agent specify a set \mathcal{P} of probabilities on (Θ, \mathcal{F}) . It represents their (incomplete) prior knowledge on the elements of \mathcal{F} ; its elements may be informed by the collected data, thus giving the analysis an (imprecise) empirical Bayes flavor [Casella, 1985]. Then, compute $\overline{P}(A) = \sup_{P \in \mathcal{P}} P(A)$, for all $A \in \mathcal{F}$, and consider $\mathcal{P}^{\text{co}} := \text{core}(\overline{P})$, assumed nonempty and weak*-closed. It represents the agent's initial beliefs. We assume that every $P \in \mathcal{P}^{\text{co}}$ has density p with respect to some sigma-finite dominating measure μ on (Θ, \mathcal{F}) , that is, $p = dP/d\mu$. We require the agent's beliefs to be represented by the core for two main reasons. The first, mathematical, one is to ensure that the upper probability is compatible with the belief set. The second, philosophical, one is the following. In Bayesian statistics, the agent selects a specific prior to encapsulate their initial beliefs. Berger [1984] points out how such choice is oftentimes arbitrary, and posits the *dogma of ideal precision* (DIP). It states that in any problem there is an *ideal probability model* P_T which is precise, but which may not be precisely known. To overcome this shortcom-

ing, the agent should specify a finite collection $\{P_s\}_{s \in \mathcal{S}}$ of "plausible" prior distributions, and compute the posterior for each P_s . Notice how this corresponds to selecting a finite number of elements from the core of \overline{P}_S , where $\overline{P}_S(A) = \sup_{s \in \mathcal{S}} P_s(A)$, for all $A \in \mathcal{F}$. A criticism to the DIP was brought forward by Walley. In [Walley, 1991, Section 2.10.4.(c)], he claims how given an upper probability \overline{P} , there is no cogent reason for which the agent should choose a specific P_T that is dominated by \overline{P} , or – for that matter – a collection $\{P_s\}_{s \in \mathcal{S}}$ of "plausible" probabilities. Because the core considers all (countably additive) probability measures that are dominated by \overline{P} , it is the perfect instrument to address Walley's criticism [Caprio and Gong, 2023].

Let the agent compute \overline{P}_θ , the upper probability associated with \mathcal{L}_θ , and consider $\mathcal{L}_\theta^{\text{co}} := \text{core}(\overline{P}_\theta)$, assumed nonempty and weak*-closed. It represents the set of plausible likelihoods. As Grünwald and van Ommen [2017] point out, accounting for ambiguity around the true data generating process is crucial, as Bayesian inference may suffer from inconsistency issues if carried out using a misspecified likelihood.

Let

$$\mathcal{L} := \left\{ L = \frac{dP_\theta}{d\nu}, P_\theta \in \mathcal{L}_\theta^{\text{co}} \right\} \subset \mathcal{Y} \quad (3)$$

be the set of pdf/pmf associated with the elements of $\mathcal{L}_\theta^{\text{co}}$. Let also $\overline{L}(\theta) := \sup_{L \in \mathcal{L}} L(\theta)$ and $\underline{L}(\theta) := \inf_{L \in \mathcal{L}} L(\theta)$, for all $\theta \in \Theta$. Call

$$\mathcal{P}_y^{\text{co}} := \left\{ P_y \in \Delta(\Theta, \mathcal{F}) : \begin{aligned} \frac{dP_y}{d\mu} = p(\theta | y) &= \frac{L(\theta)p(\theta)}{\int_\Theta L(\theta)p(\theta)d\theta}, \\ p &= \frac{dP}{d\mu}, P \in \mathcal{P}^{\text{co}}, L = \frac{dP_\theta}{d\nu}, P_\theta \in \mathcal{L}_\theta^{\text{co}} \end{aligned} \right\}$$

the class of posterior probabilities when the prior is in \mathcal{P}^{co} and the likelihood is in $\mathcal{L}_\theta^{\text{co}}$, and let $\overline{P}_y(A) = \sup_{P_y \in \mathcal{P}_y^{\text{co}}} P_y(A)$, for all $A \in \mathcal{F}$. Then, the following is a generalization of Bayes' theorem in [Wasserman and Kadane, 1990, Section 2], and is an extension of [Caprio et al., 2023, Theorem 7]. We prove it in Section 4.

Theorem 1 (Bayes' theorem for upper probabilities). *Suppose $\mathcal{P}^{\text{co}}, \mathcal{L}_\theta^{\text{co}}$ are nonempty and weak*-closed. Then for all $A \in \mathcal{F}$,*

$$\overline{P}_y(A) \leq \frac{\sup_{P \in \mathcal{P}^{\text{co}}} \int_\Theta \overline{L}(\theta) \mathbb{1}_A(\theta) P(d\theta)}{c} \quad (4)$$

$$\leq \frac{\int_0^\infty \overline{P}(\{\theta \in \Theta : \overline{L}(\theta) \mathbb{1}_A(\theta) > t\}) dt}{c'}, \quad (5)$$

provided that the ratios are well defined. Here $\mathbb{1}_A$ denotes the indicator function for

$A \in \mathcal{F}$, $\mathbf{c} := \sup_{P \in \mathcal{P}^{co}} \int_{\Theta} \bar{L}(\theta) \mathbb{1}_A(\theta) P(d\theta) + \inf_{P \in \mathcal{P}^{co}} \int_{\Theta} \underline{L}(\theta) \mathbb{1}_{A^c}(\theta) P(d\theta)$, and

$$\mathbf{c}' := \underbrace{\int_0^\infty \bar{P}(\{\theta \in \Theta : \bar{L}(\theta) \mathbb{1}_A(\theta) > t\}) dt}_{\text{upper Choquet integral of } \bar{L} \mathbb{1}_A} + \underbrace{\int_0^\infty \underline{P}(\{\theta \in \Theta : \underline{L}(\theta) \mathbb{1}_{A^c}(\theta) > t\}) dt}_{\text{lower Choquet integral of } \underline{L} \mathbb{1}_{A^c}}.$$

In addition, if \bar{P} is concave, then the inequalities in (4) and (5) are equalities, for all $A \in \mathcal{F}$.

This result is particularly appealing. Under Assumption 1, if the prior upper probability (PUP) is concave and the prior and likelihood sets \mathcal{P}^{co} , \mathcal{L}_θ^{co} are nonempty and weak*-closed, then the agent can perform a (generalized) Bayesian update of the PUP by carrying out only one operation. This is the case even when the agent faces ambiguity around the true data generating process so that a set of likelihoods is needed. The posterior lower probability is obtained immediately via the conjugacy property $\underline{P}_y(A) = 1 - \bar{P}_y(A^c)$.

Corollary 1.1. *Retain the assumptions in Theorem 1. If \mathcal{L}_θ^{co} is a singleton, we retrieve Bayes' theorem in [Wasserman and Kadane, 1990, Section 2].*

Corollary 1.1 tell us that when there is no ambiguity around the likelihood, Theorem 1 recovers Wasserman and Kadane's classical Bayes' theorem. Given the straightforward nature of the proof, we omit it here. We also have the following lemma, that is proved in Section 4.

Lemma 2 (Preserved concavity). *Suppose \mathcal{P}^{co} , \mathcal{L}_θ^{co} are nonempty and weak*-closed. Then, if \bar{P} is concave, we have that \bar{P}_y is concave as well.*

This lemma is important because it tells us that the generalized Bayesian update of Theorem 1 preserves concavity, and so it can be applied to successive iterations. If at time t the PUP is concave, then the PUP at time $t + 1$ – that is, the posterior upper probability at time t – will be concave too. Necessary and sufficient conditions for a generic upper probability to be concave are given in [Marinacci and Montrucchio, 2004, Section 5].

In the future, we plan to forego Assumption 1 and use the techniques developed in Troffaes and de Cooman [2014] to generalize our results to the case in which the elements of \mathcal{L} are unbounded and not necessarily non-negative. We also intend to extend our results to the case in which the elements of \mathcal{Y} are \mathbb{R}^d -valued, for some $d \in \mathbb{N}$. We suspect this is a less demanding endeavor since we do not use specific properties of \mathbb{R} in our proofs.

As mentioned earlier, a natural application of our results is model predictive control (MPC). MPC is a method that

is used to control a process while satisfying a set of constraints [Rakovi and Levine, 2018]. Typically, when the process is stable (or at least stable for the past k time steps, for some $k \geq 0$) and if the scholar decides to take the Bayesian approach, they proceed as follows. They specify a Normal likelihood (the distribution of the control inputs) centered at the objective function of the process, and a Normal prior on the parameter of such function. In this framework, if ambiguity enters the picture, then our results become relevant.

The importance of addressing prior ambiguity was discussed at length in section 2.2. The reasons why accounting for likelihood ambiguity is important are as follows. First, as pointed out earlier, we may run into inconsistency issues if we perform Bayesian analysis using a misspecified likelihood [Grünwald and van Ommen, 2017]. Second, a (precise) Normal likelihood is a good choice only if stability of the process is ensured. MPC procedures are used in the process industries in chemical plants, oil refineries, power system balancing models, and power electronics. These are all safety-critical applications where accounting for possible sudden unexpected instabilities is of paramount importance.

The scholar may specify a class of ϵ -contaminated truncated Normal priors and a class of η -contaminated truncated Normal likelihoods, and use Theorem 1 to compute the upper posterior. Notice that the Normals need to be truncated in light of Assumption 1. This requirement is not too stringent, and – as pointed out earlier in this section – our future work will allow us do without it.

3 CONCLUSION

In this paper, we present a new Bayes' theorem for upper probabilities that extends the one in [Wasserman and Kadane, 1990, Section 2], and [Caprio et al., 2023, Theorem 7]. In the future, we plan to generalize Theorem 1 by letting go of Assumption 1, and to apply it to an MPC problem and to other fields of engineering, and ML and AI. For example, we intend to use it to overcome the computational bottleneck of step 2 of the algorithm that computes the posterior set in an imprecise Bayesian neural network procedure [Caprio et al., 2023]. There, an element-wise application of Bayes' rule for all the extreme elements of the prior and likelihood sets is performed. As we can see, this is a combinatorial task that can potentially be greatly simplified in light of Theorem 1, conveying a computationally cheaper algorithm.

4 PROOFS

Proof of Theorem 1. Assume that \mathcal{P}^{co} , \mathcal{L}_θ^{co} are nonempty and weak*-closed. Pick any $A \in \mathcal{F}$. Recall that we can

rewrite the usual Bayes' updating rule as

$$\begin{aligned} P_y(A) &= \frac{\int_{\Theta} L(\theta) \mathbb{1}_A(\theta) P(d\theta)}{\int_{\Theta} L(\theta) \mathbb{1}_A(\theta) P(d\theta) + \int_{\Theta} L(\theta) \mathbb{1}_{A^c}(\theta) P(d\theta)} \\ &= \frac{1}{1 + \frac{\int_{\Theta} L(\theta) \mathbb{1}_{A^c}(\theta) P(d\theta)}{\int_{\Theta} L(\theta) \mathbb{1}_A(\theta) P(d\theta)}}, \end{aligned}$$

which is maximized when

$$\frac{\int_{\Theta} L(\theta) \mathbb{1}_{A^c}(\theta) P(d\theta)}{\int_{\Theta} L(\theta) \mathbb{1}_A(\theta) P(d\theta)}$$

is minimized. But

$$\frac{\int_{\Theta} L(\theta) \mathbb{1}_{A^c}(\theta) P(d\theta)}{\int_{\Theta} L(\theta) \mathbb{1}_A(\theta) P(d\theta)} \geq \frac{\inf_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \underline{L}(\theta) \mathbb{1}_{A^c}(\theta) P(d\theta)}{\sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \overline{L}(\theta) \mathbb{1}_A(\theta) P(d\theta)},$$

which proves the inequality in (4). The inequality in (5) is true because

$$\begin{aligned} &\inf_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \underline{L}(\theta) \mathbb{1}_{A^c}(\theta) P(d\theta) \\ &\geq \int_0^{\infty} \underline{P}(\{\theta \in \Theta : \underline{L}(\theta) \mathbb{1}_{A^c}(\theta) > t\}) dt \end{aligned}$$

and

$$\begin{aligned} &\sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \overline{L}(\theta) \mathbb{1}_A(\theta) P(d\theta) \\ &\leq \int_0^{\infty} \overline{P}(\{\theta \in \Theta : \overline{L}(\theta) \mathbb{1}_A(\theta) > t\}) dt. \end{aligned}$$

Assume now that \overline{P} is concave. By [Wasserman and Kadane, 1990, Lemma 1], we have that there exists $\check{P} \in \mathcal{P}^{\text{co}}$ such that

$$\sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} L(\theta) \mathbb{1}_A(\theta) P(d\theta) = \int_{\Theta} L(\theta) \mathbb{1}_A(\theta) \check{P}(d\theta), \quad (6)$$

for all $L \in \mathcal{L}$. In addition, by [Wasserman and Kadane, 1990, Lemma 4], we have that for all $Y \in \mathcal{Y}$ and all $\epsilon > 0$, there exists a non-negative, upper semi-continuous function $h \leq Y$ such that

$$\begin{aligned} \left[\sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} Y(\theta) P(d\theta) \right] - \epsilon &< \sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} h(\theta) P(d\theta) \\ &\leq \sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} Y(\theta) P(d\theta). \end{aligned} \quad (7)$$

Let now $Y = \overline{L} \mathbb{1}_A$. Notice that since $\mathcal{L}_{\theta}^{\text{co}}$ is weak*-compact (as a result of [Walley, 1991, Section 3.6.1]), by (3) so is \mathcal{L} . This implies that $\underline{L}, \overline{L} \in \mathcal{L}$, since a compact set always contains its boundary, so $\underline{L}, \overline{L} \in \mathcal{Y}$ as well, and in turn $\underline{L} \mathbb{1}_{A^c}, \overline{L} \mathbb{1}_A \in \mathcal{Y}$. Fix then any $L \in \mathcal{L}$ and put $h = L \mathbb{1}_A$. It is immediate to see that h is non-negative and upper semi-continuous. Then, by (7), we have that for all $\epsilon > 0$

$$\begin{aligned} \left[\sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \overline{L}(\theta) \mathbb{1}_A(\theta) P(d\theta) \right] - \epsilon &< \\ \sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} L(\theta) \mathbb{1}_A(\theta) P(d\theta) &\leq \sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \overline{L}(\theta) \mathbb{1}_A(\theta) P(d\theta). \end{aligned} \quad (8)$$

Combining (6) and (8), we obtain

$$\begin{aligned} &\left[\sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \overline{L}(\theta) \mathbb{1}_A(\theta) P(d\theta) \right] - \epsilon \\ &< \int_{\Theta} L(\theta) \mathbb{1}_A(\theta) \check{P}(d\theta) \leq \sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \overline{L}(\theta) \mathbb{1}_A(\theta) P(d\theta), \end{aligned} \quad (9)$$

for all $L \in \mathcal{L}$.

Pick now any $\epsilon > 0$ and put

$$\begin{aligned} k &:= \sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \overline{L}(\theta) \mathbb{1}_A(\theta) P(d\theta) \\ &+ \inf_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \underline{L}(\theta) \mathbb{1}_{A^c}(\theta) P(d\theta) > 0. \end{aligned}$$

Choose any $L \in \mathcal{L}$ and $\delta \in (0, \epsilon k)$. By (9) we have that

$$\left[\sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \overline{L}(\theta) \mathbb{1}_A(\theta) P(d\theta) \right] - \delta < \int_{\Theta} L(\theta) \mathbb{1}_A(\theta) \check{P}(d\theta) \quad (10)$$

and that

$$\left[\inf_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \underline{L}(\theta) \mathbb{1}_{A^c}(\theta) P(d\theta) \right] + \delta > \int_{\Theta} L(\theta) \mathbb{1}_{A^c}(\theta) \check{P}(d\theta). \quad (11)$$

The inequality in (10) comes from the fact that the first inequality in (9) holds for all $\epsilon > 0$, and – given how k is defined – we have that $\delta > 0$. The inequality in (11) is obtained by re-deriving (6), (7), (8), and (9) for the infimum of set \mathcal{P}^{co} rather than the supremum. In that case, we simply substitute sup with inf, \overline{L} with \underline{L} , $\mathbb{1}_A$ with $\mathbb{1}_{A^c}$, “ $-\epsilon$ ” with “ $+\epsilon$ ”, and we reverse the inequalities.

Recall that $\mathbf{c} := \sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \overline{L}(\theta) \mathbb{1}_A(\theta) P(d\theta) + \inf_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \underline{L}(\theta) \mathbb{1}_{A^c}(\theta) P(d\theta)$, and define

$$\mathbf{d} := \int_{\Theta} L(\theta) \mathbb{1}_A(\theta) \check{P}(d\theta) + \int_{\Theta} L(\theta) \mathbb{1}_{A^c}(\theta) \check{P}(d\theta).$$

Then we have,

$$\begin{aligned} \check{P}_y(A) &= \frac{\int_{\Theta} L(\theta) \mathbb{1}_A(\theta) \check{P}(d\theta)}{\mathbf{d}} \\ &\geq \frac{[\sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \overline{L}(\theta) \mathbb{1}_A(\theta) P(d\theta)] - \delta}{\mathbf{c} + \delta - \delta} \\ &= \frac{\sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \overline{L}(\theta) \mathbb{1}_A(\theta) P(d\theta)}{\mathbf{c}} - \frac{\delta}{k} \\ &> \frac{\sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \overline{L}(\theta) \mathbb{1}_A(\theta) P(d\theta)}{\mathbf{c}} - \epsilon. \end{aligned}$$

Since this holds for all $\epsilon > 0$, we have that

$$\sup_{P_y \in \mathcal{P}_y^{\text{co}}} P_y(A) = \frac{\sup_{P \in \mathcal{P}^{\text{co}}} \int_{\Theta} \overline{L}(\theta) \mathbb{1}_A(\theta) P(d\theta)}{\mathbf{c}},$$

concluding the proof of inequality (4) being an equality when \overline{P} is concave. Inequality (5) being an equality when

\overline{P} is concave comes immediately from [Wasserman and Kadane, 1990, Lemma 4], and the fact that $\overline{L}\mathbb{1}_A, \underline{L}\mathbb{1}_{A^c} \in \mathcal{Y}$, as pointed out above. \square

Proof of Lemma 2. In their works Walley [1981], Wasserman and Kadane [1990], the authors show that concave upper probabilities are closed with respect to the generalized Bayes' rule. In particular, this means that, if we let $\mathbf{b} := \sup_{P \in \mathcal{P}^\infty} \int_{\Theta} L(\theta) \mathbb{1}_A(\theta) P(d\theta) + \inf_{P \in \mathcal{P}^\infty} \int_{\Theta} L(\theta) \mathbb{1}_{A^c}(\theta) P(d\theta)$, for any fixed $A \in \mathcal{F}$, if \overline{P} is concave, then for all $L \in \mathcal{L}$

$$\overline{P}_y(A) = \frac{\sup_{P \in \mathcal{P}^\infty} \int_{\Theta} L(\theta) \mathbb{1}_A(\theta) P(d\theta)}{\mathbf{b}} \quad (12)$$

is concave. But since $\mathcal{L}_\theta^{\text{co}}$ is weak*-compact (as a consequence of [Walley, 1991, Section 3.6.1]), by (3) so is \mathcal{L} . This implies that $\underline{L}, \overline{L} \in \mathcal{L}$, since a compact set always contains its boundary. Call then $L' = \overline{L}\mathbb{1}_A + \underline{L}\mathbb{1}_{A^c}$. It is immediate to see that $L' \in \mathcal{L}$. Then, by (12) we have that if we call $\mathbf{b}' := \sup_{P \in \mathcal{P}^\infty} \int_{\Theta} L'(\theta) \mathbb{1}_A(\theta) P(d\theta) + \inf_{P \in \mathcal{P}^\infty} \int_{\Theta} L'(\theta) \mathbb{1}_{A^c}(\theta) P(d\theta)$, it follows that

$$\begin{aligned} \overline{P}_y(A) &= \frac{\sup_{P \in \mathcal{P}^\infty} \int_{\Theta} L'(\theta) \mathbb{1}_A(\theta) P(d\theta)}{\mathbf{b}'} \\ &= \frac{\sup_{P \in \mathcal{P}^\infty} \int_{\Theta} \overline{L}(\theta) \mathbb{1}_A(\theta) P(d\theta)}{\mathbf{c}} \end{aligned}$$

is concave, concluding the proof. \square

Author Contributions

Michele Caprio and Yusuf Sale contributed equally to this paper.

Acknowledgements

Michele Caprio would like to acknowledge partial funding by the Army Research Office (ARO MURI W911NF2010080). Yusuf Sale is supported by the DAAD programme Konrad Zuse Schools of Excellence in Artificial Intelligence, sponsored by the Federal Ministry of Education and Research.

References

- Augustin, Thomas and Coolen, Frank P.A. and De Cooman, Gert and Troffaes, Matthias C.M., editor. *Introduction to imprecise probabilities*. Wiley Series in Probability and Statistics. John Wiley and Sons, 2014.
- James O. Berger. The robust Bayesian viewpoint. In Joseph B. Kadane, editor, *Robustness of Bayesian Analyses*. Amsterdam : North-Holland, 1984.

Michele Caprio and Ruobin Gong. Dynamic imprecise probability kinematics. *Proceedings of Machine Learning Research*, 215:72–83, 2023.

Michele Caprio, Souradeep Dutta, Radoslav Ivanov, Kuk Jang, Vivian Lin, Oleg Sokolsky, and Insup Lee. Imprecise Bayesian neural networks. *Available at arXiv:2302.09656*, 2023.

George Casella. An introduction to empirical Bayes data analysis. *The American Statistician*, 39(2):83–87, 1985.

Simone Cerreia-Vioglio, Fabio Maccheroni, and Massimo Marinacci. Ergodic theorems for lower probabilities. *Proceedings of the American Mathematical Society*, 144: 3381–3396, 2015.

Fabio Gagliardi Cozman. Credal networks. *Artificial Intelligence*, 120:199–233, 2000.

Bruno de Finetti. *Theory of Probability*, volume 1. New York : Wiley, 1974.

Bruno de Finetti. *Theory of Probability*, volume 2. New York : Wiley, 1975.

Stefan Depeweg, Jose-Miguel Hernandez-Lobato, Finale Doshi-Velez, and Steffen Udfluft. Decomposition of uncertainty in bayesian deep learning for efficient and risk-sensitive learning. In *International Conference on Machine Learning*, pages 1184–1193. PMLR, 2018.

Daniel Ellsberg. Risk, ambiguity, and the Savage axioms. *The Quarterly Journal of Economics*, 75(4):643–669, 1961.

Larry G Epstein and Tan Wang. Intertemporal Asset Pricing Under Knightian Uncertainty. *Econometrica*, 62(2): 283–322, 1994.

Raffaella Giacomini and Toru Kitagawa. Robust Bayesian inference for set-identified models. *Econometrica*, 89(4): 1519–1556, 2021.

Itzhak Gilboa and Massimo Marinacci. Ambiguity and the Bayesian paradigm. In Daron Acemoglu, Manuel Arellano, and Eddie Dekel, editors, *Advances in Economics and Econometrics, Tenth World Congress*, volume 1. Cambridge : Cambridge University Press, 2013.

Ruobin Gong and Xiao-Li Meng. Judicious judgment meets unsettling updating: dilation, sure loss, and Simpson's paradox. *Statistical Science*, 36(2):169–190, 2021.

Peter Grünwald and Thijs van Ommen. Inconsistency of Bayesian Inference for Misspecified Linear Models, and a Proposal for Repairing It. *Bayesian Analysis*, 12(4): 1069 – 1103, 2017.

- Peter J. Huber and Elvezio M. Ronchetti. *Robust statistics*. Wiley Series in Probability and Statistics. Hoboken, New Jersey : Wiley, 2nd edition, 2009.
- Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning: an introduction to concepts and methods. *Machine Learning*, 3 (110):457–506, 2021.
- Sanyam Kapoor, Wesley J Maddox, Pavel Izmailov, and Andrew Gordon Wilson. On uncertainty, tempering, and data augmentation in bayesian classification. *arXiv preprint arXiv:2203.16481*, 2022.
- Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? *Advances in neural information processing systems*, 30, 2017.
- Peter Klibanoff and Eran Hanany. Updating preferences with multiple priors. *Theoretical Economics*, 2(3):261–298, 2007.
- Massimo Marinacci and Luigi Montrucchio. Introduction to the mathematics of ambiguity. In Itzhak Gilboa, editor, *Uncertainty in economic theory: a collection of essays in honor of David Schmeidler's 65th birthday*. London : Routledge, 2004.
- Saa V. Rakovi and William S. Levine. *Handbook of Model Predictive Control*. Birkhäuser Basel, 1st edition, 2018.
- Walter Rudin. *Principles of Mathematical Analysis*. New york : McGraw-Hill, 3rd edition, 1976.
- Yusuf Sale, Michele Caprio, and Eyke Höllermeier. Is the volume of a credal set a good measure for epistemic uncertainty? In *Uncertainty in Artificial Intelligence*, pages 1795–1804. PMLR, 2023.
- Lewis Smith and Yarin Gal. Understanding measures of uncertainty for adversarial example detection. *arXiv preprint arXiv:1803.08533*, 2018.
- Matthias C.M. Troffaes and Gert de Cooman. *Lower Previsions*. Chichester, United Kingdom : John Wiley and Sons, 2014.
- Peter Walley. Coherent lower (and upper) probabilities. Technical report, University of Warwick, Coventry, 1981.
- Peter Walley. *Statistical Reasoning with Imprecise Probabilities*, volume 42 of *Monographs on Statistics and Applied Probability*. London : Chapman and Hall, 1991.
- Larry A. Wasserman and Joseph B. Kadane. Bayes' theorem for Choquet capacities. *The Annals of Statistics*, 18 (3):1328–1339, 1990.
- Lisa Wimmer, Yusuf Sale, Paul Hofman, Bernd Bischl, and Eyke Hüllermeier. Quantifying aleatoric and epistemic uncertainty in machine learning: Are conditional entropy and mutual information appropriate measures? In *Uncertainty in Artificial Intelligence*, pages 2282–2292. PMLR, 2023.