

# The principle of maximum entropy and a problem in probability kinematics

Stefan Lukits

Received: 22 December 2012 / Accepted: 19 August 2013 / Published online: 5 September 2013  
© Springer Science+Business Media Dordrecht 2013

**Abstract** Sometimes we receive evidence in a form that standard conditioning (or Jeffrey conditioning) cannot accommodate. The principle of maximum entropy (MAX-ENT) provides a unique solution for the posterior probability distribution based on the intuition that the information gain consistent with assumptions and evidence should be minimal. Opponents of objective methods to determine these probabilities prominently cite van Fraassen's Judy Benjamin case to undermine the generality of MAXENT. This article shows that an intuitive approach to Judy Benjamin's case supports MAXENT. This is surprising because based on independence assumptions the anticipated result is that it would support the opponents. It also demonstrates that opponents improperly apply independence assumptions to the problem.

**Keywords** Judy Benjamin · Principle of maximum entropy · Coarsening at random · Full employment theorem · Probability kinematics · Epistemic entrenchment

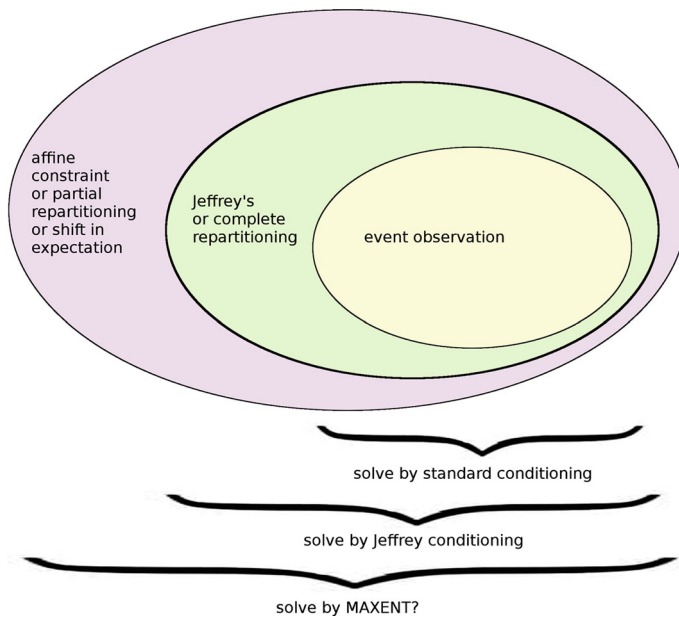
## 1 Introduction

Probability kinematics is the field of inquiry asking how we should update a probability distribution in the light of evidence. If the evidence comes as an event, it is relatively uncontroversial to use conditional probabilities (call this standard conditioning). Sometimes, however, the evidence may not relate the certainty of an event but a reassessment of its uncertainty or its probabilistic relation to other events (see Jeffrey 1965, 153ff), expressible in a shift in expectation (see Hobson 1971). Jeffrey conditionalization can deal with some of these cases, but not with all of them (see Fig. 1). Bas van Fraassen has provided an example for a case in which we cannot apply Jeffrey

---

S. Lukits (✉)

University of British Columbia, 1866 Main Mall Buchanan E370, Vancouver, BC V6T 1Z1, Canada  
e-mail: saiserit@streetgreek.com



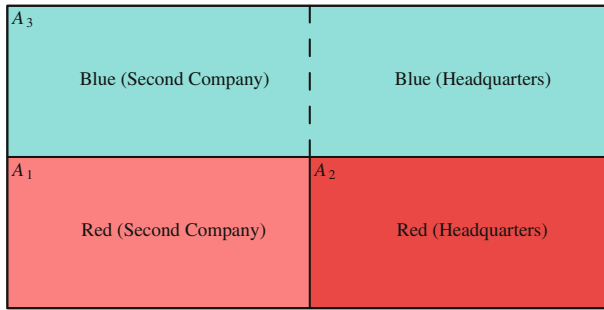
**Fig. 1** Information that leads to unique solutions for probability updating using MAXENT must come in the form of an affine constraint (a constraint in the form of an informationally closed and convex space of probability distributions consistent with the information). All information that can be processed by standard conditioning or Jeffrey conditioning comes in the form of an affine constraint. The solutions of MAXENT are consistent with the solutions of standard conditioning and Jeffrey conditioning where the latter two are applicable

conditionalization. The example is from the 1980 comedy film *Private Benjamin* (see [Van Fraassen 1981](#)), in which Goldie Hawn portrays a Jewish-American woman (Judy Benjamin) who joins the U.S. Army.

In van Fraassen's story based on the movie, Judy Benjamin is on an assignment and lands in a place where she is not sure of her location. She is on team Blue. Because of the map, initially the probability of being in Blue territory equals the probability of being in Red territory, and the probability of being in the Red Second Company area equals the probability of being in the Red Headquarters area. Her commanders then inform Judy by radio that in case she is in Red territory, her chance of being in the Red Headquarters area is three times the chance of being in the Red Second Company area. The question is what Judy's appropriate response is to this new evidence.

We cannot apply standard conditioning, because there is no immediately obvious event space in which we can condition on an event of which we are certain. [Grove and Halpern \(1997\)](#) have offered a proposal for constructing such event spaces and then conditioning on the event that Judy Benjamin receives the information that she receives from her commanders. They admit, however, that the construction of such spaces (sometimes called retrospective conditioning) is an exercise in filling in missing details and supplying information not contained in the original problem.

If we assume that the attempt fails to define an event on which Judy Benjamin could condition her probabilities, we are left with two possibilities. Her new information



**Fig. 2** Judy Benjamin's map. *Blue territory* ( $A_3$ ) is friendly and does not need to be divided into a Headquarters and a Second Company area. (Color figure online)

(it is three times as likely to land in  $A_2$  than to land in  $A_1$ , see Fig. 2 and the details of the problem in the next section) may mean that we have a redistribution of a complete partition of the probabilities. This is called Jeffrey conditioning and calls for Jeffrey's rule. Jeffrey's rule is contested in some circles, but we will for this project accept its validity in probability kinematics. We will see in what follows that some make the case that Jeffrey conditioning is the correct way to solve the Judy Benjamin problem. For reasons provided in the body of the paper their case is implausible.

The third possibility to solve this problem (after standard conditioning and Jeffrey conditioning) is to consult a highly contested updating procedure: the principle of maximum entropy (MAXENT for short). MAXENT can be applied to any situation in which we have a completely quantified probability distribution and an affine constraint (we will explain the nature of affine constraints in more detail later). If our new evidence is the observation of an event (or simply certainty about an event that we did not have previously), then the event provides an affine constraint and can be used for updating by means of standard conditioning. If our new evidence is a redistribution of probabilities where we can apply Jeffrey's rule, then the redistribution provides an affine constraint and can be used for updating by means of Jeffrey's rule. These two possibilities, however, do not exhaust affine constraints. The Judy Benjamin problem illustrates the third possibility where the affine constraint only redistributes some groups of probabilities and leaves open the question how this will affect the probabilities not included in this redistribution.

Advocates of MAXENT claim that in this case the probabilities should be adjusted so that they are minimally affected (we make this precise by using information theory) while at the same time conforming to the constraint. Opponents of this view grant that MAXENT is an important tool of probability kinematics. Noting that some results of MAXENT are difficult to accept (such as in the Judy Benjamin case), however, they urge us to embrace a more pluralistic, situation-specific methodology.

Joseph Halpern, for example, writes in *Reasoning About Uncertainty* that “there is no escaping the need to understand the details of the application” (Halpern 2003, p. 423) and concludes that MAXENT is a valuable tool, but should be used with care (see Grove and Halpern 1997, p. 110), explicitly basing his remark on the counterintuitive behaviour of the Judy Benjamin problem. Diaconis and Zabell state “that any claims

to the effect that maximum-entropy revision is the only correct route to probability revision should be viewed with considerable caution” (Diaconis and Zabell 1982, p. 829). “Great caution” (1994, p. 456) is also what Colin Howson and Allan Franklin advise about the more basic claim that the updated probabilities provided by MAXENT are as like the original probabilities as it is possible to be given the constraints imposed by the data.

In the same vein, Igor Douven and Jan-Willem Romeijn agree with Richard Bradley that “even Bayes’ rule ‘should not be thought of as a universal and mechanical rule of updating, but as a technique to be applied in the right circumstances, as a tool in what Jeffrey terms *the art of judgment*.’ In the same way, determining and adapting the weights ... may be an art, or a skill, rather than a matter of calculation or derivation from more fundamental epistemic principles” (Douven and Romeijn 2009, p. 16) (for the Bradley quote see Bradley 2005, p. 362).

What is lacking in the literature is a response by MAXENT advocates to the counterintuitive behaviour of the cases repeatedly quoted by their adversaries. This is especially surprising as we are not dealing with an array of counter-examples but only a handful, the Judy Benjamin problem being prime among them. In Halpern’s textbook, for example, the reasoning is as follows: MAXENT is a promising candidate which delivers unique updated probability distributions; but, unfortunately, there is counterintuitive behaviour in one specific case, the Judy Benjamin case (see Halpern 2003, pp. 110, 119); therefore, we must abide by the eclectic principle of considering not only MAXENT, but also lower and upper probabilities, Dempster–Shafer belief functions, possibility measures, ranking functions, relative likelihoods, and so forth. The human inquirer is the final arbiter between these conditionalization methods.

At the heart of our investigation are two incompatible but independently plausible intuitions regarding Judy’s choice of updated probabilities for her location. We will undermine the notion that MAXENT’s solution for the Judy Benjamin problem is counterintuitive. The intuition that MAXENT’s solution for the Judy Benjamin problem violates (call it T1) is based on fallacious independence and uniformity assumptions. There is another powerful intuition (call it T2) that conflicts with T1 and obeys MAXENT. Therefore, Halpern does not give us sufficient grounds for the eclecticism advocated throughout his book. We will show that another intuitive approach, the powerset approach, lends significant support to the solution provided by MAXENT for the Judy Benjamin problem, especially in comparison to intuition T1, many of whose independence and uniformity assumptions it shares.

We have no proof that MAXENT is the only rationally defensible objective method to update probabilities given an affine constraint. The literature outlines many of the ‘nice properties’ of MAXENT. It seamlessly generalizes standard conditioning and Jeffrey’s rule where they are applicable (see Caticha and Giffin 2006). It underlies the entropy concentration phenomenon described in Jaynes’ standard work *Probability Theory: the Logic of Science*, which contains other arguments in favour of MAXENT (some of which you may recognize by family resemblance in the rest of this paper). Entropy concentration refers to the unique property of the MAXENT solution to have other distributions which obey the affine constraint cluster around it. Shore and Johnson have shown that under certain rationality assumptions MAXENT provides unique solutions to problems of probability updating (see Shore and Johnson 1980). When

used to make predictions, posterior probabilities provided by MAXENT result in mini-max optimal decisions (see [Topsøe 1979](#); [Walley 1991](#); [Grünwald 2000](#)), and they are optimal under a logarithmic scoring rule.

Despite all of these nice properties, we want the reader to follow us in a more simple line of argument. When new evidence is provided to us, it is rational to adjust our beliefs minimally in light of it. We do not want to draw more information from the new evidence than necessary. There are numerous problems that need addressing. What do we mean by rationality? What are the semantics of the word ‘minimal’? What are the formal properties of such posterior probabilities? Are they unique? Are they compatible with other intuitive methods of updating? Are there counter-intuitive examples that would encourage us to give up on this line of thought rather than live with its consequences? Given some decent answers to these questions, however, we feel that MAXENT cuts a good figure as a first pass to provide objective solutions to these types of problems, and the burden on opponents who usually deny that there are such objective solutions exist grows heavy.

The distinctive contribution of this paper is to show why the reasoning of the opponents of MAXENT in the Judy Benjamin case is flawed. They make independence assumptions that on closer inspection do not hold up. We provide a number of scenarios consistent with the information in the problem which violate these independence assumptions. That does not mean that the information given in the problem suggests these scenarios, it only means that we are not entitled to make those independence assumptions. That, in turn, does not privilege the MAXENT solution, although MAXENT does not lean on independence assumptions that other solutions illegitimately make. MAXENT, however, confronts us with a much stronger claim than merely providing a passable or useful solution to the Judy Benjamin problem: it claims much greater generality and, to use a term abjured by many formal epistemologists, it claims objectivity. This claim must be motivated elsewhere, and the nature of its normativity is a matter of debate (for a pragmatic approach see [Caticha 2012](#)). We are only showing that opponents cannot declare an easy victory by pulling out old Judy Benjamin.

There is a long-standing disagreement between (mostly) philosophers on the one hand and (mostly) physicists on the other hand. The philosophers claim that updating probabilities is irreducibly accompanied by thoughtful deliberation with the choice between different updating procedures depending on individual problems. The physicists claim that problems are ill-posed if they do not contain the information necessary to let a non-arbitrary, objective procedure (such as MAXENT) arrive at a unique updated probability distribution. In the literature, Judy Benjamin serves as an example widely taken to count in favour of the philosophers. It is taken to support what I shall call the full employment theorem of probability kinematics.

The full employment theorem of probability kinematics claims that MAXENT is only one of many different strategies to update probabilities. In order to decide which strategy is the most appropriate for your problem you need a resident formal epistemologist to do the thinking and weigh the intuitions for you. For a fee, of course. Thus formal epistemologists will always be fully employed. (E.T. Jaynes makes similar observations when he derisively talks about the statistician–client relationship as one between a doctor and his patient, see Jaynes and Bretthorst (1998, pp. 492 and 506).) There is an analogous full employment theory in computer science about writing computer

programs which has been formally proven to be true. Our contention is that the full employment theorem is neither proven nor even plausible in probability kinematics.

## 2 Two intuitions

There are two pieces of information relevant to Judy Benjamin when she decides on her updated probability assignment. We will call them (MAP) and (HDQ). As in Fig. 2,  $A_1$  is the Red Second Company area,  $A_2$  is the Red Headquarters area,  $A_3$  is Blue territory. Judy presumably wants to be in Blue territory, but if she is in Red territory, she would prefer their Second Company area (where enemy soldiers are not as well-trained as in the Headquarters area).

(MAP) Judy has no idea where she is. Because of the map, her probability of being in Blue territory equals the probability of being in Red territory, and the probability of being in the Red Second Company area equals the probability of being in the Red Headquarters area.

(HDQ) Her commanders inform Judy that in case she is in Red territory, her chance of being in their Headquarters area is three times the chance of being in their Second Company area.

In formal terms (sloppily writing  $A_i$  for the event of Judy being in  $A_i$ ),

$$2 \cdot P(A_1) = 2 \cdot P(A_2) = P(A_3) \quad (\text{MAP})$$

$$\vartheta = P(A_2|A_1 \cup A_2) = \frac{3}{4} \quad (\text{HDQ})$$

(HDQ) is partial information because in contrast to the kind of evidence we are used to in Bayes' formula (such as 'an even number was rolled'), and to the kind of evidence needed for Jeffrey's rule (where a partition of the whole event space and its probability redistribution is required, not only  $A_1 \cup A_2$ , but see here the objections in Douven and Romeijn (2009)), the scenario suggests that Bayesian conditionalization and Jeffrey's rule are inapplicable. We are interested in the most defensible updated probability assignment(s) and will express them in the form of a normalized odds vector  $(q_1, q_2, q_3)$ , following Van Fraassen (1981).  $q_i$  is the updated probability  $Q(A_i)$  that Judy Benjamin is in  $A_i$ . Let  $P$  be the probability distribution prior to the new observation and  $p_i$  the individual 'prior' probabilities. These probabilities are not to be confused with prior probabilities that precede any kind of information. In the spirit of probability update, or probability kinematics, we will for the rest of the article refer to prior probabilities as probabilities prior to an observation and the subsequent update. The  $q_i$  sum to 1 (this differs from van Fraassen's canonical odds vector, which is proportional to the normalized odds vector but has 1 as its first element). We define

$$t = \frac{\vartheta}{1 - \vartheta},$$

$t$  is the factor by which (HDQ) indicates that Judy's chance of being in  $A_2$  is greater than being in  $A_1$ . In Judy's particular case,  $t = 3$  and  $\vartheta = 0.75$ . Two intuitions guide the way people think about Judy Benjamin's situation.

**T1** (HDQ) does not refer to Blue territory and should not affect  $P(A_3) : q_3 = p_3 (= 0.50)$ .

There is another, conflicting intuition (due to Peter Williams via personal communication with van Fraassen, see van Fraassen (1981, p. 379)):

**T2** If the value of  $\vartheta$  approaches 1 (in other words,  $t$  approaches infinity) then  $q_3$  should approach  $2/3$  as the problem reduces to one of ordinary conditioning. (HDQ) would turn into 'if you are in Red territory you are almost certainly in the Red Headquarters area.' Considering (MAP),  $q_3$  should approach  $2/3$ .

Continuity considerations pose a contradiction to T1. (These considerations are strong enough that Luc Bovens uses them as an assumption to solve Adam Elga's Sleeping Beauty problem by parity of reasoning in Bovens (2010).) To parse these conflicting intuitions, we will introduce several methods to provide  $G$ , the function that maps  $\vartheta$  to the appropriate normalized updated odds vector  $(q_1, q_2, q_3)$ .

The first method is extremely simple and accords with intuition T1:  $G_{\text{ind}}(\vartheta) = (0.5(1 - \vartheta), 0.5\vartheta, 0.5)$ . In Judy's particular case with  $t = 3$  the normalized odds vector is (ind stands for independent):

$$G_{\text{ind}}(0.75) = (0.125, 0.375, 0.500).$$

Both Grove and Halpern (1997) and Douven and Romeijn (2009) make a case for this distribution. Grove and Halpern use standard conditioning on the event of the message being transmitted to Judy. Douven and Romeijn use Jeffrey's rule (because they believe that T1 is in this case so strong that  $Q(A_3) = P(A_3)$  is as much of a constraint as (MAP) and (HDQ), yielding a Jeffrey partition).

T1, however, conflicts with the symmetry requirements outlined in Fraassen (1986). Van Fraassen introduces various updating methods which do not conflict with those symmetry requirements, the most notable of which is MAXENT. Shore and Johnson have already shown that, given certain assumptions (which have been heavily criticized, e.g. in Uffink (1996)), MAXENT produces the unique updated probability assignment according with these assumptions. The minimum information discrimination theorem of Kullback and Leibler (see for example Csiszár 1967, Sect. 3) demonstrates how Shannon's entropy and the Kullback–Leibler Divergence formula can provide the least informative updated probability assignment (with reference to the prior probability assignment) obeying the constraint posed by the evidence. The idea is to define a space of probability distributions, make sure that the constraint identifies a closed, convex subset in this space, and then determine which of the distributions in the closed, convex subset is least distant from the prior probability distribution in terms of information (using the minimum information discrimination theorem). It is necessary for the uniqueness of this least distant distribution that the subset be closed and convex (in other words, that the constraint be affine, see Csiszár (1967)).

For Judy Benjamin, MAXENT suggests the following normalized odds vector:

$$G_{\max}(0.75) \approx (0.117, 0.350, 0.533). \quad (1)$$

The updated probability of being on Blue territory ( $A_3$ ) has increased from 50% to approximately 53%. Grove and Halpern find this result “highly counterintuitive” (Grove and Halpern 1997, p. 2). Van Fraassen summarizes the worry:

It is hard not to speculate that the dangerous implications of being in the enemy’s Headquarters area are causing Judy Benjamin to indulge in wishful thinking, her indulgence becoming stronger as her conditional estimate of the danger increases. (Van Fraassen 1981, p. 379)

There are two ways in which we can arrive at result (1). We may use Jaynes’ constraint rule and find the updated probability distribution that is both least informative with respect to Shannon’s entropy and in accordance with the constraint (using Dempster’s Rule of Combination, which together with the constraint rule is equivalent to the principle of minimum cross-entropy, see Cover and Thomas (2006, p. 409, exercise 12.2).). Alternatively, if circumstances are favourable (as they are in Judy Benjamin’s case), we may use the Kullback–Leibler Divergence and differentiate it to obtain where it is minimal.

The constraint rule has the advantage of providing results when the derivative of the Kullback–Leibler Divergence is difficult to find. This not being the case for Judy, we go the easier route of the second method and provide a more general justification for the constraint rule in Appendix A, together with its application to the Judy Benjamin case.

The Kullback–Leibler Divergence is

$$D(Q, P) = \sum_{i=1}^m q_i \log_2 \frac{q_i}{p_i}.$$

We fill in the explicit details from Judy Benjamin’s situation and differentiate the expression to obtain the minimum (by setting the derivative to 0).

$$\frac{\partial}{\partial q_1} (q_1 \log_2(4q_1) + tq_1 \log_2(4tq_1) + (1 - (t+1)q_1) \log_2 2(1 - (t+1)q_1)) = 0.$$

The resulting expression for  $G_{\max}$  is

$$G_{\max}(\vartheta) = \left( \frac{C}{1 + Ct + C}, t \frac{C}{1 + Ct + C}, 1 - (t+1) \frac{C}{1 + Ct + C} \right),$$

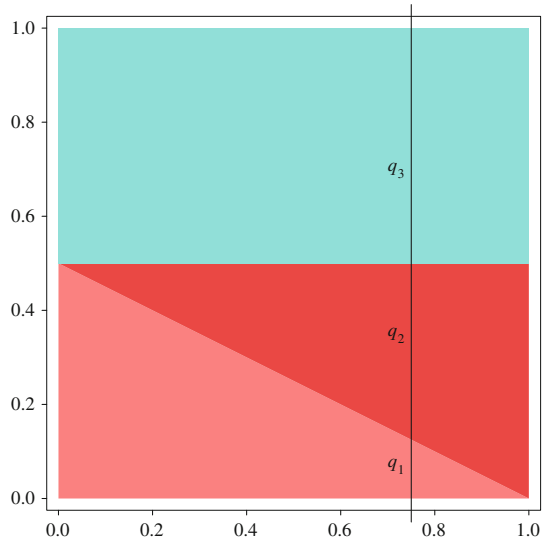
where

$$C = 2^{-\frac{t \log_2 t + t + 1}{1+t}}.$$

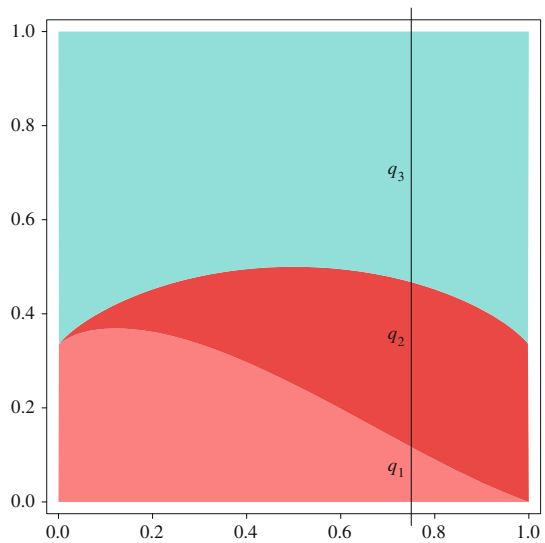
Figures 3 and 4 show in diagram form the distribution of  $(q_1, q_2, q_3)$  depending on the value of  $\vartheta$  (between 0 and 1), respectively following intuition T1 and MAXENT.



**Fig. 3** Judy Benjamin's updated probability assignment according to intuition T1.  $0 < \vartheta < 1$  forms the *horizontal axis*, the *vertical axis* shows the updated probability distribution (or the normalized odds vector)  $(q_1, q_2, q_3)$ . The *vertical line* at  $\vartheta = 0.75$  shows the specific updated probability distribution  $G_{\text{ind}}(0.75)$  for the Judy Benjamin problem



**Fig. 4** Judy Benjamin's updated probability assignment using MAXENT.  $0 < \vartheta < 1$  forms the *horizontal axis*, the *vertical axis* shows the updated probability distribution (or the normalized odds vector)  $(q_1, q_2, q_3)$ . The *vertical line* at  $\vartheta = 0.75$  shows the specific updated probability distribution  $G_{\text{max}}(0.75)$  for the Judy Benjamin problem



Notice that in accordance with intuition T2, MAXENT provides a result where  $q_3 \rightarrow 2/3$  for  $\vartheta$  approaching 0 or 1.

### 3 Epistemic entrenchment

Consider two future events  $A$  and  $B$ . You have partial belief in whether they will occur and assign probabilities to them. Then you learn that  $A$  entails  $B$ . How does this information affect your probability assignment for event  $A$ ? If  $A$  is causally independent of

*B* then your updated probability for it should equal the original probability. For example, whether Sarah and Marian have sundowners at the Westcliff hotel tomorrow may initially not depend on the weather at all, but if they learn that there will be a wedding and the hotel's indoor facilities will be closed to the public, then rainfall tomorrow implies no sundowners. Learning this conditional does not affect the probability of the antecedent (rainfall tomorrow), because the antecedent is causally independent of the consequent.

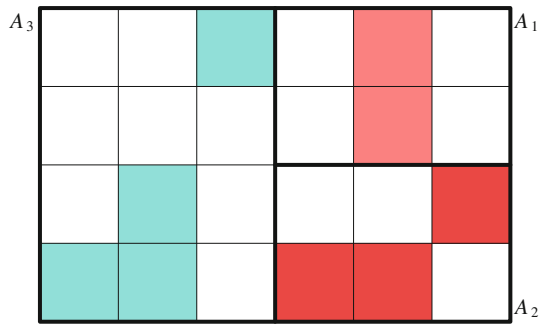
Here is another example (the two examples are from Douven and Romeijn 2009). A jeweler has been robbed, and Kate has reason to assume that Henry might be the robber. Kate knows that he is not capable of actually injuring another person, but he may very well engage in robbery. When Kate hears from the investigator that the robber also shot the jeweler she concludes that Henry is not the robber. She has learned a conditional and adjusted the probability for the antecedent. The reason for this is that Kate was epistemically entrenched to uphold her belief in Henry's nonviolent nature. Updating probabilities upon learning a conditional depends on epistemic entrenchments.

In Judy Benjamin's case, (HDQ) is also a conditional. If Judy is in Red territory, she is more likely to be in the Headquarters area. According to MAXENT, the updated probability for the antecedent of this conditional is slightly lowered (to about 47 %). It appears that MAXENT preempts our epistemic entrenchments and *nolens volens* assigns a certain degree of confirmation to the antecedent of a learned conditional. This degree of confirmation depends on the causal dependency of the antecedent on the consequent. The compatibility of epistemic entrenchments and MAXENT is material for another paper, but in this section we will focus on the independence assumptions that are improperly imported into the Judy Benjamin case by detractors of MAXENT. We will be particularly critical of Douven and Romeijn, who hold that the Judy Benjamin case is a case for Adam's conditioning, where the antecedent is left alone in the updating of probabilities.

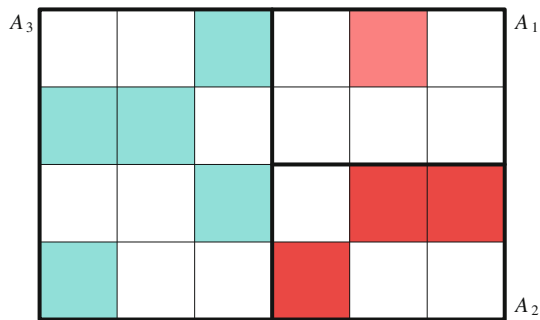
Even though T1 is an understandably strong intuition, it does not take into account that the information given to Judy by her commanders may indeed be dependent on whether she is in Blue or in Red territory. To underline this objection to intuition T1 consider three scenarios, any of which may form the basis of the partial information provided by her commanders.

- I Judy is dropped off by a pilot who flips two coins. If the first coin lands H, then Judy is dropped off in Blue territory, otherwise in Red territory. If the second coin lands H, she is dropped off in the Headquarters area, otherwise in the Second Company area. Judy's commanders find out that the second coin is biased  $\vartheta : 1 - \vartheta$  toward H with  $\vartheta = 0.75$ . The normalized odds vector is  $G_I(0.75) = (0.125, 0.375, 0.500)$  and agrees with T1, because the choice of Blue or Red is completely independent from the choice of the Red Headquarters area or the Red Second Company area.
- II The pilot randomly lands in any of the four quadrants and rolls a die. If she rolls an even number, she drops off Judy. If not, she takes her to another (or the same, the choice happens with replacement) randomly selected quadrant to repeat the procedure. Judy's commanders find out, however, that for  $A_1$ , the pilot requires a six to drop off Judy, not just an even number. The normalized odds vector in this scenario is  $G_{II}(0.75) = (0.1, 0.3, 0.6)$  and does not accord with T1.

**Fig. 5** This choice of rectangles is not a candidate because the number of rectangles in  $A_2$  is not a  $t$ -multiple of the number of rectangles in  $A_1$ , here with  $s = 2, t = 3$  as in scenario III



**Fig. 6** This choice of rectangles is a candidate because the number of rectangles in  $A_2$  is a  $t$ -multiple of the number of rectangles in  $A_1$ , here with  $s = 2, t = 3$  as in scenario III



**III** Judy's commanders have divided the map into 24 congruent rectangles,  $A_3$  into 12, and  $A_1$  and  $A_2$  into 6 rectangles each (see Figs. 5, 6). They have information that the only subsets of the 24 rectangles in which Judy Benjamin may be located are such that they contain three times as many  $A_2$  rectangles than  $A_1$  rectangles. The normalized odds vector in this scenario is  $G_{\text{III}}(0.75) \approx (0.108, 0.324, 0.568)$  (evaluating almost 17 million subsets).

I–III demonstrate the contrast between scenarios when independence is true and when it is not. Douven and Romeijn's capital mistake in their paper is that they assume that the Judy Benjamin problem is analogous to their example of Sarah and sundowners at the Westcliff (see Douven and Romeijn 2009, p. 7). Sarah, however, knows that whether it rains or not is independent of her activity the next night, whereas in Judy Benjamin we have no evidence of such independence, as scenario II makes clear. This is not to say that scenario II is the scenario that pertains in Judy Benjamin's case. It only says that there is no natural assumption in Judy Benjamin's case that the probabilities are independent of each other in light of the new evidence, for scenario II is perfectly natural (whether it is true or not is a completely different question) and reveals how dependence is consistent with the information that Judy Benjamin receives.

Douven and Romeijn's strong independence claim relying on intuition T1 leads them to apply Jeffrey's rule to the Judy Benjamin problem with the additional constraint  $q_3 = p_3$ . They claim that in most cases "the learning of a conditional is or would be irrelevant to one's degree of belief for the conditional's antecedent ... the learning

of the relevant conditional should intuitively leave the probability of the antecedent unaltered” (Douven and Romeijn 2009, p. 9).

This, according to Douven and Romeijn, is the usual epistemic entrenchment and applies in full force to the Judy Benjamin problem. They give an example where the epistemic entrenchment could go the other way and leave the consequent rather than the antecedent unaltered (Kate and Henry, see Douven and Romeijn 2009, p. 13). The idea of epistemic entrenchment is at odds with MAXENT and seems to imply just what the full employment theorem claims: judgments so framed “will depend on the judgmental skills of the agent, typically acquired not in the inductive logic class but by subject specific training” (Bradley 2005, p. 349). To pursue the relation between the semantics of conditionals, MAXENT, and the full employment theorem would take us too far afield at present and shall be undertaken elsewhere. For the Judy Benjamin problem, it is not clear why Douven and Romeijn think that the way the problem is posed implies a strong epistemic entrenchment for Adams conditioning (Adams conditioning is the kind of conditioning that will leave the antecedent alone). Scenarios II–III provide realistic alternatives where Adams conditioning is inappropriate.

Judy Benjamin may also receive (HDQ) because her informers have found out that Red Headquarters troops have occupied the entire Blue territory ( $q_1 = 3p_1, q_2 = p_2, q_3 = 0$ , the epistemic entrenchment is with respect to  $q_2$ ); because they have found out that Blue troops have occupied two-thirds of the Red Second Company area ( $q_1 = p_1, q_2 = (1/3)p_2, q_3 = (4/3)p_3$ , the epistemic entrenchment is with respect to  $q_1$ ); or because they have found out that Red Headquarters troops have taken over half of the Red Second Company area ( $q_1 = (1/2)p_1, q_2 = (3/2)p_2, q_3 = p_3$ , the epistemic entrenchment is with respect to  $q_3$  and what Douven and Romeijn take to be an assumption in the wording of the problem). There is nothing in the problem that supports Douven and Romeijn’s narrowing of the options. The table reiterates these options, with the third, shaded line representing intuition (T1) and the epistemic entrenchment defended by Douven and Romeijn.

Epistemic entrenchment	$q_1$	$q_2$	$q_3$
with respect to $A_1$	1/4	3/4	0
with respect to $A_2$	1/12	1/4	2/3
with respect to $A_3$	1/8	3/8	1/2

#### 4 Coarsening at random

Another at first blush forceful argument that MAXENT’s solution for the Judy Benjamin problem is counterintuitive has to do with coarsening at random, or CAR for short. CAR involves using more naive (or coarse) event spaces in order to arrive at solutions to probability updating problems. The mechanics are spelled out in Grünwald and Halpern (2003). Grünwald and Halpern see a parallel between the Judy Benjamin problem and Martin Gardner’s Three Prisoners problem (see Gardner 1959, 180f). In the Three Prisoners problem, three men (A, B, and C) are under sentence of death when the governor decides to pardon one of them. The warden of the prison knows which of the three men is pardoned, but none of the men do. In a private conversation,

A says to the warden, Tell me the name of one of the others who will be executed—it will not give anything away whether I will be executed or not. The warden agrees and tells A that B will be executed. For the puzzling consequences, see the wealth of literature on the Three Prisoners problem or the Monty Hall problem.

According to Grünwald and Halpern, for problems of this kind (Judy Benjamin, Three Prisoners, Monty Hall) there are naive and sophisticated spaces to which we can apply probability updates. If A uses the naive space, for example, he comes to the following conclusion: of the three possibilities that (A,B), (A,C), or (B,C) are executed, the warden's information excludes (A,C). (A,B) and (B,C) are left over, and because A has no information about which one of these is true his chance of not being executed is 0.5. His chance of survival has increased from one third to one half.

Grünwald and Halpern show, correctly, that the application of the naive space is illegitimate because the CAR condition does not hold. More generally, Grünwald and Halpern show that updating on the naive space rather than the sophisticated space is legitimate for event type observations always when the set of observations is pairwise disjoint or, when the events are arbitrary, only when the CAR condition holds. For Jeffrey type observations, there is a generalized CAR condition which applies likewise. For affine constraints on which we cannot use Jeffrey conditioning (or, a fortiori, standard conditioning) MAXENT “essentially never gives the right results” (Grünwald and Halpern 2003, p. 243).

Grünwald and Halpern conclude that “working with the naive space, while an attractive approach, is likely to give highly misleading answers” (p. 246), especially in the application of MAXENT to naive spaces as in the Judy Benjamin case “where applying [MAXENT] leads to paradoxical, highly counterintuitive results” (p. 245). For the Three Prisoners problem, Jaynes' constraint rule would supposedly proceed as follows: the vector of prior probabilities for (A,B), (A,C), and (B,C) is  $(1/3, 1/3, 1/3)$ . The constraint is that the probability of (A,C) is zero, and a simple application of the constraint rule yields  $(1/2, 0, 1/2)$  for the vector of updated probabilities. The CAR condition for the naive space does not hold, therefore the result is misleading.

By analogy, using the constraint rule on the naive space for the Judy Benjamin problem yields  $(0.117, 0.350, 0.533)$ , but as the CAR condition fails in even the simplest settings for affine constraints (“CAR is (roughly speaking) guaranteed *not* to hold except in ‘degenerate’ situations” (p. 251), emphasis in the original), it certainly fails for the Judy Benjamin problem, for which constructing a sophisticated space is complicated (see Grove and Halpern (1997), where the authors attempt such a construction by retrospective conditioning).

The analogy, however, is misguided. The constraint rule has been formally shown to generalize Jeffrey conditioning, which in turn has been shown to generalize standard conditioning (the authors admit as much in Grünwald and Halpern 2003, p. 262). We can solve both the Monty Hall problem and the Three Prisoners problem by standard conditioning, not using the naive space, but simply using the correct space for the probability update. For the Three Prisoners problem, for example, the warden will say either ‘B’ or ‘C’ in response to A's inquiry. Because A has no information that would privilege either answer the probability that the warden says ‘B’ and the probability that the warden says ‘C’ equal each other and therefore equal 0.5. Here is the difference

between using the naive space and using the correct space, but either way using standard conditional probabilities:

$$\begin{aligned}
 &P(\text{'A is pardoned'} | \text{'B will be executed'}) = \\
 &\quad \frac{P(\text{'A is pardoned'})}{P(\text{'A is pardoned'}) + P(\text{'C is pardoned'})} = \frac{1}{2} \text{ (incorrect)} \\
 &P(\text{'A is pardoned'} | \text{'warden says B will be executed'}) = \\
 &\quad \frac{P(\text{'A is pardoned' and 'warden says B will be executed'})}{P(\text{'warden says B will be executed'})} = \frac{1/6}{1/2} = \frac{1}{3} \text{ (correct).}
 \end{aligned}$$

Why is the first equation incorrect and the second one correct? Information theory gives us the right answer: in the first equation, we are conditioning on a watered down version of the evidence (watered down in a way that distorts the probabilities because we are not ‘coarsening at random’). ‘Warden says B will be executed’ is sufficient but not necessary for ‘B will be executed.’ The former proposition is more informative than the latter proposition (its probability is lower). Conditioning on the latter proposition leaves out relevant information contained in the wording of the problem.

Because MAXENT always agrees with standard conditioning, MAXENT gives the correct result for the Three Prisoners problem. For the Judy Benjamin problem, there is no defensible sophisticated space and no watering down of the evidence in what Grünwald and Halpern call the ‘naive’ space. The analogy between the Three Prisoners problem and the Judy Benjamin problem as it is set up by Grünwald and Halpern fails because of this crucial difference. A successful criticism would be directed at the construction of the ‘naive’ space: this is what we just accomplished for the Three Prisoners problem. There is no parallel procedure for the Judy Benjamin problem. The ‘naive’ space is all we have, and MAXENT is the appropriate tool to deal with this lack of information.

## 5 The powerset approach

In this section, we will focus on scenario III and consider what happens when the grain of the partition becomes finer. We call this the powerset approach. Two remarks are in order. On the one hand, the powerset approach has little independent appeal. The reason behind using MAXENT is that we want our evidence to have just the right influence on our updated probabilities, i.e. neither over-inform nor under-inform. There is no corresponding reason why we should update our probabilities using the powerset approach. On the other hand, what the powerset approach does is lend support to another approach. In this task, it is persuasive because it tells us what would happen if we were to divide the event space into infinitesimally small, uniformly weighed, and independent ‘atomic’ bits of information.

In the process of arriving at the formal result, the powerset approach resembles an empirical experiment. We are making many assumptions favouring T1, but when the result comes in it supports T2 in astonishingly non-trivial ways. The powerset approach provides support for MAXENT against T1 because it combines the assumptions grounding T1 with a limit construction and still yields a solution that

closely approximates the one generated by MAXENT rather than the one generated by T1.

On its own the powerset approach is just what Grünwald and Halpern call a naive space, for which CAR does not hold. Hence the powerset approach will not give us a precise solution for the problem, although it may with some plausibility guide us in the right direction—especially if despite all its independence and uniformity assumptions it significantly disagrees with intuition T1.

Let us assume a partition of the Blue and Red territories into sets of equal measure (this is the division into rectangles of scenario III). (MAP) dictates that the number of sets covering  $A_3$  equals the number of sets covering  $A_1 \cup A_2$ . Initially, any subset of this partition is a candidate for Judy Benjamin to consider. The constraint imposed by (HDQ) is that now we only consider subsets for which there are three times as many partition sets (or rectangles, although we are not necessarily limiting ourselves to rectangles) in  $A_2$  as there are in  $A_1$ . In Figs. 5 and 6 there are diagrams of two subsets. One of them (Fig. 5) is not a candidate, the other one (Fig. 6) is.

Let  $X$  be the random variable that corresponds to the ratio of the number of partition elements (rectangles) that are in  $A_3$  to the total number of partition elements (rectangles) for a randomly chosen candidate. We would now anticipate that the expectation of  $X$  (which we will call  $EX$ ) gives us an indication of the updated probability that Judy is in  $A_3$  (so  $EX \approx q_3$ ). The powerset approach is often superior to the uniformity approach (Grove and Halpern use uniformity, with all the necessary qualifications): if you have played Monopoly, you will know that the frequencies for rolling a 2, a 7, or a 10 with two dice tend to conform more closely to the binomial distribution (based on a powerset approach) rather than to the uniform distribution with  $P(\text{rolling } i) = 1/11$  for  $i = 2, \dots, 12$ .

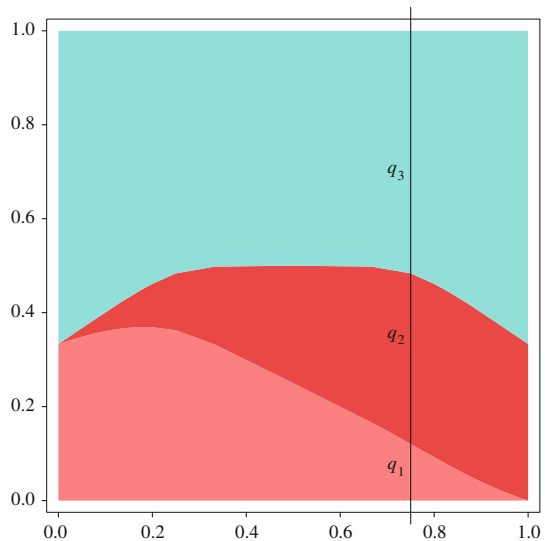
Appendix B provides a formula for the powerset approach (corresponding to  $G_{\text{ind}}$  and  $G_{\text{pws}}$ ), giving us  $q_3$  dependent on  $t$ . Notice that this formula is for  $t = 2, 3, 4, \dots$ . For  $t = 1$  use the Chu-Vandermonde identity to find that

$$EX_{12} = (t+1) \frac{\sum_{i=1}^s i \binom{ts}{i} \binom{ts}{ti}}{\sum_{i=0}^s \binom{ts}{i} \binom{ts}{ti}} = (t+1) \frac{s}{2}$$

and consequently  $EX = 1/2$ , as one would expect. For  $t = 1/2, 1/3, 1/4, \dots$  we can simply reverse the roles of  $A_1$  and  $A_2$ . These results give us  $G_{\text{pws}}$  and a graph of the normalized odds vector (see Fig. 7), a bit bumpy around the middle because the  $t$ -values are discrete and farther apart in the middle, as  $t = \vartheta/(1-\vartheta)$ . Comparing the graphs of the normalized odds vector under Grove and Halpern's uniformity approach ( $G_{\text{ind}}$ ), Jaynes' MAXENT approach ( $G_{\text{max}}$ ), and the powerset approach suggested in this paper ( $G_{\text{pws}}$ ), it is clear that the powerset approach supports MAXENT.

Going through the calculations, it seems at many places that the powerset approach should give its support to Grove and Halpern's uniformity approach in keeping with intuition T1. It is unexpected to find out that in the mathematical analysis the parameters converge to a non-trivial factor and do not tend to negative or positive infinity. Most surprisingly, the powerset approach, prima facie unrelated to an approach using information, supports the idea that a set of events about which nothing is known (such

**Fig. 7** Judy Benjamin's updated probability assignment according to the powerset approach.  $0 < \vartheta < 1$  forms the horizontal axis, the vertical axis shows the updated probability distribution (or the normalized odds vector)  $(q_1, q_2, q_3)$ . The vertical line at  $\vartheta = 0.75$  shows the specific updated probability distribution  $G_{\text{pws}}$  for the Judy Benjamin problem



as  $A_3$ ) gains in probability in the updated probability distribution compared to the set of events about which something is known (such as  $A_1$  and  $A_2$ ), even if it is only partial information. Unless independence is specified, as in Sarah and sundowners at the Westcliff, the area of ignorance gains compared to the area of knowledge.

We now have several ways to characterize Judy's updated probabilities and updated probabilities following upon partial information in general. Only one of them, the uniformity approach, violates van Fraassen, Hughes, and Harman's five symmetry requirements in Fraassen (1986) and intuition T2. The uniformity approach, however, is the only one that satisfies intuition T1, an intuition which most people have when they first hear the story.

Two arguments attenuate the position of the uniformity approach in comparison with the others. First, T1 rests on an independence assumption which is not reflected in the problem. Although there is no indication that what Judy's commanders tell her is in any way dependent on her probability of being in Blue territory, it is not excluded either (see scenarios II and III earlier in this paper). MAXENT takes this uncertainty into consideration. Second, when we investigate the problem using the powerset approach it turns out that a division into equally probable, independent, and increasingly fine bits of information supports not intuition T1 but rather intuition T2. MAXENT, for now, is vindicated. We need to look for full employment not by cleverly manipulating prior probabilities, but by making fresh observations, designing better experiments, and partitioning the theory space more finely.

## Appendix A: Jaynes' constraint rule

This appendix provides a concise but comprehensive summary of Jaynes' constraint rule not easily obtainable in the literature (although the constraint rule is a



straightforward application of Lagrange multipliers). Jaynes applies it to the Brandeis Dice Problem (see [Jaynes 1989](#), p. 243), but does not give a mathematical justification.

Let  $f$  be a probability distribution on a finite space  $x_1, \dots, x_m$  that fulfills the constraint

$$\sum_{i=1}^m r(x_i) f(x_i) = \alpha. \quad (2)$$

An affine constraint can always be expressed by assigning a value to the expectation of a probability distribution (see [Hobson 1971](#)). In Judy Benjamin's case, for example, let  $r(x_1) = 0, r(x_2) = 1 - \vartheta, r(x_3) = -\vartheta$  and  $\alpha = 0$ . Because  $f$  is a probability distribution it fulfills

$$\sum_{i=1}^m f(x_i) = 1. \quad (3)$$

We want to maximize Shannon's entropy, given the constraints (2) and (3),

$$-\sum_{i=1}^m f(x_i) \ln(x_i). \quad (4)$$

We use Lagrange multipliers to define the functional

$$J(f) = -\sum_{i=1}^m f(x_i) \ln f(x_i) + \lambda_0 \sum_{i=1}^m f(x_i) + \lambda_1 \sum_{i=1}^m r(x_i) f(x_i)$$

and differentiate it with respect to  $f(x_i)$

$$\frac{\partial J}{\partial f(x_i)} = -\ln(f(x_i)) - 1 + \lambda_0 + \lambda_1 r(x_i). \quad (5)$$

Set (5) to 0 to find the necessary condition to maximize (4)

$$g(x_i) = e^{\lambda_0 - 1 + \lambda_1 r(x_i)}.$$

This is the Gibbs distribution. We still need to do two things: (a) show that the entropy of  $g$  is maximal, and (b) show how to find  $\lambda_0$  and  $\lambda_1$ . (a) is shown in Theorem 12.1.1 in [Cover and Thomas \(2006\)](#) and there is no reason to copy it here.

For (b), let

$$\begin{aligned} \lambda_1 &= -\beta \\ Z(\beta) &= \sum_{i=1}^m e^{-\beta r(x_i)} \\ \lambda_0 &= 1 - \ln(Z(\beta)). \end{aligned}$$

To find  $\lambda_0$  and  $\lambda_1$  we introduce the constraint

$$-\frac{\partial}{\partial \beta} \ln(Z(\beta)) = \alpha.$$

To see how this constraint gives us  $\lambda_0$  and  $\lambda_1$ , Jaynes' solution of the Brandeis Dice Problem (see [Jaynes 1989](#), p. 243) is a helpful example. We are, however, interested in a general proof that this choice of  $\lambda_0$  and  $\lambda_1$  gives us the probability distribution maximizing the entropy. That  $g$  so defined maximizes the entropy is shown in (a). We need to make sure, however, that with this choice of  $\lambda_0$  and  $\lambda_1$  the constraints (2) and (3) are also fulfilled (a standard result in the application of Lagrange multipliers).

First, we show

$$\begin{aligned} \sum_{i=1}^m g(x_i) &= \sum_{i=1}^m e^{\lambda_0 - 1 + \lambda_1 r(x_i)} = e^{\lambda_0 - 1} \sum_{i=1}^m e^{\lambda_1 r(x_i)} \\ &= e^{-\ln(Z(\beta))} Z(\beta) = 1. \end{aligned}$$

Then, we show, by differentiating  $\ln(Z(\beta))$  using the substitution  $x = e^{-\beta}$

$$\begin{aligned} \alpha &= -\frac{\partial}{\partial \beta} \ln(Z(\beta)) = -\frac{1}{\sum_{i=1}^m x^{r(x_i)}} \left( \sum_{i=1}^m r(x_i) x^{r(x_i)-1} \right) (-x) \\ &= \frac{\sum_{i=1}^m r(x_i) x^{r(x_i)}}{\sum_{i=1}^m x^{r(x_i)}}. \end{aligned}$$

And, finally,

$$\begin{aligned} \sum_{i=1}^m r(x_i) g(x_i) &= \sum_{i=1}^m r(x_i) e^{\lambda_0 - 1 + \lambda_1 r(x_i)} = e^{\lambda_0 - 1} \sum_{i=1}^m r(x_i) e^{\lambda_1 r(x_i)} \\ &= e^{\lambda_0 - 1} \sum_{i=1}^m r(x_i) x^{r(x_i)} = \alpha e^{\lambda_0 - 1} \sum_{i=1}^m x^{r(x_i)} = \alpha e^{\lambda_0 - 1} \sum_{i=1}^m e^{-\beta r(x_i)} \\ &= \alpha Z(\beta) e^{\lambda_0 - 1} = \alpha Z(\beta) e^{-\ln(Z(\beta))} = \alpha. \end{aligned}$$

Filling in the variables from Judy Benjamin's scenario gives us result (1). The lambdas are:

$$\lambda_0 = 1 - \ln \left( \sum_{i=1}^m e^{\lambda_1 r(x_i)} \right) \quad \lambda_1 = \ln \vartheta - \ln(1 - \vartheta).$$

We combine the normalized odds vector (0.16, 0.48, 0.36) following from these lambdas using Dempster's Rule of Combination with ([MAP](#)) and get result (1).

## 6 The powerset approach formalized

Let us assume a partition  $\{B_i\}_{i=1,\dots,4n}$  of  $A_1 \cup A_2 \cup A_3$  into sets that are of equal measure  $\mu$  and whose intersection with  $A_i$  is either the empty set or the whole set itself (this is the division into rectangles of scenario III). (MAP) dictates that the number of sets covering  $A_3$  equals the number of sets covering  $A_1 \cup A_2$ . For convenience, we assume the number of sets covering  $A_1$  to be  $n$ . Let  $\mathcal{C}$ , a subset of the powerset of  $\{B_i\}_{i=1,\dots,4n}$ , be the collection of sets which agree with the constraint imposed by (HDQ), i.e.

$$C \in \mathcal{C} \text{ iff } C = \{C_j\} \text{ and } t\mu\left(\bigcup C_j \cap A_1\right) = \mu\left(\bigcup C_j \cap A_2\right).$$

In Figs. 5 and 6 there are diagrams of two elements of the powerset of  $\{B_i\}_{i=1,\dots,4n}$ . One of them (Fig. 5) is not a member of  $\mathcal{C}$ , the other one (Fig. 6) is.

The binomial distribution dictates the expectation  $EX$  of  $X$ , using simple combinatorics. In this case we require, again for convenience, that  $n$  be divisible by  $t$  and the ‘grain’ of the partition be  $s = n/t$ . Remember that  $t$  is the factor by which (HDQ) indicates that Judy’s chance of being in  $A_2$  is greater than being in  $A_1$ . In Judy’s particular case,  $t = 3$  and  $\vartheta = 0.75$ . We introduce a few variables which later on will help for abbreviation:

$$n = ts \quad 2m = n \quad 2j = n - 1 \quad T = t^2 + 1.$$

$EX$ , of course, depends both on the grain of the partition and the value of  $t$ . It makes sense to make it independent of the grain by letting the grain become increasingly finer and by determining  $EX$  as  $s \rightarrow \infty$ . This cannot be done for the binomial distribution, as it is notoriously uncomputable for large numbers (even with a powerful computer things get dicey around  $s = 10$ ). But, equally notorious, the normal distribution provides a good approximation of the binomial distribution and will help us arrive at a formula for  $G_{\text{pws}}$  (corresponding to  $G_{\text{ind}}$  and  $G_{\text{max}}$ ), determining the value  $q_3$  dependent on  $\vartheta$  as suggested by the powerset approach.

First, we express the random variable  $X$  by the two independent random variables  $X_{12}$  and  $X_3$ .  $X_{12}$  is the number of partition elements in the randomly chosen  $C$  which are either in  $A_1$  or in  $A_2$  (the random variable of the number of partition elements in  $A_1$  and the random variable of the number of partition elements in  $A_2$  are decisively not independent, because they need to obey (HDQ));  $X_3$  is the number of partition elements in the randomly chosen  $C$  which are in  $A_3$ . A relatively simple calculation shows that  $EX_3 = n$ , which is just what we would expect (either the powerset approach or the uniformity approach would give us this result):

$$EX_3 = 2^{-2n} \sum_{i=0}^{2n} i \binom{2n}{i} = n \left( \text{use } \binom{n}{k} = \frac{n}{k} \binom{n-1}{k-1} \right).$$

The expectation of  $X$ ,  $X$  being the random variable expressing the ratio of the number of sets covering  $A_3$  and the number of sets covering  $A_1 \cup A_2 \cup A_3$ , is

$$EX = \frac{EX_3}{EX_{12} + EX_3} = \frac{n}{EX_{12} + n}.$$

If we were able to use uniformity and independence,  $EX_{12} = n$  and  $EX = 1/2$ , just as Grove and Halpern suggest (although their uniformity approach is admittedly less crude than the one used here). Will the powerset approach concur with the uniformity approach, will it support the principle of maximum entropy, or will it make another suggestion on how to update the prior probabilities? To answer this question, we must find out what  $EX_{12}$  is, for a given value  $t$  and  $s \rightarrow \infty$ , using the binomial distribution and its approximation by the normal distribution.

Using combinatorics,

$$EX_{12} = (t+1) \frac{\sum_{i=1}^s i \binom{ts}{i} \binom{ts}{ti}}{\sum_{i=0}^s \binom{ts}{i} \binom{ts}{ti}}.$$

Let us call the numerator of this fraction NUM and the denominator DEN. According to the de Moivre–Laplace Theorem,

$$\text{DEN} = \sum_{i=0}^s \binom{ts}{i} \binom{ts}{ti} \approx 2^{2n} \sum_{i=0}^s \int_{i-\frac{1}{2}}^{i+\frac{1}{2}} \mathcal{N}\left(\frac{n}{2}, \frac{n}{4}\right)(i) \mathcal{N}\left(\frac{n}{2}, \frac{n}{4}\right)(ti) di,$$

where

$$\mathcal{N}(\mu, \sigma^2)(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right).$$

Substitution yields

$$\text{DEN} \approx 2^{2n} \frac{1}{\pi m} \sum_{i=0}^s \int_{i-\frac{1}{2}}^{i+\frac{1}{2}} \exp\left(-\frac{(x-m)^2}{m} - \frac{t^2(x-\frac{m}{t})^2}{m}\right) dx.$$

Consider briefly the argument of the exponential function:

$$-\frac{(x-m)^2}{m} - \frac{t^2(x-\frac{m}{t})^2}{m} = -\frac{t^2}{m}(a''x^2 + b''x + c'') = -\frac{t^2}{m}(a''(x-h'')^2 + k'')$$

with (the double prime sign corresponds to the simple prime sign for the numerator later on)

$$\begin{aligned} a'' &= \frac{1}{t^2}T & b'' &= (-2m)\frac{1}{t^2}(t+1) & c'' &= 2m^2\frac{1}{t^2} \\ h'' &= -b''/2a'' & k'' &= a''h''^2 + b''h'' + c''. \end{aligned}$$

Consequently,

$$\text{DEN} \approx 2^{2n} \exp\left(-\frac{t^2}{m}k''\right) \sqrt{\frac{1}{\pi a'' m t^2}} \int_{-\infty}^{s+\frac{1}{2}} \mathcal{N}\left(h'', \frac{m}{2a''t^2}\right) dx.$$

And, using the error function for the cumulative density function of the normal distribution,

$$\text{DEN} \approx 2^{2n-1} \sqrt{\frac{1}{\pi a'' m t^2}} \exp\left(-\frac{k''t^2}{m}\right) (1 - \text{erf}(w'')) \quad (6)$$

with

$$w'' = \frac{t\sqrt{a''}\left(s + \frac{1}{2} - h''\right)}{\sqrt{m}}.$$

We proceed likewise with the numerator, although the additional factor introduces a small complication:

$$\begin{aligned} \text{NUM} &= \sum_{i=1}^s i \binom{ts}{i} \binom{ts}{ti} = \sum_{i=1}^s s \binom{ts}{i} \binom{ts-1}{ti-1} \\ &\approx s 2^{2n-1} \sum_{i=1}^s \mathcal{N}\left(m, \frac{m}{2}\right)(i) \mathcal{N}\left(j, \frac{j}{2}\right)(ti-1). \end{aligned}$$

Again, we substitute and get

$$\text{NUM} \approx s 2^{2n-1} \left(\pi \sqrt{mj}\right)^{-1} \sum_{i=1}^{s-1} \int_{i-\frac{1}{2}}^{i+\frac{1}{2}} \exp\left(a'(x-h')^2 + k'\right),$$

where the argument for the exponential function is

$$-\frac{1}{mj} \left( j(x-m)^2 + mt^2 \left( x - \frac{j+1}{t} \right)^2 \right)$$

and therefore

$$\begin{aligned} a' &= j + mt^2 & b' &= 2j(1-m) + 2mt(t-j) & c' &= j(1-m)^2 + m(t-j-1)^2 \\ h' &= -b'/2a' & k' &= a'h'^2 + b'h' + c'. \end{aligned}$$

Using the error function,

$$\text{NUM} \approx 2^{2n-2} \frac{s}{\sqrt{\pi a'}} \exp\left(-\frac{k'}{mj}\right) (1 + \text{erf}(w')) \quad (7)$$

with

$$w' = \frac{\sqrt{a'} \left(s - \frac{1}{2} - h'\right)}{\sqrt{mj}}.$$

Combining (6) and (7),

$$\begin{aligned} EX_{12} &= (t+1) \frac{\text{NUM}}{\text{DEN}} \\ &\approx \frac{1}{2}(t+1) \sqrt{\frac{Tts}{Tts-1}} s e^{\alpha_{t,s}} \end{aligned}$$

for large  $s$ , because the arguments for the error function  $w'$  and  $w''$  escape to positive infinity in both cases (NUM and DEN) so that their ratio goes to 1. The argument for the exponential function is

$$\alpha_{t,s} = -\frac{k'}{mj} + \frac{k''t^2}{m}$$

and, for  $s \rightarrow \infty$ , goes to

$$\alpha_t = \frac{1}{2} T^{-2} (2t^3 - 3t^2 + 4t - 5).$$

Notice that, for  $t \rightarrow \infty$ ,  $\alpha_t$  goes to 0 and

$$EX = \frac{n}{EX_{12} + n} \rightarrow \frac{2}{3}$$

in accordance with intuition T2.

## References

- Bovens, L. (2010). Judy Benjamin is a sleeping beauty. *Analysis*, 70(1), 23–26.
- Bradley, R. (2005). Radical probabilism and Bayesian conditioning. *Philosophy of Science*, 72(2), 342–364.
- Caticha, A. (2012). *Entropic inference and the foundations of physics*. Sao Paulo: Brazilian Chapter of the International Society for Bayesian Analysis.
- Caticha, A., & Giffin, A. (2006). Updating probabilities. In *MaxEnt 2006, the 26th international workshop on Bayesian inference and maximum entropy methods*.
- Cover, T., & Thomas, J. (2006). *Elements of information theory* (Vol. 6). Hoboken, NJ: Wiley.
- Csiszár, I. (1967). Information-type measures of difference of probability distributions and indirect observations. *Studia Scientiarum Mathematicarum Hungarica*, 2, 299–318.

- Diaconis, P., & Zabell, S. (1982). Updating subjective probability. *Journal of the American Statistical Association*, 77, 822–830.
- Douven, I., & Romeijn, J. (2009). A new resolution of the Judy Benjamin problem. *CPNSS Working Paper*, 5(7), 1–22.
- Gardner, M. (1959). Mathematical games. *Scientific American*, 201(4), 174–182.
- Grove, A., & Halpern, J. (1997). Probability update: Conditioning vs. cross-entropy. In *Proceedings of the thirteenth conference on uncertainty in artificial intelligence*, Citeseer, Providence, RI.
- Grünwald, P. (2000). Maximum entropy and the glasses you are looking through. *Proceedings of the sixteenth conference on uncertainty in artificial intelligence* (pp. 238–246). Burlington: Morgan Kaufmann Publishers.
- Grünwald, P., & Halpern, J. (2003). Updating probabilities. *Journal of Artificial Intelligence Research*, 19, 243–278.
- Halpern, J. Y. (2003). *Reasoning about uncertainty*. Cambridge, MA: MIT Press.
- Hobson, A. (1971). *Concepts in statistical mechanics*. New York: Gordon and Beach.
- Howson, C., & Franklin, A. (1994). Bayesian Conditionalization and Probability Kinematics. *The British Journal for the Philosophy of Science*, 45(2), 451–466.
- Jaynes, E. (1989). *Papers on probability, statistics and statistical physics*. Dordrecht: Springer.
- Jaynes, E., & Bretthorst, G. (1998). *Probability theory: The logic of science*. Cambridge, UK: Cambridge University Press.
- Jeffrey, R. (1965). *The logic of decision*. New York: McGraw-Hill.
- Shore, J., & Johnson, R. (1980). Axiomatic derivation of the principle of maximum entropy and the principle of minimum cross-entropy. *IEEE Transactions on Information Theory*, 26(1), 26–37.
- Topsøe, F. (1979). Information-theoretical optimization techniques. *Kybernetika*, 15(1), 8–27.
- Uffink, J. (1996). The constraint rule of the maximum entropy principle. *Studies in History and Philosophy of Science Part B: Studies in History and Philosophy of Modern Physics*, 27(1), 47–79.
- Van Fraassen, B. (1981). A problem for relative information minimizers in probability kinematics. *The British Journal for the Philosophy of Science*, 32(4), 375–379.
- Van Fraassen, B. (1986). A problem for relative information minimizers, continued. *The British Journal for the Philosophy of Science*, 37(4), 453–463.
- Walley, P. (1991). *Statistical reasoning with imprecise probabilities*. London: Chapman and Hall.