

Maximum Entropy and Probability Kinematics Constrained by Conditionals

Stefan Lukits

Abstract

Two open questions of inductive reasoning are solved: (1) does the principle of maximum entropy (PME) give a solution to the obverse Majerník problem; and (2) is Wagner correct when he claims that Jeffrey's updating principle (JUP) contradicts PME? Majerník shows that PME provides unique and plausible marginal probabilities, given conditional probabilities. The obverse problem posed here is whether PME also provides such conditional probabilities, given certain marginal probabilities. The theorem developed to solve the obverse Majerník problem demonstrates that in the special case introduced by Wagner PME does not contradict JUP, but elegantly generalizes it and offers a more integrated approach to probability updating.

1 Introduction

Here is some bibliography help: [5], [9], [14], [13], [20], [16], [3], [11], [19], [23].

Sometimes, when we reason inductively, outcomes that are observed have entailment relationships with partitions of the possibility space that pose updating challenges which Jeffrey conditioning cannot meet. As we will see, it is not difficult to resolve these challenges by generalizing Jeffrey conditioning. There are claims in the literature that the principle of maximum entropy, from now on PME, conflicts with this generalization. We will show under which conditions this conflict obtains. Since proponents of PME are unlikely to subscribe to these conditions, the position of PME in the larger debate over inductive logic and reasoning is not undermined.

In Section 1, I will introduce the obverse Majerník problem and sketch how it ties in with two natural generalizations of Jeffrey conditioning: Wagner

conditioning and the PME. In Section 2, I will introduce Jeffrey conditioning in a notation that will later help us to solve the obverse Majerník problem. In Section 3, I will introduce Wagner conditioning and show how it naturally generalizes Jeffrey conditioning. In Section 4, finally, I will show that PME does so as well under conditions that are straightforward to accept for proponents of PME. This solves the obverse Majerník problem and makes Wagner conditioning unnecessary as a generalization of Jeffrey conditioning, since the PME seamlessly incorporates it. The conclusion summarizes my claims and briefly refers to epistemological consequences.

In his paper “Marginal Probability Distribution Determined by the Maximum Entropy Method” (see [15]), Vladimír Majerník asks the following question: If we had two partitions of an event space and knew all the conditional probabilities (any conditional probability of one event in the first partition conditional on another event in the second partition), would we be able to calculate the marginal probabilities for the two partitions? The answer is yes, if we commit ourselves to PME:

[PME] Keep the information entropy of your probability distribution maximal within the constraints that the evidence provides (in the synchronic case), or your cross-entropy minimal (in the diachronic case).

For Majerník’s question, PME provides us with a unique and plausible answer (see Majerník’s paper). We may also be interested in the obverse question: if the marginal probabilities of the two partitions were given, would we similarly be able to calculate the conditional probabilities? The answer is yes: given PME, Theorems 2.2.1. and 2.6.5. in *Elements of Information Theory* (see [2]) reveal that the joint probabilities are the product of the marginal probabilities. Once the joint probabilities and the marginal probabilities are available, it is trivial to calculate the conditional probabilities.

There is an older problem by Carl Wagner (see [21]), which can be cast in similar terms. If we were given some of the marginal probabilities in an updating problem as well as some logical relationships between the two partitions, would we be able to calculate the remaining marginal probabilities? This problem is best understood by example (see Wagner’s *Linguist* problem in section 3). Wagner solves it with a natural generalization of Jeffrey conditioning, which we will call Wagner conditioning. It is not based on PME, but on what we may call Jeffrey’s updating principle, or JUP for short:

[JUP] In a diachronic updating process, keep the ratio of probabilities con-

stant as long as they are unaffected by the constraints that the evidence poses.

Richard Jeffrey made this principle famous when he introduced probability kinematics (see [12]), an updating method which operates on uncertain evidence. As is the case for PME, there is a debate whether updating on evidence by rational agents is bound by JUP (for a defence see [18]; for detractors see [8]).

Our interest in this article is the relationship between PME and JUP, both of which are updating principles. Wagner contends that his natural generalization of Jeffrey conditioning, based on JUP, contradicts PME. Among formal epistemologists, there is a widespread view that, while PME is a generalization of Jeffrey conditioning, it is an inappropriate updating method in certain cases and does not enjoy the generality of Jeffrey conditioning. Wagner’s claims support this view inasmuch as Wagner conditioning is based on the relatively plausible JUP and naturally generalizes Jeffrey conditioning, but according to Wagner it contradicts PME, which gives wrong results in these cases.

This paper resists Wagner’s conclusions and shows that PME generalizes both Jeffrey conditioning and Wagner conditioning, providing a much more integrated approach to probability updating. This integrated approach also gives a coherent answer to the obverse Majernik problem posed above. A more epistemologically oriented companion paper shows that Wagner’s argument about the contradiction between JUP and PME is specious.

2 Jeffrey Conditioning

Richard Jeffrey proposes an updating method for cases in which the evidence is uncertain, generalizing standard probabilistic conditioning. I will present this method in unusual notation, anticipating using my notation to solve Wagner’s *Linguist* problem and to give a general solution for the obverse Majernik problem. Let Ω be an event space with finitely many elements and $\{\theta_j\}_{j=1,\dots,n}$ a partition of Ω . Let κ be an $m \times n$ matrix for which each column contains exactly one 1, otherwise 0. Let $P = P_{\text{prior}}$ and $\hat{P} = P_{\text{posterior}}$. Then $\{\omega_i\}_{i=1,\dots,m}$, for which

$$\omega_i = \bigcup_{j=1, \dots, n} \theta_{ij}^*, \quad (1)$$

is likewise a partition of Ω (the ω s are basically a more coarsely grained partition than the θ s). $\theta_{ij}^* = \emptyset$ if $\kappa_{ij} = 0$, $\theta_{ij}^* = \theta_j$ otherwise. Let β be the vector of prior probabilities for $\{\theta_j\}_{j=1, \dots, n}$ ($P(\theta_j) = \beta_j$) and $\hat{\beta}$ the vector of posterior probabilities ($\hat{P}(\theta_j) = \hat{\beta}_j$); likewise for α and $\hat{\alpha}$ corresponding to the prior and posterior probabilities for $\{\omega_i\}_{i=1, \dots, m}$, respectively.

A Jeffrey-type problem is when β and $\hat{\alpha}$ are given and we are looking for $\hat{\beta}$. A mathematically more concise characterization of a Jeffrey-type problem is the triple $(\kappa, \beta, \hat{\alpha})$. The solution, using Jeffrey conditioning, is

$$\hat{\beta}_j = \beta_j \sum_{i=1}^m \frac{\kappa_{ij} \hat{\alpha}_i}{\sum_{\kappa_{il}=1} \beta_l} \text{ for all } j = 1, \dots, n. \quad (2)$$

I will soon introduce an example which makes this notation more perspicuous. The notation is at first glance off-putting, but we will in the following take full advantage of it to present a generalization where the ω_i do not range over the θ_j .

3 Wagner Conditioning

Carl Wagner uses JUP (explained in more detail in [22]) to solve a problem which cannot be solved by Jeffrey conditioning. Here is the narrative (call this the *Linguist* problem):

You encounter the native of a certain foreign country and wonder whether he is a Catholic northerner (θ_1), a Catholic southerner (θ_2), a Protestant northerner (θ_3), or a Protestant southerner (θ_4). Your prior probability p over these possibilities (based, say, on population statistics and the judgment that it is reasonable to regard this individual as a random representative of his country) is given by $p(\theta_1) = 0.2, p(\theta_2) = 0.3, p(\theta_3) = 0.4$, and $p(\theta_4) = 0.1$. The individual now utters a phrase in his native tongue which, due to the aural similarity of the phrases in question, might be a traditional Catholic

piety (ω_1), an epithet uncomplimentary to Protestants (ω_2), an innocuous southern regionalism (ω_3), or a slang expression used throughout the country in question (ω_4). After reflecting on the matter you assign subjective probabilities $u(\omega_1) = 0.4, u(\omega_2) = 0.3, u(\omega_3) = 0.2$, and $u(\omega_4) = 0.1$ to these alternatives. In the light of this new evidence how should you revise p ? (See [21, 252], and [17, 197].)

Let us call a problem of this type a Wagner-type problem. It is an instance of the more general obverse Majernik problem where partitions are given with logical relationships between them as well as some marginal probabilities. Wagner-type problems seek as a solution missing marginals, while obverse Majernik problems seek the conditional probabilities as well, both of which we will eventually provide using PME.

Wagner’s solution for such problems (from now on Wagner conditioning) rests on JUP and a formal apparatus established by Arthur Dempster (see [4]), which is quite different from our notational approach. Wagner legitimately calls his solution a “natural generalization of Jeffrey conditioning” (see [21, 250]). There is, however, another natural generalization of Jeffrey conditioning, E.T. Jaynes’ principle of maximum entropy. PME does not rest on JUP, but rather claims that one should keep one’s entropy maximal within the constraints that the evidence provides (in the synchronic case) and one’s cross-entropy minimal (in the diachronic case). Some distinguish between MAXENT, the synchronic rule, and *Infomin*, the diachronic rule, but I have shown elsewhere that the two are compatible and both follow PME (see also [22]).

It turns out that PME elegantly generalizes Jeffrey conditioning and therefore absorbs JUP on the more narrow domain of problems that we can solve using Jeffrey conditioning (for a proof see [1]). Wagner’s contention is that on the wider domain of problems where we must use Wagner conditioning, JUP and PME contradict each other. We are now in the awkward position of being confronted with two plausible intuitions, JUP and PME, and it appears that we have to let one of them go. Wagner adduces other conceptual problems for PME (e.g. Bas van Fraassen’s *Judy Benjamin* problem and Abner Shimony’s Lagrange multiplier problem, see [6]) to reinforce his conclusion that PME is not a principle on which we should rely in general.

We will see that Wagner’s conclusion is incorrect. JUP and PME are compatible. Wagner’s formal apparatus, although inspiring, is unnecessary and ad hoc, as the much more integrated maximum entropy approach seamlessly

generalizes JUP. There are now two distinctive tasks at hand. One is to show how Wagner construes a contradiction between JUP and PME and where this construction is misleading. I will do this in a more epistemological companion paper, because Wagner’s mistake is more epistemological in nature than formal, attributing implausible assumptions to adherents of PME. The other more general and more formal task, which we will pursue here, is to show how PME generalizes Jeffrey conditioning and Wagner conditioning to boot.

4 A Natural Generalization of Jeffrey and Wagner Conditioning

To achieve the second task, we use the notation that we have already introduced for Jeffrey conditioning. We can characterize Wagner-type problems analogously to Jeffrey-type problems by a triple $(\kappa, \beta, \hat{\alpha})$. $\{\theta_j\}_{j=1,\dots,n}$ and $\{\omega_i\}_{i=1,\dots,m}$ now refer to independent partitions of Ω , i.e. (1) need not be true. Besides the marginal probabilities $P(\theta_j) = \beta_j, \hat{P}(\theta_j) = \hat{\beta}_j, P(\omega_i) = \alpha_i, \hat{P}(\omega_i) = \hat{\alpha}_i$, we therefore also have joint probabilities $m_{ij} = P(\omega_i \cap \theta_j)$ and $\hat{m}_{ij} = \hat{P}(\omega_i \cap \theta_j)$.

Given the specific nature of Wagner-type problems, there are a few constraints on the triple $(\kappa, \beta, \hat{\alpha})$. The last row $(m_{mj})_{j=1,\dots,n}$ is special because it represents the probability of ω_m , which is the negation of the events deemed possible after the observation. In the *Linguist* problem, for example, ω_5 is the event (initially highly likely, but impossible after the observation of the native’s utterance) that the native does not make any of the four utterances. The native may have, after all, uttered a typical Buddhist phrase, asked where the nearest bathroom was, complimented your fedora, or chosen to be silent. κ will have all 1s in the last row. Let $\hat{\kappa}_{ij} = \kappa_{ij}$ for $i = 1, \dots, m-1$ and $j = 1, \dots, n$; and $\hat{\kappa}_{mj} = 0$ for $j = 1, \dots, n$. $\hat{\kappa}$ equals κ except that its last row are all 0s, and $\hat{\alpha}_m = 0$. Otherwise the 0s are distributed over κ (and equally over $\hat{\kappa}$) so that no row and no column has all 0s, representing the logical relationships between the ω_i s and the θ_j s. We set $P(\omega_m) = x$ ($\hat{P}(\omega_m) = 0$), where x depends on your prior knowledge. Fortunately, the value of x cancels out nicely and will play no further role. For convenience, we define $\zeta = (0, \dots, 0, 1)^\top$ with $\zeta_m = 1$ and $\zeta_i = 0$ for $i \neq m$.

The best way to visualize such a problem is by providing the joint probability matrix $M = (m_{ij})$ together with the marginals α and β in the last column/row, here for example as for the *Linguist* problem with $m = 5$ and

$n = 4$,

$$\begin{bmatrix} m_{11} & m_{12} & 0 & 0 & \alpha_1 \\ m_{21} & m_{22} & 0 & 0 & \alpha_2 \\ 0 & m_{32} & 0 & m_{34} & \alpha_3 \\ m_{41} & m_{42} & m_{43} & m_{44} & \alpha_4 \\ m_{51} & m_{52} & m_{53} & m_{54} & x \\ \beta_1 & \beta_2 & \beta_3 & \beta_4 & 1.00 \end{bmatrix}. \quad (3)$$

The $m_{ij} \neq 0$ where $\kappa_{ij} = 1$. Ditto, mutatis mutandis, for $\hat{M}, \hat{\alpha}, \hat{\beta}$. To make this a little less abstract, Wagner's *Linguist* problem is characterized by the triple $(\kappa, \beta, \hat{\alpha})$,

$$\kappa = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} \text{ and } \hat{\kappa} = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad (4)$$

$$\beta = (0.2, 0.3, 0.4, 0.1)^\top \text{ and } \hat{\alpha} = (0.4, 0.3, 0.2, 0.1, 0)^\top. \quad (5)$$

Wagner's solution, based on JUP, is

$$\hat{\beta}_j = \beta_j \sum_{i=1}^{m-1} \frac{\hat{\kappa}_{ij} \hat{\alpha}_i}{\sum_{l=1}^m \hat{\kappa}_{il} \beta_l} \text{ for all } j = 1, \dots, n. \quad (6)$$

In numbers,

$$\hat{\beta}_j = (0.3, 0.6, 0.04, 0.06)^\top. \quad (7)$$

The posterior probability that the native encountered by the linguist is a northerner, for example, is 34%. Wagner's notation is completely different

and never specifies or provides the joint probabilities, but I hope the reader appreciates both the analogy to (2) underlined by this notation as well as its efficiency in delivering a correct PME solution for us (compared to Wagner's incorrect PME solution, which is to some extent misleadingly suggested by Wagner's Dempsterian setup). That (6) follows from JUP is well-documented in Wagner's article.

For the PME solution for this problem, we will not use (6) or JUP, but maximize the entropy for the joint probability matrix M and then minimize the cross-entropy between the prior probability matrix M and the posterior probability matrix \hat{M} . The PME solution, despite its seemingly different ancestry in principle, formal method, and assumptions, agrees with (6). This completes our argument.

What follows may only be accessible to PME cognoscenti, since it involves the Lagrange multiplier method (see [7, 327ff], and [10, 244]). Others may want to skip to the Conclusion. To maximize the Shannon entropy of M and minimize the Kullback-Leibler divergence between \hat{M} and M , consider the Lagrangian functions:

$$\begin{aligned} \Lambda(m_{ij}, \mu) = & \\ & \sum_{\kappa_{ij}=1} m_{ij} \log m_{ij} + \sum_{j=1}^n \mu_j \left(\beta_j - \sum_{\kappa_{kj}=1} m_{kj} \right) + \\ & \lambda_m \left(x - \sum_{j=1}^n m_{mj} \right) \end{aligned} \tag{8}$$

and

$$\begin{aligned} \hat{\Lambda}(\hat{m}_{ij}, \hat{\lambda}) = & \\ & \sum_{\hat{\kappa}_{ij}=1} \hat{m}_{ij} \log \frac{\hat{m}_{ij}}{m_{ij}} + \sum_{i=1}^m \hat{\lambda}_i \left(\hat{\alpha}_i - \sum_{\hat{\kappa}_{il}=1} \hat{m}_{il} \right). \end{aligned} \tag{9}$$

For the optimization, we set the partial derivatives to 0, which results in

$$M = rs^\top \circ \kappa \quad (10)$$

$$\hat{M} = \hat{r}s^\top \circ \hat{\kappa} \quad (11)$$

$$\beta = S\kappa^\top r \quad (12)$$

$$\hat{\alpha} = \hat{R}\kappa s \quad (13)$$

where $r_i = e^{\zeta_i \lambda_m}$, $s_j = e^{-1-\mu_j}$, $\hat{r}_i = e^{-1-\hat{\lambda}_i}$ represent factors arising from the Lagrange multiplier method. The operator \circ is the entry-wise Hadamard product in linear algebra. r, s, \hat{r} are the vectors containing the r_i, s_j, \hat{r}_i , respectively. R, S, \hat{R} are the diagonal matrices with $R_{il} = r_i \delta_{il}$, $S_{kj} = s_j \delta_{kj}$, $\hat{R}_{il} = \hat{r}_i \delta_{il}$ (δ is Kronecker delta).

Note that

$$\frac{\beta_j}{\sum_{\hat{\kappa}_{il}=1} \beta_l} = \frac{s_j}{\sum_{\hat{\kappa}_{il}=1} s_l} \text{ for all } (i, j) \in \{1, \dots, m-1\} \times \{1, \dots, n\}. \quad (14)$$

(13) implies

$$\hat{r}_i = \frac{\hat{\alpha}_i}{\sum_{\hat{\kappa}_{il}=1} s_l} \text{ for all } i = 1, \dots, m-1. \quad (15)$$

Consequently,

$$\hat{\beta}_j = s_j \sum_{i=1}^{m-1} \frac{\hat{\kappa}_{ij} \hat{\alpha}_i}{\sum_{\kappa_{il}=1} s_l} \text{ for all } j = 1, \dots, n. \quad (16)$$

(16) gives us the same solution as (6), taking into account (14). Therefore, Wagner conditioning and PME agree.

5 Conclusion

Wagner-type problems (but not obverse Majernik-type problems) can be solved using JUP and Wagner’s ad hoc method. Obverse Majernik-type problems, and therefore all Wagner-type problems, can also be solved using PME and its established and integrated formal method. What at first blush looks like serendipitous coincidence, namely that the two approaches deliver the same result, reveals that JUP is safely incorporated in PME. Not to gain information where such information gain is unwarranted and to process all the available and relevant information is the intuition at the foundation of PME. My results show that this more fundamental intuition generalizes the more specific intuition that ratios of probabilities should remain constant unless they are affected by observation or evidence. Wagner’s argument that PME conflicts with JUP is ineffective because, as my more epistemological companion paper demonstrates, it rests on assumptions that advocates of PME naturally reject.

6 References

- [1] Caticha, Ariel, and Adom Giffin. “Updating Probabilities.” In *Max-Ent 2006, the 26th International Workshop on Bayesian Inference and Maximum Entropy Methods*. 2006.
- [2] Cover, T.M., and J.A. Thomas. *Elements of Information Theory*, volume 6. Hoboken, NJ: Wiley, 2006.
- [3] Debbah, Mérouane, and Ralf Müller. “MIMO Channel Modeling and the Principle of Maximum Entropy.” *IEEE Transactions on Information Theory* 51, 5: (2005) 1667–1690.
- [4] Dempster, Arthur. “Upper and Lower Probabilities Induced by a Multivalued Mapping.” *The Annals of Mathematical Statistics* 38, 2: (1967) 325–339.
- [5] Kampé de Fériet, J., and B. Forte. “Information et probabilité.” *Comptes rendus de l’Académie des sciences A* 265: (1967) 110–114.
- [6] Friedman, Kenneth, and Abner Shimony. “Jaynes’s Maximum Entropy Prescription and Probability Theory.” *Journal of Statistical Physics* 3, 4: (1971) 381–384.

- [7] Guiaşu, Silviu. *Information Theory with Application*. New York: McGraw-Hill, 1977.
- [8] Howson, Colin, and Allan Franklin. “Bayesian Conditionalization and Probability Kinematics.” *The British Journal for the Philosophy of Science* 45, 2: (1994) 451–466.
- [9] Ingarden, R. S., and K. Urbanik. “Information Without Probability.” *Colloquium Mathematicum* 9: (1962) 131–150.
- [10] Jaynes, E.T. “Where Do We Stand on Maximum Entropy.” In *The Maximum Entropy Formalism*, edited by R.D. Levine, and M. Tribus, Cambridge, MA: MIT, 1978, 15–118.
- [11] ———. “Optimal Information Processing and Bayes’s Theorem: Comment.” *The American Statistician* 42, 4: (1988) 280–281.
- [12] Jeffrey, Richard. *The Logic of Decision*. New York, NY: Gordon and Breach, 1965.
- [13] Khinchin, Aleksandr. *Mathematical Foundations of Information Theory*. New York: Dover Publications, 1957.
- [14] Kolmogorov, Andrey. “Logical Basis for Information Theory and Probability Theory.” *IEEE Transactions on Information Theory* 14, 5: (1968) 662–664.
- [15] Majerník, Vladimír. “Marginal Probability Distribution Determined by the Maximum Entropy Method.” *Reports on Mathematical Physics* 45, 2: (2000) 171–181.
- [16] Palmieri, Francesco, and Domenico Ciuonzo. “Objective Priors from Maximum Entropy in Data Classification.” *Information Fusion* 14, 2: (2013) 186–198.
- [17] Spohn, Wolfgang. *The Laws of Belief: Ranking Theory and Its Philosophical Applications*. Oxford, 2012.
- [18] Teller, Paul. “Conditionalization and Observation.” *Synthese* 26, 2: (1973) 218–258.
- [19] Van Fraassen, Bas, R.I.G. Hughes, and Gilbert Harman. “A Problem for Relative Information Minimizers, Continued.” *The British Journal for the Philosophy of Science* 37, 4: (1986) 453–463.

- [20] Veveakis, Emmanuil, and Klaus Regenauer-Lieb. “Review of Extremum Postulates.” *Current Opinion in Chemical Engineering* 7: (2015) 40–46.
- [21] Wagner, Carl. “Generalized Probability Kinematics.” *Erkenntnis* 36, 2: (1992) 245–257.
- [22] ———. “Probability Kinematics and Commutativity.” *Philosophy of Science* 69, 2: (2002) 266–278.
- [23] Zellner, Arnold. “Optimal Information Processing and Bayes’s Theorem.” *The American Statistician* 42, 4: (1988) 278–280.