

PHILOSOPHIES OF PROBABILITY: OBJECTIVE BAYESIANISM AND ITS CHALLENGES

Jon Williamson

Draft of May 19, 2005.

In A. Irvine (ed.): Handbook of the Philosophy of Mathematics,
Volume 9 of the Handbook of the Philosophy of Science, Elsevier.

ABSTRACT

This chapter presents an overview of the major interpretations of probability followed by an outline of the objective Bayesian interpretation and a discussion of the key challenges it faces.

CONTENTS

§1	INTRODUCTION	3
I	FRAMEWORKS FOR PROBABILITY	3
§2	VARIABLES	3
§3	EVENTS	4
§4	SENTENCES	5
II	INTERPRETATIONS OF PROBABILITY	6
§5	INTERPRETATIONS AND DISTINCTIONS	6
§6	FREQUENCY	7
§7	PROPENSITY	8
§8	CHANCE	9
§9	BAYESIANISM	10
§10	CHANCE AS ULTIMATE BELIEF	11
§11	APPLYING PROBABILITY	12

III	OBJECTIVE BAYESIANISM	12
§12	SUBJECTIVE AND OBJECTIVE BAYESIANISM	12
§13	OBJECTIVE BAYESIANISM OUTLINED	13
§14	CHALLENGES	15
§15	MOTIVATION	15
§16	LANGUAGE DEPENDENCE	17
§17	COMPUTATION	19
§18	QUALITATIVE KNOWLEDGE	21
§19	INFINITE DOMAINS	23
§20	FULLY OBJECTIVE PROBABILITY	25
§21	PROBABILITY LOGIC	28
§22	CONCLUSION	33

§1

INTRODUCTION

The concept of probability motivates two key questions.

First, how is probability to be defined? Probability was axiomatised in the first half of the 20th century;¹ this axiomatisation has by now become well entrenched, and in fact the only leeway these days is with regard to the type of domain on which probability functions are defined. Part I introduces three types of domain: variables (§2), events (§3), and sentences (§4).

Second, how is probability to be applied? In order to know how probability can be applied we need to know what probability means: how probabilities can be measured and how probabilistic predictions say something about the world. Part II discusses the predominant interpretations of probability: the frequency (§6), propensity (§7), chance (§§8, 10), and Bayesian interpretations (§9).

In Part III, we shall focus on one interpretation of probability, objective Bayesianism, and look more closely at some of the challenges that this interpretation faces.

PART I

FRAMEWORKS FOR PROBABILITY

§2

VARIABLES

The most basic framework for probability involves defining a probability function relative to a finite set V of variables, each of which takes finitely many possible values. I shall write $v@V$ to indicate that v is an assignment of values to V .

A *probability function* on V is a function p that maps each assignment $v@V$ to a non-negative real number and which satisfies *additivity*:

$$\sum_{v@V} p(v) = 1.$$

This restriction forces each probability $p(v)$ to lie in the unit interval $[0, 1]$.

The *marginal probability function* on $U \subseteq V$ induced by probability function p on V is a probability function q on U which satisfies:

$$q(u) = \sum_{v@V, v \sim u} p(v)$$

for each $u@U$, and where $v \sim u$ means that v is consistent with u , i.e. u and v assigns the same values to $U \cap V = U$. The marginal probability function q on U is uniquely determined by p . Marginal probability functions are usually thought of as extensions of p and denoted by the same letter p . Thus p can be construed as a function that maps each $u@U \subseteq V$ to a non-negative real number. p can be further extended to assign numbers to conjunctions tu of assignments where $t@T \subseteq V, u@U \subseteq V$: if $t \sim u$ then tu is an assignment to

¹(Kolmogorov, 1933)

$T \cup U$ and $p(tu)$ is the marginal probability awarded to $tu @ (T \cup U)$; if $t \not\sim u$ then $p(tu)$ is taken to be 0.

A *conditional probability function* induced by p is a function r from pairs of assignments of subsets of V to non-negative real numbers which satisfies (for each $t @ T \subseteq V, u @ U \subseteq V$):

$$r(t|u)p(u) = p(tu),$$

$$\sum_{t @ T} r(t|u) = 1,$$

Note that $r(t|u)$ is not uniquely determined by p when $p(u) = 0$. If $p(u) \neq 0$ and the first condition holds, then the second condition, $\sum_{t @ T} r(t|u) = 1$, also holds. Again, r is often thought of as an extension of p and is usually denoted by the same letter p .

Consider an example. Take $V = \{A, B\}$ is a domain of variables, where A signifies *age of vehicle* taking possible values *less than 3 years, 3-10 years* and *greater than 10 years*, and B signifies *breakdown in the last year* taking possible values *yes* and *no*. An assignment $b @ B$ is of the form $B = \textit{yes}$ or $B = \textit{no}$. The assignments $a @ A$ are most naturally written $A < 3, 3 \leq A \leq 10$ and $A > 10$. According to the above definition a probability function p on V assigns a non-negative real number to each assignment of the form ab where $a @ A$ and $b @ B$, and these numbers must sum to 1. For instance,

$$p(A < 3 \cdot B = \textit{yes}) = 0.05$$

$$p(A < 3 \cdot B = \textit{no}) = 0.1$$

$$p(3 \leq A \leq 10 \cdot B = \textit{yes}) = 0.2$$

$$p(3 \leq A \leq 10 \cdot B = \textit{no}) = 0.2$$

$$p(A > 10 \cdot B = \textit{yes}) = 0.35$$

$$p(A > 10 \cdot B = \textit{no}) = 0.1.$$

This function p can be extended to assignments of subsets of V , yielding $p(A > 10) = p(A > 10 \cdot B = \textit{yes}) + p(A > 10 \cdot B = \textit{no}) = 0.35 + 0.1 = 0.45$ for example, and to conjunctions of assignments in which case inconsistent assignments are awarded probability 0, e.g. $p(B = \textit{yes} \cdot B = \textit{no}) = 0$. The function p can then be extended to yield conditional probabilities and in this example the probability of a breakdown conditional on age greater than 10 years, $p(B = \textit{yes} | A > 10)$, is $p(A > 10 \cdot B = \textit{yes}) / p(A > 10) = 0.35 / 0.45 \approx 0.78$.

§3

EVENTS

While the definition of probability over assignments to variables is straightforward, simplicity is gained at the expense of generality. By moving from variables to abstract events we can capture generality. The main definition proceeds as follows.²

²Billingsley (1979) provides a good introduction to the theory behind this approach.

Abstract *events* are construed as subsets of an *outcome space* Ω , which represents the possible outcomes of an experiment or observation. For example, if the age of a vehicle were observed, the outcome space might be $\Omega = \{0, 1, 2, \dots\}$, and $\{0, 1, 2\} \subseteq \Omega$ represents the event that the vehicle's age is less than three years.

An *event space* \mathcal{F} is a set of subsets of Ω . \mathcal{F} is a *field* if it contains Ω and is closed under the formation of complements and finite unions; it is a σ -*field* if it is also closed under the formation of countable unions.

A *probability function* is a function p from a field \mathcal{F} to the non-negative real numbers that satisfies *countable additivity*:

- if $E_1, E_2, \dots \in \mathcal{F}$ partition Ω (i.e. $E_i \cap E_j = \emptyset$ for $i \neq j$ and $\bigcup_{i=1}^{\infty} E_i = \Omega$) then $\sum_{i=1}^{\infty} p(E_i) = 1$.

In particular, $p(\Omega) = 1$. The triple (Ω, \mathcal{F}, p) is called a *probability space*.

The variable framework is captured by letting Ω contain all assignments to V and taking \mathcal{F} to be the set of all subsets of Ω , which corresponds to the set of disjunctions of assignments to V . Given variable $A \in V$, the function that maps $v \in V$ to the value that v assigns to A is called a *simple random variable* in the event framework.

§4

SENTENCES

Logicians tend to define probability over logical languages.³ The simplest such framework is based around the propositional calculus, as follows.

A *propositional variable* is a variable which takes two possible values, *true* or *false*. A set L of propositional variables is called a *propositional language*. The sentences S_L of L include the propositional variables, together with the *negation* $\neg\theta$ of each sentence $\theta \in S_L$ (which is true iff θ is false) and each *implication* of the form $\theta \rightarrow \phi$ for $\theta, \phi \in S_L$ (which is true iff θ is false or both θ and ϕ are true). The *conjunction* $\theta \wedge \phi$ is defined to be $\neg(\theta \rightarrow \neg\phi)$ and is true iff both θ and ϕ are true; the *disjunction* $\theta \vee \phi$ is defined to be $\neg\theta \rightarrow \phi$ and is true iff either θ or ϕ are true. An assignment l of values to L *models* sentence θ , written $l \models \theta$, if θ is true under l . A sentence θ is a *tautology*, written $\models \theta$, if it is true whatever the values of the propositional variables in θ , i.e. if each assignment to L models θ .

A *probability function* is then a function p from a set S_L of sentences to the non-negative real numbers that satisfies *additivity*:

- if $\theta_1, \dots, \theta_n \in S_L$ satisfy $\models \neg(\theta_i \wedge \theta_j)$ for $i \neq j$ and $\models \theta_1 \vee \dots \vee \theta_n$ then $\sum_{i=1}^n p(\theta_i) = 1$.

If the language L is finite then the sentence framework can be mapped to the variable framework. $V = L$ is a finite set of variables each of which takes finitely many values. A sentence $\theta \in S_V$ can be identified with the set of assignments v of values to V which model θ . p thus maps sets of assignments, and in particular individual assignments, to real numbers. p is additive because of additivity on sentences. Hence p induces a probability function over assignments to V .

³See Paris (1994).

The sentence framework can also be mapped to the event framework. Let Ω contain all assignments to L , and let \mathcal{F} be the field of sets of the form $\{l : l \models \theta\}$ for $\theta \in S_L$.⁴ By defining $p(\{l : l \models \theta\}) = p(\theta)$ we get a probability function.⁵

PART II INTERPRETATIONS OF PROBABILITY

§5

INTERPRETATIONS AND DISTINCTIONS

The definitions of probability given in Part I are purely formal. In order to apply the formal concept of probability we need to know how probability is to be interpreted. The standard interpretations of probability will be presented in the next few sections.⁶ These interpretations can be categorised according to the stances they take on three key distinctions:

SINGLE-CASE / REPEATABLE A variable is *single-case* (or *token-level*) if it can only be assigned a value once. It is *repeatable* (or *repeatably instantiatable* or *type-level*) if it can be assigned values more than once. For example, variable A standing for *age of car with registration AB01 CDE on January 1st 2005* is single-case because it can only ever take one value (assuming the car in question exists). If however A stands for *age of vehicles selected at random in London in 2005* then A is repeatable: it gets reassigned a value each time a new vehicle is selected.⁷

MENTAL / PHYSICAL Probabilities are *mental* (or *epistemological*⁸ or *personalist*) if they are interpreted as features of an agent's mental state, otherwise they are *physical* (or *aleatory*⁹).

SUBJECTIVE / OBJECTIVE Probabilities are *subjective* (or *agent-relative*) if two agents with the same background knowledge can disagree as to a probability value and yet neither of them be wrong. Otherwise they are *objective*.¹⁰

There are four main interpretations of probability: the frequency theory (§6), the propensity theory (§7), chance (§8) and Bayesianism (§9).

⁴These sets are called *cylinder sets* when L is infinite—see Billingsley (1979, p. 27).

⁵This depends on the fact that every probability function on the field of cylinders which is *finitely additive* (i.e. which satisfies $\sum_{i=1}^n p(E_i) = 1$ for partition E_1, \dots, E_n of Ω) is also countably additive. See Billingsley (1979, Theorem 2.3).

⁶For a more detailed exposition of the interpretations see Gillies (2000).

⁷'Single-case variable' is clearly an oxymoron because the value of a single-case variable does not vary. The value of a single-case variable may not be known, however, and one can still think of the variable as taking a range of possible values.

⁸(Gillies, 2000)

⁹(Hacking, 1975)

¹⁰Warning: some authors, such as Popper (1983, §3.3) and Gillies (2000, p. 20), use the term 'objective' for what I call 'physical'. However their terminology has the awkward consequence that the interpretation of probability commonly known as 'objective Bayesianism' (described in Part III) does not get classed as 'objective'.

§6

FREQUENCY

The *frequency* interpretation of probability was propounded by Venn¹¹ and Reichenbach¹² and developed in detail by Richard von Mises.¹³ Von Mises' theory can be formulated in our framework as follows. Given a set V of repeatable variables one can repeatedly determine the values of the variables in V and write down the observations as assignments to V . For example, one could repeatedly select cars and determine their age and whether they broke down in the last year, writing down $A < 3 \cdot B = no$, $A < 3 \cdot B = yes$, $A > 10 \cdot B = yes$, and so on. Under the assumption that this process of measurement can be repeated ad infinitum, we generate an infinite sequence of assignments $\mathcal{V} = (v_1, v_2, v_3, \dots)$ called a *collective*.

Let $|v|_{\mathcal{V}}^n$ be the number of times assignment v occurs in the first n places of \mathcal{V} , and let $freq_{\mathcal{V}}^n(v)$ be the frequency of v in the first n places of \mathcal{V} , i.e.

$$freq_{\mathcal{V}}^n(v) = \frac{|v|_{\mathcal{V}}^n}{n}.$$

Von Mises noted two things. First, these frequencies tend to stabilise as the number n of observations increases. Von Mises hypothesised that

AXIOM OF CONVERGENCE $freq_{\mathcal{V}}^n(v)$ tends to a fixed limit as $n \rightarrow \infty$, denoted by $freq_{\mathcal{V}}(v)$.

Second, gambling systems tend to be ineffective. A gambling system can be thought of as function for selecting places in the sequence of observations on which to bet, on the basis of past observations. Thus a *place selection* is a function $f(v_1, \dots, v_n) \in 0, 1$, such that if $f(v_1, \dots, v_n) = 0$ then no bet is to be placed on the $n + 1$ -st observation and if $f(v_1, \dots, v_n) = 1$ then a bet is to be placed on the $n + 1$ -st observation. So betting according to a place selection gives rise to a sub-collective \mathcal{V}_f of \mathcal{V} consisting of the places of \mathcal{V} on which bets are placed. In practice we can only use a place selection function if it is simple enough for us to compute its values: if we cannot decide whether $f(v_1, \dots, v_n)$ is 0 or 1 then it is of no use as a gambling system. According to Church's thesis a function is computable if it belongs to the class of functions known as *recursive functions*.¹⁴ Accordingly we define a *gambling system* to be a recursive place selection. A gambling system is said to be effective if we are able to make money in the long run when we place bets according to the gambling system. Assuming that stakes are set according to frequencies of \mathcal{V} , a gambling system f can only be effective if the frequencies of \mathcal{V}_f differ to those of \mathcal{V} : if $freq_{\mathcal{V}_f}(v) > freq_{\mathcal{V}}(v)$ then betting *on* v will be profitable in the long run; if $freq_{\mathcal{V}_f}(v) < freq_{\mathcal{V}}(v)$ then betting *against* v will be profitable. We can then explicate von Mises' second observation as follows:

AXIOM OF RANDOMNESS Gambling systems are ineffective: if \mathcal{V}_f is determined by a recursive place selection f , then for each v , $freq_{\mathcal{V}_f}(v) = freq_{\mathcal{V}}(v)$.

¹¹(Venn, 1866)

¹²(Reichenbach, 1935)

¹³(von Mises, 1928, 1964)

¹⁴(Church, 1936)

Given a collective \mathcal{V} we can then define—following von Mises—the probability of v to be the frequency of v in \mathcal{V} :

$$p(v) =_{df} \text{freq}_{\mathcal{V}}(v).$$

Clearly $\text{freq}_{\mathcal{V}}(v) \geq 0$. Moreover $\sum_{v \in V} |v|_{\mathcal{V}}^n = n$ so $\sum_{v \in V} \text{freq}_{\mathcal{V}}^n(v) = 1$ and, taking limits, $\sum_{v \in V} \text{freq}_{\mathcal{V}}(v) = 1$. Thus p is indeed a well-defined probability function.

Suppose we have a statement involving probability function p on V . If we also have a collective \mathcal{V} on V then we can interpret the statement to be saying something about the frequencies of \mathcal{V} , and as being true or false according to whether the corresponding statement about frequencies is true or false respectively. This is the frequency interpretation of probability. The variables in question are repeatable, not single-case, and the interpretation is physical, relative to a collective of potential observations, not to the mental state of an agent. The interpretation is objective, not subjective, in the sense that once the collective is fixed then so too are the probabilities: if two agents disagree as to what the probabilities are, then at most one of the agents is right.

§7

PROPENSITY

Karl Popper initially adopted a version of von Mises' frequency interpretation,¹⁵ but later, with the ultimate goal of formulating an interpretation of probability applicable to single-case variables, developed what is called the *propensity* interpretation of probability.¹⁶ The propensity theory can be thought of as the frequency theory together with the following law:¹⁷

AXIOM OF INDEPENDENCE If collectives \mathcal{V}_1 and \mathcal{V}_2 on V are generated by the same repeatable experiment (or repeatable conditions) then for all assignments v to V , $\text{freq}_{\mathcal{V}_1}(v) = \text{freq}_{\mathcal{V}_2}(v)$.

In other words frequency, and hence probability, attaches to repeatable experiment rather than a collective, in the sense that frequencies do not vary with collectives generated by the same repeatable experiment. The repeatable experiment is said to have a propensity for generating the corresponding frequency distribution.

In fact, despite Popper's intentions, the propensity theory interprets probability defined over repeatable variables, not single-case variables. If for example V consists of repeatable variables A and B , where A stands for *age of vehicles selected at random in London in 2005* and B stands for *breakdown in the last year of vehicles selected at random in London in 2005*, then V determines a repeatable experiment, namely the selection of vehicles at random in London in 2005, and thus there is a natural propensity interpretation. Suppose on the

¹⁵(Popper, 1934, Chapter VIII)

¹⁶(Popper, 1959; Popper, 1983, Part II)

¹⁷Popper (1983, pp. 290 and 355). It is important to stress that the axioms of this section and the last had a different status for Popper than they did for von Mises. Von Mises used the frequency axioms as part of an operationalist definition of probability, but Popper was not an operationalist. See Gillies (2000, Chapter 7) on this point. Gillies also argues in favour of a propensity interpretation.

other hand that V contains single-case variables A and B , standing for *age of car with registration AB01 CDE on January 1st 2005* and *breakdown in last year of car with registration AB01 CDE on January 1st 2005*. Then V defines an experiment, namely the selection of car AB01 CDE on January 1st 2005, but this experiment is not repeatable and does not generate a collective—it is a single case. The car in question might be selected by several different repeatable experiments, but these repeatable experiments need not yield the same frequency for an assignment v , and thus the probability of v is not determined by V . (This is known as the *reference class problem*: we do not know from the specification of the single case how to uniquely determine a repeatable experiment which will fix probabilities.) In sum the propensity theory is, like the frequency theory, an objective, physical interpretation of probability over repeatable variables.

§8

CHANCE

The question remains as to whether one can develop a viable objective interpretation of probability over single-case variables—such a concept of probability is often called *chance*.¹⁸ We saw that frequencies are defined relative to a collective and propensities are defined relative to a repeatable experiment; however a single-case variable does not determine a unique collective or repeatable experiment and so neither approach allows us to attach probabilities directly to single-case variables. What then does fix the chances of a single-case variable? The view finally adopted by Popper was that the ‘whole physical situation’ determines probabilities.¹⁹ The physical situation might be thought of as ‘the complete situation of the universe (or the light-cone) at the time’,²⁰ the complete history of the world up till the time in question,²¹ or ‘a complete set of (nominally and/or causally) relevant conditions . . . which happens to be instantiated in that world at that time’.²² Thus the chance, on January 1st 2005, of car with registration AB01 CDE breaking down in the subsequent year, is fixed by the state of the universe at that date, or its entire history up till that date, or all the relevant conditions instantiated at that date. However the chance-fixing ‘complete situation’ is delineated, these three approaches associate a unique chance-fixer with a given single-case variable. (In contrast, the frequency / propensity theories do not associate a unique collective / repeatable experiment with a given single-case variable.) Hence we can interpret the probability of an assignment to the single-case variable as the chance of the assignment holding, as determined by its chance-fixer.

Further explanation is required as to how one can measure probabilities under the chance interpretation. Popper’s line is this: if the chance-fixer is a set of relevant conditions, and these conditions are repeatable then the conditions determine a propensity and that can be used to measure the chance.²³ Thus

¹⁸Note that some authors use ‘propensity’ to cover a physical chance interpretation as well as the propensity interpretation discussed above.

¹⁹(Popper, 1990, p. 17)

²⁰(Miller, 1994, p. 186)

²¹Lewis (1980, p. 99). See §§10, 20.

²²(Fetzer, 1982, p. 195)

²³(Popper, 1990, p. 17)

if the set of conditions relevant to car AB01 CDE breaking down that hold on January 1st 2005 also hold for other cars at other times, then the chance of AB01 CDE breaking down in the next year can be equated with the frequency with which cars satisfying the same set of conditions break down in the subsequent year. The difficulty with this view is that it is hard to determine all the chance-fixing relevant conditions, and there is no guarantee that enough individuals will satisfy this set of conditions for the corresponding frequency to be estimable.

§9

BAYESIANISM

The *Bayesian* interpretation of probability also deals with probability functions defined over single-case variables. But in this case the interpretation is mental rather than physical: probabilities are interpreted as an agent's rational degrees of belief.²⁴ Thus for an agent, $p(B = \textit{yes}) = q$ if and only if the agent believes that $B = \textit{yes}$ to degree q and this ascription of degree of belief is rational in the sense outlined below. An agent's degrees of belief are construed as a guide to her actions: she believes $B = \textit{yes}$ to degree q if and only if she is prepared to place a bet of qS on $B = \textit{yes}$, with return S if $B = \textit{yes}$ turns out to be true. Here S is an unknown stake, which may be positive or negative, and q is called a *betting quotient*. An agent's *belief function* is the function that maps an assignment to the agent's degree of belief in that assignment.

An agent's betting quotients are called *coherent* if one cannot choose stakes for her bets that force her to lose money whatever happens. (Such a set of stakes is called a *Dutch book*.) It is not hard to see that a coherent belief function is a probability function. First $q \geq 0$, for otherwise one can set S to be negative and the agent will lose whatever happens: she will lose $qS > 0$ if the assignment on which she is betting turns out to be false and will lose $(q - 1)S > 0$ if it turns out to be true. Moreover $\sum_{v \in V} q_v = 1$, where q_v is the betting quotient on assignment v , for otherwise if $\sum_v q_v > 1$ we can set each $S_v = S > 0$ and the agent will lose $(\sum_v q_v - 1)S > 0$ (since exactly one of the v will turn out true), and if $\sum_v q_v < 1$ we can set each $S_v = S < 0$ to ensure positive loss.

Coherence is taken to be a necessary condition for rationality. For an agent's degrees of belief to be rational they must be coherent, and hence they must be probabilities. *Subjective Bayesianism* is the view that coherence is also sufficient for rationality, so that an agent's belief function is rational if and only if it is a probability function. This interpretation of probability is subjective because it depends on the agent as to whether $p(v) = q$. Different agents can choose different probabilities for v and their belief functions will be equally rational. *Objective Bayesianism*, discussed in detail in Part III, imposes further rationality constraints on degrees of belief—not just coherence. The aim of objective Bayesianism is to constrain degree of belief in such a way that only one value for $p(v)$ will be deemed rational on the basis of an agent's background knowledge. Thus objective Bayesian probability varies as background knowledge varies but two agents with the same background knowledge must adopt the same probabilities as their rational degrees of belief.

²⁴This interpretation was developed in Ramsey (1926) and de Finetti (1937). See Howson and Urbach (1989) and Earman (1992) for recent expositions.

Note that many Bayesians claim that an agent should update her degrees of belief by *Bayesian conditionalisation*: her new degrees of belief should be her old degrees of belief conditional on new knowledge, $p_{t+1}(v) = p_t(v|u)$ where u represents the knowledge that the agent has learned between time t and time $t+1$. In cases where $p_t(v|u)$ is harder to quantify than $p_t(u|v)$ and $p_t(v)$ this conditional probability may be calculated using *Bayes' theorem*: $p(v|u) = p(u|v)p(v)/p(u)$, which holds for any probability function p . 'Bayesianism' is variously used to refer to the Bayesian interpretation of probability, the endorsement of Bayesian conditionalisation or the use of Bayes' theorem.

§10

CHANCE AS ULTIMATE BELIEF

The question still remains as to whether one can develop a viable notion of chance, i.e. an objective single-case interpretation of probability. While the Bayesian interpretations are single-case, they either define probability relative to the whimsy of an agent (subjective Bayesianism) or relative to an agent's background knowledge (objective Bayesianism). Is there a probability of my car breaking down in the next year, where this probability does not depend on me or my knowledge?

Bayesians typically have two ways of tackling this question.

Subjective Bayesians tend to argue that although degrees of belief may initially vary widely from agent to agent, if agents update their degrees of belief by Bayesian conditionalisation then their degrees of belief will converge in the long run: chances are these long run degrees of belief. Bruno de Finetti developed such an argument to explain the apparent existence of physical probabilities.²⁵ He showed that prior degrees of beliefs converge to frequencies under the assumption of *exchangeability*: given an infinite sequence of single-case variables A_1, A_2, \dots which take the same possible values, an agent's degrees of belief are *exchangeable* if the degree of belief $p(v)$ she gives to assignment v to a finite subset of variables depends only on the values in v and not the variables in v —for example $p(a_1^1 a_2^0 a_3^1) = p(a_3^0 a_4^1 a_5^1)$ since both assignments assign two 1s and one 0. Suppose the actual observed assignments are a_1, a_2, \dots and let \mathcal{V} be the collective of such values (which can be thought of as arising from a single repeatable variable A). De Finetti showed that $p(a_n | a_1 \dots a_{n-1}) \rightarrow \text{freq}_{\mathcal{V}}(a)$ as $n \rightarrow \infty$, where a is the assignment to A of the value that occurs in a_n . The chance of a_n is then identified with $\text{freq}_{\mathcal{V}}(a)$. The trouble with de Finetti's account is that since degrees of belief are subjective there is no reason to suppose exchangeability holds. Moreover, a single-case variable A_n can occur in several sequences of variables, each with a different frequency distribution (the reference class problem again), in which case the chance distribution of A_n is ill-defined. Haim Gaifman and Marc Snir took a slightly different approach, showing that as long as agents give probability 0 to the same assignments and the evidence that they observe is unrestricted, then their degrees of belief must converge.²⁶ Again, the problem here is that there is no reason to suppose that agents will give probability 0 to the same assignments. One might try to provide such

²⁵(de Finetti, 1937; Gillies, 2000, pp. 69–83)

²⁶(Gaifman and Snir, 1982, §2)

a guarantee by bolstering subjective Bayesianism with a rationality constraint that says that agents must be *undogmatic*, i.e. they must only give probability 0 to logically impossible assignments. But this is not a feasible strategy in general, since this constraint is inconsistent with the constraint that degrees of belief be probabilities: in the more general event or sentence frameworks the laws of probability force some logical possibilities to be given probability 0.²⁷

Objective Bayesians have another recourse open to them: objective Bayesian probability is fixed by an agent's background knowledge, and one can argue that chances are those degrees of belief fixed by some suitable all-encompassing background knowledge. Thus the problem of producing a well-defined notion of chance is reducible to that of developing an objective Bayesian interpretation of probability. I shall call this the *ultimate belief* notion of chance to distinguish it from physical notions such as Popper's (§8), and discuss this approach in §20.

§11

APPLYING PROBABILITY

In sum, there are four key interpretations of probability: frequency and propensity interpret probability over repeatable variables while chance and Bayesianism deal with single-case variables; frequency and propensity are physical interpretations while Bayesianism is mental and chance can be either mental or physical; all the interpretations are objective apart from Bayesianism which can be subjective or objective.

Having chosen an interpretation of probability, one can use the probability calculus to draw conclusions about the world. Typically, having made an observation $u@U \subseteq V$, one determines the conditional probability $p(t|u)$ to tell us something about $t@T \subseteq (V \setminus U)$: a frequency, propensity, chance or degree of belief.

PART III

OBJECTIVE BAYESIANISM

§12

SUBJECTIVE AND OBJECTIVE BAYESIANISM

In Part II we saw that probabilities can either be interpreted physically—as frequencies, propensities or physical chances—or they can be interpreted mentally, with Bayesians arguing that an agent's degrees of belief ought to satisfy the axioms of probability. Many Bayesians are strict subjectivists, holding that there are no rational constraints on degrees of belief other than the requirement that they be probabilities.²⁸ Thus subjective Bayesians maintain that one may give probability 0—or indeed any value between 0 and 1—to a coin toss yielding heads, even if one knows that the coin is symmetrical and has yielded heads in roughly half of all its previous tosses. The chief criticism of strict subjectivism is that practical applications of probability tend to demand objectivity; in science

²⁷See Gaifman and Snir (1982, Theorem 3.7), for example.

²⁸(de Finetti, 1937)

some beliefs are considered more rational than others on the basis of available evidence. This motivates an alternative position, objective Bayesianism, which posits further constraints on degrees of belief, and which would only deem the agent to be rational in this case if she gave a probability of a half to the toss yielding heads.²⁹

Objective Bayesianism holds that the probability of u is the degree to which an agent ought to believe u and that this degree is objectively determined by the agent's background knowledge. Versions of this view were put forward by Jakob Bernoulli,³⁰ Laplace³¹ and Keynes.³² More recently Jaynes claimed that an agent's probabilities ought to satisfy constraints imposed by background knowledge but otherwise ought to be as non-committal as possible. Moreover, Jaynes argued, this principle could be explicated using Shannon's information theory:³³ the agent's probability function should be that probability function, from all those that satisfy constraints imposed by background knowledge, that maximises entropy.³⁴ This has become known as *the maximum entropy principle* and has been taken to be the foundation of the objective Bayesian interpretation of probability by its proponents.³⁵

In the next section, I shall sketch my own version of objective Bayesianism. This version is discussed in detail in chapter 4 of Williamson (2005a). In subsequent sections we shall examine a range of important challenges that face the objective Bayesian interpretation of probability.

§13

OBJECTIVE BAYESIANISM OUTLINED

While Bayesianism requires that degrees of belief respect the axioms of probability, objective Bayesianism imposes two further norms:

EMPIRICAL An agent's knowledge of the world should constrain her degrees of belief. Thus if one knows that a coin is symmetrical and has yielded heads roughly half the time, then one's degree of belief that it will yield heads on the next throw should be roughly $\frac{1}{2}$.

LOGICAL An agent's degrees of belief should also be fixed by her lack of knowledge of the world. If the agent knows nothing about an experiment except that it has two possible outcomes, then she should award degree of belief $\frac{1}{2}$ to each outcome.

Jakob Bernoulli pointed out that where they conflict, the empirical norm should override the logical norm:

three ships set sail from port; after some time it is announced that one of them suffered shipwreck; which one is guessed to be the one that was destroyed? If I considered merely the number of ships, I

²⁹(Jaynes, 1988)

³⁰(Bernoulli, 1713)

³¹(Laplace, 1814)

³²(Keynes, 1921)

³³(Shannon, 1948)

³⁴(Jaynes, 1957)

³⁵(Rosenkrantz, 1977; Jaynes, 2003)

would conclude that the misfortune could have happened to each of them with equal chance; but because I remember that one of them had been eaten away by rot and old age more than the others, had been badly equipped with masts and sails, and had been commanded by a new and inexperienced captain, I consider that this ship, more probably than the others, was the one to perish.³⁶

One can prioritise the empirical norm over the logical norm by insisting that:

EMPIRICAL An agent's degrees of belief, represented by probability function p_β , should satisfy any constraints imposed by her background knowledge β .

LOGICAL The agent's belief function p_β should otherwise be as non-committal as possible.

The empirical norm can be explicated as follows. Background knowledge β might contain a number of considerations that bear on a degree of belief: the symmetry of a penny might incline one to degree of belief $\frac{1}{2}$ in heads, past performance (say 47 heads in a hundred past tosses) may incline one to degree of belief 0.47, the mint may report an estimate of the frequency of heads on its pennies to be 0.45, and so on. These considerations may be thought of as conflicting reports as to the probability of heads. Intuitively, any individual report, say 0.47, is compatible with the evidence, and indeed intermediary degrees of belief such as 0.48 seem reasonable. On the other hand, a degree of belief that falls outside the range of reports, say 0.9, does not seem warranted by the evidence. Thus background knowledge constrains degree of belief to lie in the smallest closed interval that contains all the reports.

As mentioned in §12, the logical norm is explicated using the maximum entropy principle: entropy is a measure of the lack of commitment of a probability function, so p_b should be the probability function, out of all those that satisfy constraints imposed by β , that has maximum entropy. Justifications of the maximum entropy principle are well known—see Jaynes (2003), Paris (1994) or Paris and Vencovská (2001) for example.

We can thus put the two norms on a more formal footing. Given a domain V of finitely many variables, each of which takes finitely many values, an agent with background knowledge β should adopt as her belief function the probability function p_β on V determined as follows:

EMPIRICAL p_β should satisfy any constraints imposed by her background knowledge β : p_β should lie in the smallest closed convex set \mathbb{P}_β of probability functions containing those probability functions that are compatible with the reports in β .³⁷

LOGICAL p_β should otherwise be as non-committal as possible: it should be a member of \mathbb{P}_β that maximises entropy $H(p) = -\sum_{v \in V} p(v) \log p(v)$.

It turns out that there is a unique entropy maximiser on a closed convex set of probability functions: the degrees of belief p_β that an agent should adopt are uniquely determined by her background knowledge β . Thus there is no room for subjective choice of degrees of belief.

³⁶(Bernoulli, 1713, §IV.II)

³⁷See Williamson (2005a, §5.3) for more detailed discussion of this norm.

§14

CHALLENGES

Objective Bayesianism has not been widely accepted, however, largely because there are a number of perceived problems with the interpretation. Several of these problems have in fact already been resolved, but other challenges remain. In the remainder of the chapter we shall explore the key challenges and assess the prospects of objective Bayesianism.

In §15 we shall see that one challenge is to motivate the adoption of a logical norm. Objective Bayesianism has also been criticised for being language dependent (§16) and for being impractical from a computational point of view (§17). Handling qualitative background knowledge poses a significant challenge (§18), as does extending objective Bayesianism to infinite event or sentence frameworks (§19). The question of whether objective Bayesianism can be used to provide an interpretation of objective chance is explored in §20, while §21 considers the application of objective Bayesianism to providing semantics for probability logic.

Jaynes points out that the maximum entropy principle is a powerful tool but warns:

Of course, it is as true in probability theory as in carpentry that introduction of more powerful tools brings with it the obligation to exercise a higher level of understanding and judgement in using them. If you give a carpenter a fancy new power tool, he *may* use it to turn out more precise work in greater quantity; or he may just cut off his thumb with it. It depends on the carpenter.³⁸

§15

MOTIVATION

The first key question concerns the motivation behind objective Bayesianism. Recall that in §12 objective Bayesianism was motivated by the need for objective probabilities in science. Many Bayesians accept this desideratum and indeed accept the empirical norm (so that degrees of belief are constrained by knowledge of frequencies, symmetries etc.) but do not go as far as admitting a logical norm. The ensuing position, according to which degrees of belief reflect background knowledge but need not be maximally non-committal, is sometimes called *empirically-based subjective probability*. It yields degrees of belief that are more objective (i.e. more highly constrained) than those of strictly subjective Bayesianism, yet not as objective as those of objective Bayesianism—there is generally still some room for subjective choice of degrees of belief. The key question is thus: what grounds are there for going beyond empirically-based subjective probability and adopting objective Bayesianism?

Current justifications of the logical norm fail to address this question. Jaynes' original justification of the maximum entropy principle ran like this: *given that degrees of belief ought to be maximally non-committal*, Shannon's information theory shows us that they are entropy-maximising probabilities.³⁹ This type of

³⁸(Jaynes, 1979, pp. 40–41 of the original 1978 lecture)

³⁹(Jaynes, 1957)

justification assumes from the outset that some kind of logical norm is desired. On the other hand, axiomatic derivations of the maximum entropy principle take the following form: *given that we need a procedure for determining degrees of belief from background knowledge*, and given various desiderata that such a procedure should satisfy, that procedure must be entropy maximisation.⁴⁰ This type of justification takes objectivity of rational degrees of belief for granted. Thus the challenge is to augment current justifications, perhaps by motivating non-committal degrees of belief or by motivating the strong objectivity of objective Bayesianism as opposed to the partial objectivity yielded by empirically-based subjective probability.

One possible approach is to argue that empirically-based subjective probability is *not objective enough* for many applications of probability. Many applications of probability follow a Bayesian statistical methodology: produce a *prior* probability function p_t , collect some evidence u , and draw predictions using the *posterior* probability function $p_{t+1}(v) = p_t(v|u)$. Now the prior function is determined before empirical evidence is available; this is matter of subjective choice for empirically-based subjectivists. However, the ensuing conclusions and predictions may be sensitive to this initial choice, rendering them subjective too. Yet such relativism is anathema in science: a disagreement between agents about a hypothesis should be arbitrated by evidence; it should be a fact of the matter, not mere whim, as to whether the evidence confirms the hypothesis.

That argument is rather inconclusive however. The proponent of empirically-based subjective probability can counter that scientists have simply over-estimated the extent of objectivity in science, and that subjectivity needs to be made explicit. Even if one grants a need for objectivity, one could argue that it is a pragmatic need: it just makes science simpler. The objective Bayesian must accept that it can not be empirical warrant that motivates the selection of a particular belief function from all those compatible with background knowledge, since all such belief functions are equally warranted by available empirical evidence. In the absence of any non-empirical justification for choosing a particular belief function, such a function can only be considered objective in a *conventional* sense. One can drive on the right or the left side of the road; but we must all do the same thing; by convention we choose the left. That does not mean that the left is objectively correct or most warranted—either side will do.

A second line of argument offers explicitly pragmatic reasons for selecting a particular belief function. If probabilities are subjective then measuring probabilities must involve elicitation of degrees of belief from agents. As developers of expert systems in AI have found, elicitation and the associated consistency-checking are prohibitively time-consuming tasks (the inability of elicitation to keep pace with the demand for expert systems is known as *Feigenbaum's bottleneck*). If a subjective approach is to be routinely applied throughout science it is clear that a similar bottleneck will be reached. On the other hand, if degrees of belief are objectively determined by background knowledge then elicitation is not required—degrees of belief are calculated by maximising entropy. Objective Bayesianism is thus to be preferred for reasons of efficiency.

Indeed many Bayesian statisticians now (often tacitly) appeal to non-committal objective priors rather than embark on a laborious process of introspection, elicitation or analysis of sensitivity of posterior to choice of prior.

⁴⁰(Paris and Vencovská, 1990; Paris, 1994; Paris and Vencovská, 2001)

A third motivating argument appeals to caution. In many applications of probability the risks attached to bold predictions that turn out wrong are high. For instance, a patient's symptoms may narrow her condition down to meningitis or 'flu, but there may be no empirical evidence—such as information about relative prevalence—to decide between the two. In this case, the risks associated with meningitis are so much higher than those associated with 'flu, that a non-committal belief function seems more appropriate as a basis for action than a belief function that gives the probability of meningitis to be zero, even though both are compatible with available information. (With a non-committal belief function one will not dismiss the possibility of meningitis, but if one gives meningitis probability zero one will disregard it.) High-risk applications thus favour cautious conclusions, non-committal degrees of belief and an objective Bayesian approach.

I argue in Williamson (2005b) that the appeal to caution is the most decisive motivation for objective Bayesianism, although pragmatic considerations play a role too.

§16

LANGUAGE DEPENDENCE

The maximum entropy principle has been criticised for being language or representation dependent: it has been argued that the principle awards the same event different probabilities depending on the way in which the problem domain is formulated.

John Maynard Keynes surveyed several purported examples of language dependence in his discussion of Laplace's Principle of Indifference.⁴¹ This latter principle advocates assigning the same probability to each of a number of possible outcomes in the absence of any knowledge which favours one outcome over the others. (Keynes added the condition that the possible outcomes must be indivisible.⁴²) The maximum entropy principle makes the same recommendation in the absence of background knowledge and so inherits any language dependence of the Principle of Indifference.

A typical example of language dependence proceeds as follows.⁴³ Suppose an agent's language can be represented by the propositional language $L = \{c\}$ with just one propositional variable c which asserts that a particular book is colourful. The agent has no background knowledge and so by the Principle of Indifference (or equally by the maximum entropy principle) assigns $p(c) = p(\neg c) = 1/2$. But now consider a second language $L' = \{r, b, g\}$ where r signifies that the book is red, b that it is blue and g that it is green. An agent with no knowledge will give $p(\pm r \wedge \pm b \wedge \pm g) = 1/8$. Now $\neg c$ is equivalent to $\neg r \wedge \neg b \wedge \neg g$, yet the former is given probability $\frac{1}{2}$ while the latter is given probability $\frac{1}{8}$. Thus the probability assignments of the Principle of Indifference and the maximum entropy principle depend on choice of language.

Paris and Vencovská (1997) offer the following resolution. They argue that the maximum entropy principle has been misapplied in this type of example: if

⁴¹(Keynes, 1921)

⁴²(Keynes, 1921, §4.21)

⁴³(Halpern and Koller, 1995, §1)

an agent refines the propositional variable c into $r \vee b \vee g$ one should consider not L' but $L'' = \{c, r, b, g\}$ and make the agent's knowledge, namely $c \leftrightarrow r \vee b \vee g$, explicit. If we do that then the probability function on L'' with maximum entropy, out of all those that satisfy the background knowledge (i.e. which assign $p(c \leftrightarrow r \vee b \vee g) = 1$), will yield a value $p(\neg c) = 1/2$. This is just the same value as that given by the maximum entropy principle on L with no background knowledge. Thus there is no inconsistency.

This resolution is all well and good if we are concerned with a single agent who refines her language. But the original problem may be construed rather differently. If *two* agents have languages L and L' respectively, and no background knowledge, then they assign two different probabilities to what we know (but they don't know) is the same proposition. There is no getting round it: probabilities generated by the maximum entropy principle depend on language as well as background knowledge.

Interestingly, language dependence in this latter multilateral sense is not confined to the maximum entropy principle. As Halpern and Koller (1995) and Paris and Vencovská (1997) point out, there is no non-trivial principle for selecting rational degrees of belief which is language-independent in the multilateral sense. More precisely, suppose we want a principle that selects a set \mathbb{O}_β of probability functions that are optimally rational on the basis of an agent's background knowledge β . If $\mathbb{O}_\beta \subseteq \mathbb{P}_\beta$, i.e. if every optimally rational probability function must satisfy constraints imposed by β , and if \mathbb{O}_β ignores irrelevant information inasmuch as $\mathbb{O}_{\beta \cup \beta'}(\theta) = \mathbb{O}_\beta(\theta)$ whenever β' involves no propositional variables in sentence θ , then the only candidate for \mathbb{O}_β that is multilaterally language independent is $\mathbb{O}_\beta = \mathbb{P}_\beta$.⁴⁴ Only empirically-based subjective probability is multilaterally language independent.

So much the better for empirically-based subjective probability and so much the worse for objective Bayesianism, one might think. But such an inference is too quick. It takes the desirability of multilateral language independence for granted. I argue in Williamson (2005a, Chapter 12) that an agent's choice of language embodies knowledge about the world:⁴⁵ knowledge about natural kinds, knowledge about which variables are relevant to which, and perhaps even knowledge about which partitions are amenable to the Principle of Indifference. For example, having dozens of words for snow in one's language says something about the environment in which one lives. Granted that language itself is a kind of background knowledge, and granted that an agent's degrees of belief should depend on her background knowledge, language independence becomes a rather dubious desideratum.

Note that while Howson (2001, p. 139) criticises the Principle of Indifference on account of its language dependence, the example he cites can be used to support the case *against* language independence as a desideratum. Howson considers two first-order languages with equality: L_1 has just a unary predicate Q while L_2 has unary Q together with two constants a and b . The explicit background knowledge β is just 'there are exactly 2 individuals', while sentence θ is 'something has the property Q '. L_1 has three models of β , which contain 0, 1 and 2 instances of Q respectively, so $p(\theta) = 2/3$. In L_2 individuals can be

⁴⁴(Halpern and Koller, 1995, Theorem 3.10)

⁴⁵Halpern and Koller (1995, §4) also suggest this tack, although they do not give their reasons. Interestingly, though, they do show in §5 that relaxing the notion of language independence leads naturally to an entropy-based approach.

distinguished by constants and thus there are eight models of β (if constants can name the same individual), six of which satisfy θ so $p(\theta) = 3/4 \neq 2/3$. While this is a good example of language dependence, the question remains whether language dependence is a problem here. As Howson himself hints, L_1 might be an appropriate language for talking about bosons, which are indistinguishable, while L_2 is more suited to talk about classical particles, which are distinguishable and thus able to be named by constants. Hence choice of language L_2 over L_1 indicates distinguishability, while conversely choice of L_1 over L_2 indicates indistinguishability. In this example, then, choice of language betokens implicit background knowledge. Of course all but the most ardent subjectivists agree that an agent's degrees of belief ought to be influenced by her background knowledge. Therefore language independence becomes an inappropriate desideratum.

In sum, while the Principle of Indifference and the maximum entropy principle have both been dismissed on the grounds of language dependence, it seems clear that some dependence on language is to be expected if degrees of belief are to adequately reflect implicit as well as explicit background knowledge. So much the better for objective Bayesianism, and so much the worse for empirically-based subjective probability which is language-invariant.

§17

COMPUTATION

There are important concerns regarding the application of objective Bayesianism. One would like to apply objective Bayesianism in artificial intelligence: when designing an artificial agent it would be very useful to have normative rules which prescribe how the agent's beliefs should change as it gathers information about its domain. However, there has seemed little prospect of fulfilling this hope, for the following reason. Maximising entropy involves finding the parameters $p(v)$ that maximise the entropy expression, but the number of such parameters is exponential in the number of variables in the domain, thus the size of the entropy maximisation problem quickly gets out of hand as the size of the domain increases. Indeed Pearl (1988, p. 468) has influentially criticised maximum entropy methods on account of their computational difficulties.

The computational problem poses a serious challenge for objective Bayesianism. However, recent techniques for more efficient entropy maximisation have largely addressed this issue. While no technique offers efficient entropy maximisation in all circumstances (entropy maximisation is an NP-complete problem), techniques exist that offer efficiency in a wide range of natural circumstances. I shall sketch my own approach here—this is developed in detail in Williamson (2005a, §§5.5–5.7).⁴⁶

Given a domain V of variables and some background knowledge β involving V which consists of a set of constraints on the agent's belief function p , one wants to find the probability function p , out of all those that satisfy the constraints in β , that maximises entropy. This can be achieved via the following procedure. First form an undirected graph on vertices V by linking

⁴⁶Maximum entropy methods have recently been applied to natural language processing, and other techniques for entropy maximisation have been tailored to that context—see Della Pietra et al. (1997) for example.

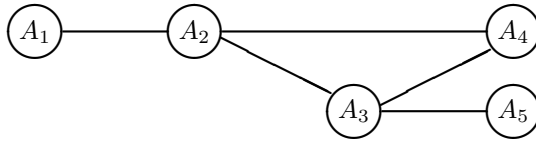


Figure 1: Example constraint graph.

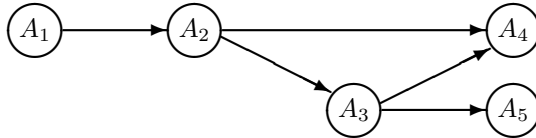


Figure 2: Example directed constraint graph.

pairs of variables that occur in the same constraint with an edge. For example, if $V = \{A_1, A_2, A_3, A_4, A_5\}$ and β contains a constraint involving A_1 and A_2 (e.g. $p(a_2|a_1) = 0.9$), a constraint involving A_2, A_3 and A_4 , a constraint involving A_3 and A_5 and a constraint involving just A_4 , then the corresponding undirected constraint graph appears in Fig. 1. The undirected constraint graph has the following crucial property: if a set Z of variables separates $X \subseteq V$ from $Y \subseteq V$ in the graph then the maximum entropy function p will render X and Y probabilistically independent conditional on Z .

Next transform the undirected constraint graph into a directed constraint graph, Fig. 2 in the case of our example.⁴⁷ The independence property ensures that the directed constraint graph can be used as a graph in a *Bayesian net* representation of the maximum entropy function p . A Bayesian net offers the opportunity of a more efficient representation of a probability function p : in order to determine p , one only needs to determine the parameters $p(a_i|par_i)$, i.e. the probability distribution of each variable conditional on its parents, rather than the parameters $p(v)$, i.e. the joint probability distribution over all the variables. Depending on the structure of the directed graph, there may be far fewer parameters in the Bayesian net representation. In the case of our example, if we suppose that each variable has two possible values then the Bayesian net representation requires 11 parameters rather than the 32 parameters $p(v)$ for each assignment v of values to V . For problems involving more variables the potential savings are very significant.

Roughly speaking, efficiency savings are greatest when each variable has few parents in the directed constraint graph, and this occurs when each constraint in β involves relatively few variables. Note that when dealing with large domains it tends to be the case that while one might make a large number of observations, each observation involves relatively few variables. For example, one might use hospital data as empirical observations pertaining to a large number of health-related variables, each department of the hospital contributing some statistics; while there might be a large number of such statistics, each statistic is likely to involve relatively few variables, namely those variables that are relevant to the department in question; such observations would yield a sparse constraint graph

⁴⁷The algorithm for this transformation is given in Williamson (2005a, §5.7).

and an efficient Bayesian net representation. Hence this method for reducing the complexity of entropy maximisation offers efficiency savings that are achievable in a wide range of natural situations.

§18

QUALITATIVE KNOWLEDGE

The maximum entropy principle has been criticised for yielding the wrong results when the agent's background knowledge contains qualitative causal information.⁴⁸ Daniel Hunter gives the following example:

The puzzle is this: Suppose that you are told that three individuals, Albert, Bill and Clyde, have been invited to a party. You know nothing about the propensity of any of these individuals to go to the party nor about any possible correlations among their actions. Using the obvious abbreviations, consider the eight-point space consisting of the events $ABC, ABC\bar{C}, A\bar{B}C$, etc. (conjunction of events is indicated by concatenation). With no constraints whatsoever on this space, MAXENT yields equal probabilities for the elements of this space. Thus $Prob(A) = Prob(B) = 0.5$ and $Prob(AB) = 0.25$, so A and B are independent. It is reasonable that A and B turn out to be independent, since there is no information that would cause one to revise one's probability for A upon learning what B does. However, suppose that the following information is presented: Clyde will call the host before the party to find out whether Al or Bill or both have accepted the invitation, and his decision to go to the party will be based on what he learns. Al and Bill, however, will have no information about whether or not Clyde will go to the party. Suppose, further, that we are told the probability that Clyde will go conditional on each combination of Al and Bill's going or not going. For the sake of specificity, suppose that these conditional probabilities are ... $[P(C|AB) = 0.1, P(C|A\bar{B}) = 0.5, P(C|\bar{A}B) = 0.5, P(C|\bar{A}\bar{B}) = 0.8]$.

When MAXENT is given these constraints ... A and B are no longer independent! But this seems wrong: the information about Clyde should not make A 's and B 's actions dependent.⁴⁹

But this counter-intuitive conclusion is attributable to a misapplication of the maximum entropy principle. The conditional probabilities are allowed to constrain the entropy maximisation process but the knowledge that Al's and Bill's decisions are causes of Clyde's decision is simply ignored. This failure to consider the qualitative causal background knowledge leads to the counter-intuitive conclusion.

Keynes himself had stressed the importance of taking qualitative knowledge into account and the difficulties that ensue if qualitative information is ignored:

Bernoulli's second axiom, that in reckoning a probability we must take everything into account, is easily forgotten in these cases of

⁴⁸(Pearl, 1988, p. 468; Hunter, 1989)

⁴⁹(Hunter, 1989, p. 91)

statistical probabilities. The statistical result is so attractive in its definiteness that it leads us to forget the more vague though more important considerations which may be, in a given particular case, within our knowledge.⁵⁰

Indeed, in the party example, the temptation is to consider only the definite probabilities and to ignore the important causal knowledge.

The party example and Keynes' advice highlight an important challenge for objective Bayesianism. In order that objective Bayesianism can be applied, all knowledge—*qualitative* as well as quantitative—must be taken into account. However, objective Bayesianism as outlined in §13 depends on background knowledge taking *quantitative* form: background knowledge must be explicated as a set of quantitative constraints on degrees of belief in order to narrow down a set of probability functions that satisfy those constraints. Thus the general challenge for objective Bayesianism is to show how qualitative knowledge can be converted into precise quantitative constraints on degrees of belief.

To some extent this challenge has already been met. In the case where qualitative knowledge takes the form of causal constraints, as in Hunter's party example above, I advocate a solution which exploits the following asymmetry of causality. Learning of a common cause of two events may warrant a change in the degrees of belief awarded to them: one may reason that if one event occurs, then this may well be because the common cause has occurred, in which case the other event is more likely—the two events become more dependent than previously thought. On the other hand, learning of a common effect would not warrant a change in degrees of belief: while the occurrence of one event may make the common effect more likely, this has no bearing on the other cause. This asymmetry motivates what I call the *Causal Irrelevance Principle*: if the agent's language contains a variable A that is known not to be a cause of any of the other variables, then her degrees of belief concerning these other variables should be the same as the degrees of belief she should adopt were she not to have A in her language (as long as any quantitative knowledge involving A is compatible with those degrees of belief). The Causal Irrelevance Principle allows one to transfer qualitative causal knowledge into quantitative constraints on degrees of belief—if domain $V = U \cup \{A\}$ then we have constraints of the form $p_{|U}^V = p^U$, i.e. the agent's belief function defined on V , when restricted to U , should be the same as the belief function defined just on U . By applying the Causal Irrelevance Principle, qualitative causal knowledge as well as quantitative information can be used to constrain the entropy maximisation process. It is not hard to see that use of the principle avoids counter-intuitive conclusions like those in Hunter's example: knowledge that Clyde's decision is a common effect of Al's and Bill's decision ensures that Al's and Bill's actions are probabilistically independent, as seems intuitively plausible. See Williamson (2005a, §5.8) for a more detailed analysis of this proposal.

Thus the challenge of handling qualitative knowledge has been met in the case of causal knowledge. Moreover, by treating logical influence analogously to causal influence one can handle qualitative logical knowledge using the same strategy (Williamson, 2005a, Chapter 11). But the challenge has not yet been met in other cases of qualitative knowledge. In particular, I claimed in §16 that choice of language implies knowledge about the domain. Clearly work remains

⁵⁰(Keynes, 1921, p. 322)

to be done to render such knowledge explicit and quantitative, so that it can play a role in the entropy maximisation process.

There is another scenario in which the challenge has not yet been fully met. Some critics of the maximum entropy principle argue that objective Bayesianism renders learning from experience impossible, as follows. The maximum entropy principle will, in the absence of background knowledge linking them, render outcomes probabilistically independent. Thus observing outcomes will not change degrees of belief in unobserved outcomes if there is no background knowledge linking them: observing a million black ravens will not shift the probability of the next raven being black from $\frac{1}{2}$ (which is the most non-committal value given only that there are two outcomes, black or not black). So, the argument concludes, there is no learning from experience. The problem with this argument is that we do have background knowledge that connects the outcomes—the qualitative knowledge that the outcomes all concern ravens—but this knowledge is mistakenly being ignored in the application of the maximum entropy principle. Qualitative knowledge should be taken into account so that learning from experience becomes possible—but how? Carnap (1952, 1971) addressed the problem, as have Paris and Vencovská (2003) more recently; however this line of work deals with rather simple frameworks (restricting attention to monadic predicates, for example) and even within the context of such frameworks it is clear that much remains to be done.

§19

INFINITE DOMAINS

The maximum entropy principle is most naturally defined on a finite domain—for example, a space of finitely many variables each of which takes finitely many values, as in §2. The question thus arises as to whether one can extend the applicability of objective Bayesianism to infinite domains. In the variable framework, one might be interested in domains with infinitely many variables, or domains of variables with an infinite range. Alternatively, one might want to apply objective Bayesianism to full generality of the mathematical framework of §3, or to infinite logical languages (§4). This challenge has been confronted, but at the expense of objectivity, as we shall now see.

There are two lines of work here, one of which proceeds as follows. Paris and Vencovská (2003) treat problems involving countable logical languages as limiting cases of finite problems. Consider a countably infinite domain $V = \{A_1, A_2, \dots\}$ of variables taking finitely many values, and schematic background knowledge β which may pertain to infinitely many variables. If $V_n = \{A_1, \dots, A_n\}$ and β_n is that part of β that involves only variables in V_n , then $p_{\beta_n}^{V_n}(u)$ can be found by maximising entropy as usual (here $u \in U \subseteq V_n$). Interestingly (see Paris and Vencovská (2003)) the limit $\lim_{n \rightarrow \infty} p_{\beta_n}^{V_n}(u)$ exists, so one can define $p_{\beta}^V(u)$ to be this limit. Paris and Vencovská (2003) show that this approach can be applied to very simple predicate languages and conjecture that it is applicable more generally to predicate logic.

In the transition from the finite to the infinite, the question arises as to whether countable additivity (introduced in §3) holds. Paris and Vencovská (2003) make no demand that this axiom hold. Indeed it seems that the type

of schematic background knowledge that they consider can not be used to express the knowledge that an infinite set of outcomes forms a partition. Thus the question of countable additivity can not be formulated in their framework. In fact, even if one were to extend the framework to formulate the question, the strategy of taking limits would be unlikely to yield probabilities satisfying countable additivity. If the only background knowledge is that E_1, \dots, E_n partition the outcome space, maximising entropy will give each event the same probability $1/n$. Taking limits will assign members of an infinite partition probability $\lim_{n \rightarrow \infty} 1/n = 0$. But then $\sum_{i=1}^{\infty} p(E_i) = 0 \neq 1$, contradicting countable additivity.

However, not only is countable additivity important from the point of view of mathematical convenience, but according to the standard betting foundations for Bayesian interpretations of probability introduced in §9, countable additivity *must* hold: an agent whose betting quotients are not countably additive can be Dutch booked.⁵¹ Once we accept countable additivity, we are forced either to concede that the strategy of taking limits has only limited applicability, or to reject the method altogether in favour of some alternative, as yet unformulated, strategy. Moreover, as argued in Williamson (1999), we are forced to accept a certain amount of subjectivity: a countably additive distribution of probabilities over a countably infinite partition must award some member of the partition more probability than some other member; but if background knowledge does not favour any member over any other then it is just a matter of subjective choice as to how one skews the distribution.

The other line of work deals with uncountably infinite domains. Jaynes (1968, §6) presents essentially the following procedure. First find a non-negative real function $m(x)$ that represents the invariances of the problem in question: if β offers nothing to favour x over y then $m(x) = m(y)$. Next, find a probability function p satisfying β that is closest to the invariance function m , in the sense that it minimises *cross-entropy* distance $d(p, m) = \int p(x) \log p(x)/m(x) dx$. It is this function that one ought to take as one's belief function p_β .⁵²

This approach generalises entropy maximisation on discrete domains. In the case of finite domains m can be taken to be a probability function, found by maximising entropy, and so the probability function p that is closest to it is just m itself. For instance, if the domain is of size n and there is no background knowledge, the invariance function m can be taken as the function that gives value $1/n$ to each member of the domain; this is a probability function so $p_\beta = m$. In the case of countably infinite domains m may not be a probability function: if β is empty then as discussed above m must award the same value, k say, to each member of the domain; however, such a function can not be a probability function since countable additivity fails; therefore one must choose a probability function closest to m . Here we minimise $d(p, m) = \sum p(x) \log p(x)/m(x) = \sum p(x) \log p(x) - \log k \sum p(x) = \sum p(x) \log p(x) - \log k$; this is minimised just when the entropy $-\sum p(x) \log p(x)$ is maximised.

By drawing this parallel with the discrete case we can see where problems arise: even if the constraints β on p are closed and convex, there may be no probability function closest to m or there may be more than one probability

⁵¹(Williamson, 1999)

⁵²Objective Bayesian statisticians have developed a whole host of techniques for obtaining invariance functions and uninformative probability functions—see e.g. Kass and Wasserman (1996). Berger and Pericchi (2001) discuss the use of such priors in statistics.

function closest to m . This latter case, non-uniqueness, means subjectivity: the agent can exercise arbitrary choice as which distribution of degrees of belief to select. Subjectivity can also enter at the first stage, choice of m , since there may be cases in which several different functions represent the invariances of a problem.⁵³

But does such subjectivity really matter? Perhaps not. Although objective Bayesianism strives for objectivity, it can hardly be blamed where little is to be found. If there is nothing to decide between two assignments of degrees of belief, then subjectivity simply does not matter. Under such a view, all the Bayesian positions—strict subjectivism, empirically-based subjective probability and objective Bayesianism—accept the fact that selection of degrees of belief can be a matter of arbitrary choice, they just draw the line in different places as to the extent of subjectivity. Strict subjectivists allow most choice, drawing the line at infringements of the axioms of probability.⁵⁴ Proponents of empirically-based subjective probability occupy a half-way house, allowing extensive choice but insisting that knowledge of frequencies as well as axioms of probabilities constrain degrees of belief. Objective Bayesians go furthest by also using logical constraints to narrow down the class of acceptable degrees of belief.

Moreover, arguably the infinite is just a tool to help us reason about the large but finite and discrete universe in which we live.⁵⁵ Just as we create infinite continuous geometries to reason about finite discrete space, we create continuous probability spaces to reason about discrete situations. In which case if subjectivity infects the infinite then we can only conclude that the infinite may not be as effective a tool as we would like for probabilistic reasoning. Such relativity merely urges caution when idealising to the infinite; it does not tell against objective Bayesianism.

§20

FULLY OBJECTIVE PROBABILITY

We see then that objectivity is a matter of degree and that while subjectivity may infect some problems, objective Bayesianism yields a high degree of objectivity. We have been focussing on what we might call *epistemic objectivity*, the extent to which an agent's degrees of belief are determined by her background knowledge. In applications of probability a high degree of epistemic objectivity is an important desideratum: disagreements as to probabilities can be attributed to differences in background knowledge; by agreeing on background knowledge consensus can be reached on probabilities.

While epistemic objectivity requires uniqueness relative to background knowledge, there are stronger grades of objectivity. In particular, the strongest grade of objectivity, *full objectivity*, i.e. uniqueness simpliciter, arouses philosophical interest. Are probabilities uniquely determined, independently of background knowledge? If two agents disagree as to probabilities must at least one of them

⁵³See Gillies (2000, pp. 37–49); Jaynes (1968, §§6–8) and Jaynes (1973). The determination of invariant measures has become an important topic in statistics—see Berger and Pericchi (2001).

⁵⁴Subjectivists usually slip in a few further constraints: e.g. known truths must be given probability 1, and degrees of belief should be updated by Bayesian conditionalisation.

⁵⁵(Hilbert, 1925)

be wrong, even if they disagree as to background knowledge? Intuitively many probabilities are fully objective: there seems to be a fact of the matter as to the probability that an atom of cobalt-60 will decay in 5 years, and there seems to be a fact of the matter as to the chance that a particular roulette wheel will yield a black on the next spin. (A qualification is needed. Chances can not be quite fully objective inasmuch as they depend on time. There might now be a probability just under 0.5 of cobalt-60 atom decaying in the next five years; after the event, if it has decayed its chance of decaying in that time-frame is 1. Thus chances need to be indexed by time.)

As indicated in §10, objective Bayesianism has the wherewithal to meet the challenge of accounting for intuitions about full objectivity. By considering some ultimate background knowledge β^* one can define fully objective probability $p^* = p_{\beta^*}$ in terms of the degrees of belief one ought to adopt if one were to have this ultimate background knowledge. This is the *ultimate belief* notion of chance.

What should be included in β^* ? Clearly it should include all information relevant to the domain at time t . To be on the safe side we can take β^* to include all facts about the universe that are determined by time t —the entire history of the universe up to and including time t . (Remember: this challenge is of philosophical rather than practical interest).

While the ultimate belief notion of chance is relatively straightforward to state, much needs to be done to show that this type of approach is viable. One needs to show that this notion can capture our intuitions about chance. Moreover, one needs to show that that account is coherent—in particular one might have concerns about circularity: if probabilistic beliefs are beliefs about probability, yet probability is defined in terms of probabilistic beliefs, then probability appears to be defined in terms of itself.

However, this apparent circularity dissolves when we examine the premisses of this circularity argument more closely. Indeed at most one premiss can be true. In our framework, ‘probability is defined in terms of probabilistic beliefs’ is true if we substitute ‘fully objective single-case probability’ or ‘chance’ for ‘probability’ and ‘degrees of belief’ for ‘probabilistic beliefs’: chance is defined in terms of degrees of belief. But then the first premiss is false. Degrees of belief are not beliefs about chance, they are partial beliefs about elements of a domain—variables, events or sentences. According to this reading ‘probabilistic’ modifies ‘belief’, isolating a type of belief; it does not specify the object of belief. On the other hand, if the first premiss is to be true and ‘probabilistic beliefs’ are construed as beliefs about probability, then the second premiss is false since chance is not here defined in terms of beliefs about probability. Thus neither reading permits the conclusion that probability is defined in terms of itself.

Note that Bayesian statisticians often consider probability distributions over probability parameters. These *can* be interpreted as degrees of belief about chances, where chances are special degrees of belief. But there is no circularity here either. This is because the degrees of belief about chances are of a higher order than the chances themselves. Consider for instance a degree of belief that a particular coin toss will yield heads. The present chance of the coin toss yielding heads can be defined using such degrees of belief. One can then go on to formulate the higher-order degree of belief that the chance of heads is 0.5. But this degree of belief is not used in the (lower order) definition of the chance itself, so there is no circularity. (One can go on to define higher and

higher order chances and degrees of belief—regress, rather than circularity, is the obvious problem.)

One can make a stronger case for circularity though. One can read the empirical norm of §13 as saying that degrees of belief ought to be set to chances where they are known.⁵⁶ Under such a reading the concept of rational degree of belief appeals to the notion of chance, yet in this section chances are being construed as special degrees of belief; circularity again. Here circularity is not an artifice of ambiguity of terms like ‘probabilistic beliefs’. However, as before, circularity does disappear under closer investigation. One way out is to claim that there are two notions of chance in play: a physical notion which is used in the empirical norm, and an ultimate belief notion which is defined in terms of degrees of belief. But this strategy would not appeal to those who find a physical notion of chance metaphysically dubious. An alternative strategy is to argue that any notion of chance in the formulation of an empirical norm is simply *eliminable*. One can substitute references to chance with references to the *indicators* of chance instead. Intuitively, symmetry considerations, physical laws and observed frequencies all provide some evidence as to chances; one can simply say that an agent’s degrees of belief should be appropriately constrained by her knowledge of symmetries, laws and frequencies. While this may lead to a rather more complicated formulation of the empirical norm, it is truer to the epistemological route to degrees of belief—the agent has direct knowledge of the indicators of chances rather than the chances themselves. Further, it shows how these indicators of chances can actually provide evidence for chances: knowledge of frequencies constrains degrees of belief, and chances are just special degrees of belief. Finally, this strategy eliminates circularity, since it shows how degrees of belief can be defined independently of chances. It does however, pose the challenge of explicating exactly how frequencies, symmetries and so on constrain degrees of belief—a challenge that (as we saw in §18) is not easy to meet.

The ultimate belief notion of chance is not quite fully objective: it is indexed by time. Moreover, if we want a notion of chance defined over infinite domains then, as the arguments of §19 show, subjectivity can creep in, for example in cases—if such cases ever arise—in which the entire history of the universe fails to differentiate between the members of an infinite partition. This mental, ultimate belief notion of chance is arguably more objective than the influential physical notion of chance put forward by David Lewis however.⁵⁷ Lewis accepts a version of the empirical norm which he calls the *Principal Principle*: knowledge of chances ought to constrain degrees of belief. However Lewis does not go on to advocate the ultimate belief notion of chance presented here: ‘chance is [not] the credence warranted by our total available evidence . . . if our total evidence came from misleadingly unrepresentative samples, that wouldn’t affect chance in any way.’⁵⁸ (Unrepresentative samples do not seem to me to be a real problem for the ultimate belief approach, because the entire history of the universe up to the time in question is likely to contain more information pertinent to an event than simply a small sample frequency—plenty of large samples of relevant events, and plenty of relevant qualitative information, for instance.) Lewis instead takes chances to be products of the best system of laws, the best way of systematising the universe. The problem is that the criteria for comparing systems of laws—a

⁵⁶See Williamson (2005a, §5.3).

⁵⁷(Lewis, 1980, 1994)

⁵⁸(Lewis, 1994, p. 475)

balance between simplicity and strength—seem to be subjective. What counts a simple for a rocket scientist may be complicated for a robot and vice versa.⁵⁹ This is not a problem that besets the ultimate belief account: as Lewis accepts, there does seem to be a fact of the matter as to how evidence should inform degrees of belief. Thus an ultimate belief notion of chance, despite being a mental rather than physical notion, suffers less from subjectivity than Lewis’ theory.

Note that Lewis’ approach also suffers from a type of circularity known as *undermining*. Because chances for Lewis are analysed in terms of laws, they depend not only on the past and present state of the universe, but also on the future of the universe: ‘present chances are given by probabilistic laws, plus present conditions to which those laws are applicable, and . . . those laws obtain in virtue of the fit of candidate systems to the whole of history.’⁶⁰ Of course, non-actual futures (i.e. series of events which differ from the way in which the universe will actually turn out) must have positive chance now, for otherwise the notion of chance would be redundant. Thus there is now a positive chance of events turning out in the future in such a way that present chances turn out differently. But this yields a paradox: present chances can not turn out differently to what they actually are. Lewis (1994) has to modify the Principal Principle to avoid a formal contradiction, but this move does not resolve the intuitive paradox. In contrast, under the ultimate belief account present chances depend on just the past and the present state of the universe, not the future, so present chances can not undermine themselves.

§21

PROBABILITY LOGIC

There are increasing demands from researchers in artificial intelligence for formalisms for normative reasoning that combine probability and logic. Purely probabilistic techniques work quite well in many areas but fail to exploit logical relationships that obtain in particular problems. Thus for example probabilistic techniques are applied widely in natural language processing,⁶¹ with some success, yet largely without exploiting logical sentence structure. On the other hand purely logical techniques take problem structure into account without being able to handle the many uncertainties inherent in practical problem solving. Thus automated proof systems for mathematical reasoning⁶² depend heavily on implementing logics but often fail to prioritise searches that are most likely to be successful. It is natural to suppose that systems which combine probability and logic will yield improved results. Formalisms that combine probability and logic would also be applicable to many new problems in bioinformatics,⁶³ from inducing protein folding from noisy relational data to forecasting toxicity from uncertain knowledge of deterministic chemical reactions in cell metabolism.

⁵⁹In response Lewis (1994, p. 479) just plays the optimism card: ‘if nature is kind to us, the problem needn’t arise.’

⁶⁰(Lewis, 1994, p. 482)

⁶¹(Manning and Schütze, 1999)

⁶²(Quaife, 1992; Schumann, 2001)

⁶³(Durbin et al., 1999)

In a *probability logic*, or *prolog* for short, probability is combined with logic in one or more of the following two ways:

EXTERNAL probabilities are defined on logical sentences rather than events,

INTERNAL logical sentences incorporate statements about probabilities.

In an *external* prolog, logical implications take the form:

$$\theta_1 : x_1, \dots, \theta_n : x_n \models \phi : y$$

Here $\theta_1, \dots, \theta_n, \phi \in S_L$ are sentences of a logical language L which does not contain probabilities and $x_1, \dots, x_n, y \in [0, 1]$ are the probabilities themselves. For example if $L = \{A_1, A_2, A_3, A_4, A_5\}$ is a propositional language on propositional variables A_1, \dots, A_5 , we might be interested in whether

$$A_1 \wedge \neg A_2 : .9, (\neg A_4 \vee A_3) \rightarrow A_2 : .2, A_5 \vee A_3 : .3, A_4 : .7 \models A_5 \rightarrow A_1 : .4$$

In an *internal* prolog, logical implications take the form:

$$\theta_1, \dots, \theta_n \models \phi$$

where $\theta_1, \dots, \theta_n, \phi \in S_{L_p}$ are sentences of a logical language L_p which contains probabilities. L_p might be a first-order language with equality containing a (probability) function p , predicates R, S, T and constants sorted into individuals a_i , events e_i and real numbers $x_i \in [0, 1]$, and we might want to know whether

$$p(e_1) = x_1 \vee R(a_3), \neg p(e_2) = x_1 \rightarrow T(a_5) \models R(a_5)$$

Note that an internal prolog might have several probability functions each with a difference interpretation.

In a *mixed* prolog, the probabilities may appear both internally and externally. A logical implication takes the form

$$\theta_1 : x_1, \dots, \theta_n : x_n \models \phi : y$$

where $\theta_1, \dots, \theta_n, \phi \in S_{L_p}$ are sentences of a logical language L_p which contains probabilities.

There are two main questions to be dealt with when providing semantics for a prolog: how are the probabilities to be interpreted? what is the meaning of the implication operator \models ?

The *standard semantics* remains neutral about the interpretation of the probabilities and deals with implication thus:

EXTERNAL $\theta_1 : x_1, \dots, \theta_n : x_n \models \phi : y$ holds if and only if every probability function p that satisfies the left-hand side (i.e., $p(\theta_1) = x_1, \dots, p(\theta_n) = x_n$) also satisfies the right-hand side (i.e. $p(\phi) = y$).

INTERNAL $\theta_1, \dots, \theta_n \models \phi$ if and only if every L_p -model of the left-hand side in which p is interpreted as a probability function is also a model of the right-hand side.

The difficulty with the standard semantics for an external logic is that of *underdetermination*. Given some premiss sentences $\theta_1, \dots, \theta_n$ and their probabilities x_1, \dots, x_n we often want to know what probability y to give to a conclusion sentence ϕ of interest. However, the standard semantics may give no answer to this question: often $\theta_1 : x_1, \dots, \theta_n : x_n \not\models \phi : y$ for any $y \in [0, 1]$, because probability functions that satisfy the left-hand side disagree as to the probability they award to ϕ on the right-hand side. The premisses underdetermine the conclusion. Consequently an alternative semantics is often preferred.

According to the *objective Bayesian semantics* for an external logic on a finite propositional language $L = \{A_1, \dots, A_N\}$, $\theta_1 : x_1, \dots, \theta_n : x_n \models \phi : y$ if and only if an agent whose knowledge is summed up by the constraints on the left-hand side (i.e. who ought to believe θ_1 to degree x_1, \dots, θ_n to degree x_n) ought to believe ϕ to degree y . As long as the constraints $\theta_1 : x_1, \dots, \theta_n : x_n$ are consistent, there will be a unique function p that maximises entropy and a unique $y \in [0, 1]$ such that $p(\phi) = y$, so there is no problem of underdetermination.

I shall briefly sketch just three of the principal proposals in this area.⁶⁴

Colin Howson put forward his account of the relationship between probability and logic in Howson (2001, 2003). Howson interprets probability as follows: ‘the agent’s probability is the odds, or the betting quotient, they currently believe fair, with the sense of ‘fair’ that there is no calculable advantage to either side of a bet at those odds.’⁶⁵ The connection with logic is forged by introducing the concept of consistency of betting quotients: a set of betting quotients is consistent if it can be extended to a single-valued function on all the propositions of a given logical language L which satisfies certain regularity properties. Howson then shows that an assignment of betting quotients is consistent if and only if it is satisfiable by a probability function.⁶⁶ Having developed a notion of consistency, Howson shows that this leads naturally to an external logic with the standard semantics: consequence is defined in terms of satisfiability by probability functions, as outlined above.⁶⁷

In Halpern (2003) Joseph Halpern studies the standard semantics for internal logics. In the propositional case, L is a propositional language extended by permitting linear combinations of probabilities $\sum_{i=1}^n a_i p_i(\psi_i) > b$ where $a_1, \dots, a_n, b \in \mathbb{R}$ and p_1, \dots, p_n are probability functions each of which represents the degrees of belief of an agent and which are defined over sentences ψ of L .⁶⁸ This language allows nesting of probabilities: for example $p_1(\neg(p_2(\theta) > 1/3)) > 1/2$ represents ‘with degree more than a half, agent 1 believes that agent 2’s degree of belief in θ is less than or equal to $\frac{1}{3}$.’ Note though that the language can not represent probabilistic independencies, which are expressed using multiplication rather than linear combination of probabilities, such as $p_1(\theta \wedge \phi) = p_1(\theta)p_1(\phi)$. Halpern provides a possible-worlds semantics for the resulting logic: given a space of possible worlds, a probability measure $\mu_{w,i}$ over this space for each possible world and agent, and a valuation function π_i for each possible world, $p_1(\psi) > 1/2$ is true at a world w if the measure $\mu_{w,1}$ of the set of possible worlds at which ψ is true is greater than half, $\mu_{w,1}(\{w' : \pi_{w'}(\psi) = 1\}) > 1/2$. Consequence is defined straightforwardly

⁶⁴Williamson (2002) presents a more comprehensive survey.

⁶⁵(Howson, 2001, 143)

⁶⁶(Howson, 2001, Theorem 1)

⁶⁷(Howson, 2001, 150)

⁶⁸(Halpern, 2003, §7.3)

in terms of satisfiability by worlds.

Halpern later extends the above propositional language to a first-order language and introduces frequency terms $\|\psi\|_X$, interpreted as ‘the frequency with which ψ holds when variables in X are repeatedly selected at random.’⁶⁹ Linear combinations of frequencies are permitted, as well as linear combinations of degrees of belief. When providing the semantics for this language, one must provide an interpretation for frequency terms, a probability measure over the domain of the language.

In Paris (1994) Jeff Paris discusses external logics in detail, in conjunction with the objective Bayesian semantics. In the propositional case, Paris proposes a number of common sense desiderata which ought to be satisfied by any method for picking out a most rational belief function for the objective Bayesian semantics, and goes on to show that the maximum entropy principle is the *only* method that satisfies these desiderata.⁷⁰ Later Paris shows how an external logic can be defined over the sentences of a first order logic—such a function is determined by its values over quantifier-free sentences.⁷¹ Paris then introduces the problem of learning from experience: what value should an agent give to $p(R(a_{n+1})|\pm R(a_1) \wedge \dots \wedge \pm R(a_n))$, that is, to what extent should she believe a new instance of R , given n observed instances?⁷² As mentioned in §§18, 19, Paris and Vencovská (2003) suggest that the maximum entropy principle may be extended to the first-order case to address this problem.

In the case of the standard semantics it is natural to look for a traditional proof theory to accompany the semantics:

EXTERNAL Given $\theta_1, \dots, \theta_n \in S_L, x_1, \dots, x_n \in [0, 1]$, find a mechanism for churning out all $\phi : y$ such that $\theta_1 : x_1, \dots, \theta_n : x_n \models \phi : y$.

INTERNAL Given $\theta_1, \dots, \theta_n \in S_{L_p}$, find a mechanism for churning out all $\phi \in S_{L_p}$ such that $\theta_1, \dots, \theta_n \models \phi$.

In a sense this is straightforward: the premisses imply the conclusion just if the conclusion follows from the premisses and the axioms of probability by deductive logic. Fagin et al. (1990) produced a traditional proof theory for the standard probabilistic semantics, for an internal propositional logic. As with propositional logic, deciding satisfiability is NP-complete. Halpern (1990) discusses a logic which allows reasoning about both degrees of belief and frequencies. In general, no complete axiomatisation is possible, though axiom systems are provided in cases where complete axiomatisation is possible. Abadi and Halpern (1994) consider first-order degree of belief and frequency logics separately, and show that they are highly undecidable. Halpern (2003) presents a general overview of this line of work.

Paris and Vencovská (1990) made a start at a traditional proof theory for a type of objective Bayesian logic, but express some scepticism as to whether the goal of a traditional proof system can be achieved.

A traditional proof theory, though interesting, is often not what is required in applications of an external logic. To reiterate, given some premiss sentences

⁶⁹(Halpern, 2003, §10.3)

⁷⁰(Paris, 1994, Theorem 7.9; Paris and Vencovská, 2001)

⁷¹(Paris, 1994, Chapter 11; Gaifman, 1964)

⁷²(Paris, 1994, Chapter 12)

$\theta_1, \dots, \theta_n$ and their probabilities x_1, \dots, x_n we often want to know what probability y to give to a conclusion sentence ϕ of interest—not to churn out all $\phi : y$ that follow from the premisses. Objective Bayesianism provides semantics for this problem, and it is an important question as to whether there is a calculus that accompanies this semantics:

OBPROGIC Given $\theta_1, \dots, \theta_n, x_1, \dots, x_n, \phi$, find y such that $\theta_1 : x_1, \dots, \theta_n : x_n \models \phi : y$.

It is known that even finding an approximate solution to this problem is NP-complete.⁷³ Hence the best one can do is to find an algorithm that is scalable in a range of natural problems, rather than tractable in every case. My own approach, presented in Williamson (2005a), deals with the propositional case but does not take the form of a traditional logical proof theory, involving axioms and rules of inference. Instead the proposal is to apply the computational methods of §17 to find a Bayesian net representation of the p that satisfies constraints $p(\theta_1) = x_1, \dots, p(\theta_n) = x_n$ and maximises entropy, and then to use the net to calculate $p(\phi)$. The advantage of using Bayesian nets is that if sufficiently sparse, they allow the efficient representation of a probability function and efficient methods for calculating marginal probabilities of that function. In this context, the net is sparse and the method scalable in cases where each sentence involves few propositional variables in comparison with the size of the language.

Consider an example. Suppose we have a propositional language $L = \{A_1, A_2, A_3, A_4, A_5\}$ and we want to find y such that

$$A_1 \wedge \neg A_2 : .9, (\neg A_4 \vee A_3) \rightarrow A_2 : .2, A_5 \vee A_3 : .3, A_4 : .7 \models A_5 \rightarrow A_1 : y$$

According to our semantics we must find p that maximises

$$H = - \sum p(\pm A_1 \wedge \pm A_2 \wedge \pm A_3 \wedge \pm A_4 \wedge \pm A_5) \log p(\pm A_1 \wedge \pm A_2 \wedge \pm A_3 \wedge \pm A_4 \wedge \pm A_5)$$

subject to the constraints,

$$p(A_1 \wedge \neg A_2) = .9, p((\neg A_4 \vee A_3) \rightarrow A_2) = .2, p(A_5 \vee A_3) = .3, p(A_4) = .7$$

One could find p by directly using numerical optimisation techniques or Lagrange multiplier methods. However, this approach would not be feasible on large languages—already we would need to optimise with respect to 2^5 parameters $p(\pm A_1 \wedge \pm A_2 \wedge \pm A_3 \wedge \pm A_4 \wedge \pm A_5)$.

Instead take the approach of §17:

STEP 1 Construct an undirected *constraint graph*, Fig. 1, by linking variables that occur in the same constraint.

As mentioned, the constraint graph satisfies a key property, namely separation in the constraint graph implies conditional independence for the entropy maximising probability function p . Thus A_2 separates A_5 from A_1 so $A_1 \perp\!\!\!\perp_p A_5 \mid A_2$, (p renders A_1 probabilistically independent of A_5 conditional on A_2).

STEP 2 Transform this into a *directed constraint graph*, Fig. 2.

⁷³(Paris, 1994, Theorem 10.6)

Now D -separation, a directed version of separation,⁷⁴ implies conditional independence for p . Having found a directed acyclic graph which satisfies this property we can construct a Bayesian net by augmenting the graph with conditional probability distributions:

STEP 3 Form a Bayesian network by determining parameters $p(A_i|par_i)$ that maximise entropy.

Here par_i are the states of parents of A_i . Thus we need to determine $p(A_1), p(A_2|\pm A_1), p(A_3|\pm A_2), p(A_4|\pm A_3 \wedge \pm A_2), p(A_5|\pm A_3)$. This can be done by reparameterising the entropy equation in terms of these conditional probabilities and then using Lagrange multiplier methods or numerical optimisation techniques. This representation of p will be efficient if the graph is sparse, that is, if each constraint sentence θ_i involves few propositional variables in comparison with the size of the language.

STEP 4 Simplify ϕ into a disjunction of mutually exclusive conjunctions $\bigvee \sigma_j$ (e.g. disjunctive normal form) and calculate $p(\phi) = \sum p(\sigma_j)$ by using standard Bayesian net algorithms to determine the marginals $p(\sigma_j)$.

In our example,

$$\begin{aligned} p(A_5 \rightarrow A_1) &= p(\neg A_5 \vee A_1) \\ &= p(\neg A_5 \wedge A_1) + p(A_5 \wedge A_1) + p(\neg A_5 \wedge \neg A_1) \\ &= p(\neg A_5|A_1)p(A_1) + p(A_5|A_1)p(A_1) + p(\neg A_5|\neg A_1)p(\neg A_1) \\ &= p(A_1) + p(\neg A_5|\neg A_1)(1 - p(A_1)) \end{aligned}$$

We thus require only two Bayesian net calculations to determine $p(A_1)$ and $p(\neg A_5|\neg A_1)$. These calculations can be performed efficiently if the graph is sparse and ϕ involves few propositional variables relative to the size of the domain.

A major challenge for the objective Bayesian approach is to see whether potentially efficient procedures can be developed for first-order predicate logic.

§22

CONCLUSION

If probability is to be applied it must be interpreted. Typically we are interested in single-case probabilities—e.g. the probability that I will live to the age of 80, the probability that my car will break down today, the probability that quantum mechanics is true. The Bayesian interpretation tells us what such probabilities mean: they are rational degrees of belief.

Subjective Bayesianism has the advantage that it is easy to justify—the Dutch book argument is all that is needed. But subjective Bayesianism does not successfully capture our intuition that many probabilities are objective.

If we move to objective Bayesianism what we gain in terms of objectivity, we pay for in terms of hard graft to address the challenges outlined above. (For this reason, many Bayesians are subjectivist in principle but tacitly objectivist in

⁷⁴(Pearl, 1988, §3.3)

practice.) These are just challenges though; none seem to present insurmountable problems. They map out an interesting and important research programme rather than reasons to abandon any hope of objectivity.⁷⁵

REFERENCES

- Abadi, M. and Halpern, J. Y. (1994). Decidability and expressiveness for first-order logics of probability. *Information and Computation*, 112(1):1–36.
- Berger, J. O. and Pericchi, L. R. (2001). Objective Bayesian methods for model selection: introduction and comparison. In Lahiri, P., editor, *Model Selection*, volume 38 of *Monograph Series*, pages 135–207. Beachwood, Ohio. Institute of Mathematical Statistics Lecture Notes.
- Bernoulli, J. (1713). *Ars Conjectandi*. cerebro.xu.edu/math/Sources/JakobBernoulli/ars.sung/ars.sung.html. Trans. Bing Sung.
- Billingsley, P. (1979). *Probability and measure*. John Wiley and Sons, New York, third (1995) edition.
- Carnap, R. (1952). *The continuum of inductive methods*. University of Chicago Press, Chicago IL.
- Carnap, R. (1971). A basic system of inductive logic part 1. In Carnap, R. and Jeffrey, R. C., editors, *Studies in inductive logic and probability*, volume 1, pages 33–165. University of California Press, Berkeley CA.
- Church, A. (1936). An unsolvable problem of elementary number theory. *American Journal of Mathematics*, 58:345–363.
- de Finetti, B. (1937). Foresight. its logical laws, its subjective sources. In Kyburg, H. E. and Smokler, H. E., editors, *Studies in subjective probability*, pages 53–118. Robert E. Krieger Publishing Company, Huntington, New York, second (1980) edition.
- Della Pietra, S., Della Pietra, V. J., and Lafferty, J. D. (1997). Inducing features of random fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(4):380–393.
- Durbin, R., Eddy, S., Krogh, A., and Mitchison, G. (1999). *Biological sequence analysis: probabilistic models of proteins and nucleic acids*. Cambridge University Press, Cambridge.
- Earman, J. (1992). *Bayes or bust?* MIT Press, Cambridge MA.
- Fagin, R., Halpern, J. Y., and Megiddo, N. (1990). A logic for reasoning about probabilities. *Information and Computation*, 87(1-2):277–291.
- Fetzer, J. H. (1982). Probabilistic explanations. *Philosophy of Science Association*, 2:194–207.

⁷⁵I am very grateful to Oxford University Press for permission to reprint material from Williamson (2005a) in Part I and Part II of this chapter.

- Gaifman, H. (1964). Concerning measures in first order calculi. *Israel Journal of Mathematics*, 2:1–18.
- Gaifman, H. and Snir, M. (1982). Probabilities over rich languages. *Journal of Symbolic Logic*, 47(3):495–548.
- Gillies, D. (2000). *Philosophical theories of probability*. Routledge, London and New York.
- Hacking, I. (1975). *The emergence of probability*. Cambridge University Press, Cambridge.
- Halpern, J. Y. (1990). An analysis of first-order logics of probability. *Artificial Intelligence*, 46:311–350.
- Halpern, J. Y. (2003). *Reasoning about uncertainty*. MIT Press, Cambridge MA.
- Halpern, J. Y. and Koller, D. (1995). Representation dependence in probabilistic inference. In Mellish, C. S., editor, *Proceedings of the 14th International Joint Conference on Artificial Intelligence (IJCAI 95)*, pages 1853–1860. Morgan Kaufmann, San Francisco CA.
- Hilbert, D. (1925). On the infinite. In Benacerraf, P. and Putnam, H., editors, *Philosophy of mathematics: selected readings*. Cambridge University Press (1983), Cambridge, second edition.
- Howson, C. (2001). The logic of Bayesian probability. In Corfield, D. and Williamson, J., editors, *Foundations of Bayesianism*, pages 137–159. Kluwer, Dordrecht.
- Howson, C. (2003). Probability and logic. *Journal of Applied Logic*, 1(3-4):151–165.
- Howson, C. and Urbach, P. (1989). *Scientific reasoning: the Bayesian approach*. Open Court, Chicago IL, second (1993) edition.
- Hunter, D. (1989). Causality and maximum entropy updating. *International Journal in Approximate Reasoning*, 3:87–114.
- Jaynes, E. T. (1957). Information theory and statistical mechanics. *The Physical Review*, 106(4):620–630.
- Jaynes, E. T. (1968). Prior probabilities. *IEEE Transactions Systems Science and Cybernetics*, SSC-4(3):227.
- Jaynes, E. T. (1973). The well-posed problem. *Foundations of Physics*, 3:477–492.
- Jaynes, E. T. (1979). Where do we stand on maximum entropy? In Levine, R. and Tribus, M., editors, *The maximum entropy formalism*, page 15. MIT Press, Cambridge MA.

- Jaynes, E. T. (1988). The relation of Bayesian and maximum entropy methods. In Erickson, G. J. and Smith, C. R., editors, *Maximum-entropy and Bayesian methods in science and engineering*, volume 1, pages 25–29. Kluwer, Dordrecht.
- Jaynes, E. T. (2003). *Probability theory: the logic of science*. Cambridge University Press, Cambridge.
- Kass, R. E. and Wasserman, L. (1996). The selection of prior distributions by formal rules. *Journal of the American Statistical Association*, 91:1343–1370.
- Keynes, J. M. (1921). *A treatise on probability*. Macmillan (1948), London.
- Kolmogorov, A. N. (1933). *The foundations of the theory of probability*. Chelsea Publishing Company (1950), New York.
- Laplace, P. S. m. d. (1814). *A philosophical essay on probabilities*. Dover (1951), New York.
- Lewis, D. K. (1980). A subjectivist’s guide to objective chance. In *Philosophical papers*, volume 2, pages 83–132. Oxford University Press (1986), Oxford.
- Lewis, D. K. (1994). Humean supervenience debugged. *Mind*, 412:471–490.
- Manning, C. D. and Schütze, H. (1999). *Foundations of statistical natural language processing*. MIT Press, Cambridge MA.
- Miller, D. (1994). *Critical rationalism: a restatement and defence*. Open Court, Chicago IL.
- Paris, J. B. (1994). *The uncertain reasoner’s companion*. Cambridge University Press, Cambridge.
- Paris, J. B. and Vencovská, A. (1990). A note on the inevitability of maximum entropy. *International Journal of Approximate Reasoning*, 4:181–223.
- Paris, J. B. and Vencovská, A. (1997). In defence of the maximum entropy inference process. *International Journal of Approximate Reasoning*, 17:77–103.
- Paris, J. B. and Vencovská, A. (2001). Common sense and stochastic independence. In Corfield, D. and Williamson, J., editors, *Foundations of Bayesianism*, pages 203–240. Kluwer, Dordrecht.
- Paris, J. B. and Vencovská, A. (2003). The emergence of reasons conjecture. *Journal of Applied Logic*, 1(3-4):167–195.
- Pearl, J. (1988). *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Morgan Kaufmann, San Mateo CA.
- Popper, K. R. (1934). *The Logic of Scientific Discovery*. Routledge (1999), London. With new appendices of 1959.
- Popper, K. R. (1959). The propensity interpretation of probability. *British Journal for the Philosophy of Science*, 10:25–42.

- Popper, K. R. (1983). *Realism and the aim of science*. Hutchinson, London.
- Popper, K. R. (1990). *A world of propensities*. Thoemmes, Bristol.
- Quaife, A. (1992). *Automated development of fundamental mathematical theories*. Kluwer, Dordrecht.
- Ramsey, F. P. (1926). Truth and probability. In Kyburg, H. E. and Smokler, H. E., editors, *Studies in subjective probability*, pages 23–52. Robert E. Krieger Publishing Company, Huntington, New York, second (1980) edition.
- Reichenbach, H. (1935). *The theory of probability: an inquiry into the logical and mathematical foundations of the calculus of probability*. University of California Press (1949), Berkeley and Los Angeles. Trans. Ernest H. Hutten and Maria Reichenbach.
- Rosenkrantz, R. D. (1977). *Inference, method and decision: towards a Bayesian philosophy of science*. Reidel, Dordrecht.
- Schumann, J. M. (2001). *Automated theorem proving in software engineering*. Springer-Verlag.
- Shannon, C. (1948). A mathematical theory of communication. *The Bell System Technical Journal*, 27:379–423 and 623–656.
- Venn, J. (1866). *Logic of chance: an essay on the foundations and province of the theory of probability*. Macmillan, London.
- von Mises, R. (1928). *Probability, statistics and truth*. Allen and Unwin, London, second (1957) edition.
- von Mises, R. (1964). *Mathematical theory of probability and statistics*. Academic Press, New York.
- Williamson, J. (1999). Countable additivity and subjective probability. *British Journal for the Philosophy of Science*, 50(3):401–416.
- Williamson, J. (2002). Probability logic. In Gabbay, D., Johnson, R., Ohlbach, H. J., and Woods, J., editors, *Handbook of the logic of argument and inference: the turn toward the practical*, pages 397–424. Elsevier, Amsterdam.
- Williamson, J. (2005a). *Bayesian nets and causality: philosophical and computational foundations*. Oxford University Press, Oxford.
- Williamson, J. (2005b). Motivating objective Bayesianism: from empirical constraints to objective probabilities. In Harper, W. L. and Wheeler, G. R., editors, *Probability and Inference: Essays in Honour of Henry E. Kyburg Jr.* Elsevier.