

PHI 533 Decision Theory

Session 2: Reflection Principles

September 25, 2002

1. Semantics/pragmatics/ logic of epistemic judgment

Hintikka/modal logic approach. Problem of the moron

"Entailment on pain of incoherence": initial characterization

New: opinions concerning my own opinions -- effect on the logic

Expression versus statement of my opinion: interpretation of "P" and "p". If $P(A) = x$ expresses my opinion, then if I express that plus anything contrary to $P(p_{\text{Now}}(A)=x) = 1$, I am in a Moore paradox situation ('pragmatic incoherence'). See further #12 below.

2. Normal and revolutionary changes in view.

Orthodox Bayesian criterion for [diachronically] coherent opinion management.

Distinguish: OK scenarios, OK policies -- difference in evaluation

Policies are framed with respect to an aim or goal. Setting aside an aim is not in itself irrational; but evaluation of a policy is w.r.t. the aim i.q.

DKL[ewis] proof (as reported by Teller) that we will follow the rule SC if we know beforehand (a) what the possible experience scenarios will be, and (b) what our posterior opinion will be after each scenario.

Question: what criteria are left when such conditions not imposed?

3. Introduction to the General Reflection Principle

General criterion of irrationality: **self-sabotage by one's own lights**

Pierone: has only yes-no degrees of belief; *The Forecasting Manual*

Principle: current belief in logical span of possible updated beliefs

Piero: subjective probabilities; *The New Forecasting Manual*

Principle: current probabilities to lie in logical span of possible updated probabilities

Formula, Notation: $P(A)$ is in the interval spanned by the numbers $p_t(A)$ such that $P(\text{my posterior opinion at } t \text{ will be } p_t) > 0$

4. Enter subjective expectation ["expectation value", "expected value"]

Pascal and Huygens: the geometry of deliberation

The Bayesian imperative (Maximize your expected gain) -- parochial!

Russian roulette example, to illustrate how it works.

Two scenarios: you are kidnapped and forced to play Russian roulette with a revolver that can hold six bullets, but actually is loaded with less.

You are asked how much you will pay to have one more bullet removed.

On first scenario, the revolver is loaded with five, on the second with one.

Question: would you want to pay the same or different amounts?

Formula, Notation:

value of v:	v_1	v_2	v_N
	X_1	X_2	X_N
probab:	p_1	p_2	p_N

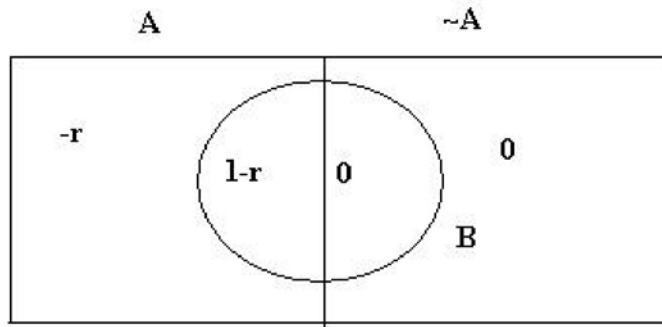
$$EXP_{\mathbf{p}}(v) = p_1v_1 + p_2v_2 + \dots + p_Nv_N = \Sigma\{ p_jv_j : j = 1, \dots, N\}$$

5. General Reflection Principle generalized

Principle: current probabilities and expectation values to lie in logical span of possible updated probabilities and expectation values respectively.

The quantity p_t . What are my expectations regarding this quantity?

6. Translation of a conditional probability into an expectation value:



Let quantity q have value $-r$ on $A \& \sim B$, have value $1-r$ on $A \& B$, and have value 0 everywhere else (i.e. on $\sim A$). Then we can deduce:¹

$$P(B|A) = r \text{ iff } EXP_{\mathbf{p}}(q) = 0$$

7. Deduction of the [Special] Reflection Principle

That is the principle: $P(A | p_t(A) = x) = x$

Consider the quantity q which (on all the [relevant] possibilities in cosmic history, why be modest) takes value:

- 0 if $p_t(A)$ is not x
- $(1-x)$ if $p_t(A) = x$ and A is true
- $-x$ if $p_t(A) = x$ and A is not true

Applying the preceding note we see that $P(A | p_t(A) = x) = x$ iff $EXP_{\mathbf{p}}(q) = 0$. Should it be? Well what can we deduce about $EXP_{\mathbf{p}_t}(q)$? If it must be 0 for any such posterior p_t then the General Reflection Principle entails that it must be 0 now as well. Will it?

8. **Corollary:** *My Current Probability Is My Expected Value Of Posterior Probability*

¹ Useful facts to remember for proofs; $P(A \& B) + P(A \& \sim B) = P(A)$;
 $P(A|B) = P(A \& B)/P(B)$, or equivalently (useful!!) that
 $P(A \& B) = P(B)P(A|B)$. (Note also that in this last equation, the order does not matter: conjunction is commutative ...)

"Expected" is meant in the sense of expectation value of course; and it does not apply to times after [mind-]death: conditional on which all possible posteriors have probability zero.

Notice that Reflection therefore does not imply that, from a later point of view, my current opinion is correct. It does imply that when I think about how I may have to change my current probability, I can see it possibly going up as much as down 'on the average'.

9. The Orthodox Bayesian agent obeys the Reflection Principle

Such an agent begins with any prior and gets a ticker tape of propositions as input, which s/he fully believes when they arrive, and conditionalizes upon (rule SC).

Mixture Principle

A *mixture* (also "convex combination", "weighted average") of two numbers or two functions w and z is the sum $aw + (1-a)z$ for some number a in $[0,1]$. If $0 < a < 1$ then this is a *proper mixture*.

An mixture of w and z lies in the closed interval spanned by w and z ; a proper mixture lies in the open interval spanned by w and z (i.e. is in between, not one of the end points w or z)

Application: Take any partition $\{X_i : i = 1, 2, \dots\}$. Then P is a mixture of $\{P(\cdot | X_i) : i = 1, 2, \dots\}$. Let $X_i = p_i^i(A)$, the i -th posterior probability for A , as foreseen now. Then $P(A)$ must lie in the interval those numbers span.

The same point can be elaborated to yield: *My expectation for quantity f is my expectation of the value of my posterior (at t) expectation value for f .* (Explored saliently by statistician Michael Goldstein; apparently in use among mathematical economists for quite some time before that.)

10. DKL was right (see #2 above)

The conditions of DKL's proof are realized in a controlled experiment of proper scientific design. There is a prior probability for each possible outcome, as well as for the hypothesis being tested, and it is known how the latter value will be updated in each possible outcome scenario.

The Reflection Principle implies that under these conditions the posterior probability is foreseen to be the result of conditionalizing on the outcome found. The proof uses the same idea as the deduction of the Special Reflection Principle from the General one (see # 7 above); details in "Conditionalization, a New Argument For".

It is crucial to this argument that it is known (or given, or fully believed) beforehand that the total input between now and then will identify one element in a given partition of the possibility space, and that the posterior will give probability one to that element.

11. Expressing doubts about my own future opinions

Initial impressions to the contrary, the Reflection Principle allows one to express strong doubts about the reliability of one's future opinions, though it also implies limits thereto. Consider the sentence:

$$(p_t(A) = x) \ \& \ A$$

If x is low, then any degree of belief in proposition expresses a doubt as to the reliability of that future opinion. So what can that degree of belief be?

$$\begin{aligned} P((p_t(A) = x) \ \& \ A) &= P(p_t(A) = x)P(A | p_t(A) = x) \\ &= x P(p_t(A) = x) \\ &\leq x \end{aligned}$$

So if I reflect on the possibility that A will seem unlikely to me tomorrow, then today it does not seem impossible to me that A is true nevertheless, but that eventuality does not seem more likely than *that*.²

Similarly, $P((p_t(A) = x) \ \& \ \sim A) \leq (1 - x)$.

12. Enter Moore sentences

A Moore sentence is one that could be true, but is not coherently believable.

Example: Peter knows that it is raining in Peking and that Paul does not believe that. He tells this to Paul. Should Paul now add "It is raining in Peking and I do not believe that" to his beliefs? In fact, if he does, he can certainly explain that each of these conjuncts is one that can be true, and that neither conflicts logically with the other. Yet his state of belief is incoherent, in the following sense: it is not possible for Paul to have only true beliefs if his beliefs include both.

Generalizing to subjective probability context: a Moore sentence is one which one can give a positive probability but cannot conditionalize on without arriving at an incoherent posterior opinion.

We have just seen examples of this, given the Reflection Principle (and therefore pertinent to orthodox Bayesian agents as well!), such as $(p_t(A) = x) \ \& \ \sim A$.

But we need to broaden this. Suppose I said "It is raining in Peking, but I do not believe that I believe that" or "It is raining in Peking, but I am not certain that I believe that". Those are equally cases of Moore sentences it would seem. That does not follow from the Reflection Principle, but it does follow from an assumption (stated at the outset - see #1) that functions in the deduction of Special from General Reflection.

13. Diagnosis of the Reflection Principle

The incoherence of a state of belief that includes belief that a given Moore sentence is true, is *not* that the content is unsatisfiable. The incoherence must therefore be explained/characterized at a different level: not semantics but pragmatics -- e.g. points regarding the use of "I".

The Reflection Principle declares certain sentences to be Moore sentences -- so it is applicable only if "P" is read in the first person, as expressing a state of opinion. No statement of a biographical fact about a person's state could be subject to such a

² For believers in objective chance, satisfying DKL's similar principle $P(A | ch(A) = x) = x$, we can add similar points about my possible probability for $((ch(A) = x) \ \& \ A)$.

constraint, for after all the statement could be true (and it can have any probability you like for someone else).

Any critique or defence of the Reflection Principle must therefore treat it as such. That our opinions may fluctuate wildly and uncontrollably due to mind-worms, drugs, illness, poor choice of breakfast cereal or life's companions -- that is all very true, but pertains mainly to what can be true by way of *statements of fact* about one's opinions.

14. Limitations of the Reflection Principle -- I: *the persistence of memory*

If $P(A) = 0$ or 1 , then I am also sure that all my posterior opinions will give 0 or 1 respectively to A .

(This comes from the proper mixing principle: 0 and 1 can here only appear as endpoints, so my current opinion cannot be a proper mixture of 0 and something else, yet be 0 .)

This is a general feature of probabilist models of opinion, and is seriously confronted only when we start looking at irreducible conditional probabilities -- not currently in our scope, to be taken up in Session 11 (currently scheduled for Dec. 4)

15. Limitations of the Reflection Principle -- II: *subjective time only*

An example due to John Collins shows that we must treat 'calendar time' as a random variable, and read the Reflection Principle as pertaining to my own 'subjective time'. Note that the subject in Collins' experiment can be assumed to have no memory loss or other impairment, and has a perfectly running personal clock from the outset of the experiment on -- but has some ignorance concerning how his personal time relates to 'objective time', while some of his opinions concern propositions about what happens as dated in 'objective time'. *Note also that he is a conditionalizer!*

Examination of a 'lost in time' example

I am a subject in a psych experiment (we suppose this for definiteness, and to contrast two points of view: the experimenter, who has perfect knowledge, and the subject's (mine) who does not.)

At the outset I am certain of all the following. It is either 10pm or 11pm (but I do not know which --let us say that each seems equally likely to me). A coin has been tossed (say, 1 hour ago) to decide between two possibilities. If the coin came up Heads then the light will be turned off in this room at midnight minus ϵ , a very small fraction.

Thus I am certain that I am in one of four possible worlds:

- w1: outset at 10pm, Heads
- w2: outset at 11pm, Heads
- w3: outset at 10pm, Tails
- w4: outset at 11pm, Tails

Let Q be the proposition [The coin came up Heads]. My initial probability for Q is $1/2$.

We will now examine my foreseen later opinions about Q at various times, and find that I am apparently violating the Reflection Principle.

Situation as seen by the experimenter who subjects me to this ordeal:

9PM	10PM	11PM	12MIDNIGHT
Heads [w1]	Exp starts	light on $p_{11}(w3)=0$	light off $p_{12}(Head)=1$
Tails [w2]	Exp starts	light on $p_{11}(w3)=0$ $p_{11}(w1,2,4)=1$	<i>light on</i> $p_{12}(w2,4)=1$ $p_{12}(Head)=0$
[w3]	Heads	Exp starts	light off $p_{12}(Heads)=1$
[w4]	Tails	Exp starts	light on $p_{12}(w3)=0$ $p_{12}(Heads)=1/3$

I see that at 12midnight I will have probability 1/3 for Heads iff I am in scenario w4, which however entails that the coin came up Tails. Thus I say: $P(\text{Heads} \mid p_{12}(\text{Heads})=1/3) = 0$ which certainly appears to violate the Reflection Principle.

Let us draw the table with all descriptions as the subject sees it:

1or 2 hrs ago	now	+1 hour	+2 hours
Heads 9pm [w1]	Exp starts 10 pm	light on $p_{+1}(w3)=0$	light off $p_{+2}(w1)=1$ $p_{+2}(Head)=1$
Tails 9pm [w2]	Exp starts 10pm	light on $p_{+1}(w3)=0$	<i>light on</i> $p_{+2}(w2,4)=1$ $p_{+2}(Head)=0$
Heads 10pm [w3]	Exp starts 11pm	light off $p_{+1}(w3)=1$	light off $p_{+2}(Heads)=1$
Tails 10pm [w4]	Exp starts 11pm	light on $p_{+1}(w3)=0$	<i>light on</i> $p_{+2}(w2,4)=1$ $p_{1a}(Head)=0$

No problem appears here; the current probability is properly related to the posteriors for +1 and +2 (see note³) I will draw from this the conclusion that the Reflection Principle should be recast in the following form:

Reflection Principle. $P(A \mid p_t A = r) = r$, where “t” is a relative time term for times in the future or present.

³ *Now* the probabilities that at time = +1, Heads will receive 1 and that Heads will receive 1/3 then are 1/4 and 3/4 respectively, "averaging out" to 1/2, the current probability.