

Coincidences and How to Think about Them

Elliott Sober

A Familiar Dialectic

The naïve see causal connections everywhere. Consider the fact that Evelyn Marie Adams won the New Jersey lottery twice. The naïve find it irresistible to think that this cannot be a coincidence. Maybe the lottery was rigged or perhaps some uncanny higher power placed its hand upon her brow. Sophisticates respond with an indulgent smile and ask the naïve to view Adams' double win within a larger perspective. Given all the lotteries there have been, it isn't at all surprising that someone would win one of them twice. No need to invent conspiracy theories or invoke the paranormal – the double win was a mere coincidence.

The naïve focus on a detailed description of the event they think needs to be explained. The *New York Times* reported Adams' good fortune and said that the odds of this happening by chance are 1 in 17 trillion; this is the probability that Adams would win both lotteries if she purchased a single ticket for each and the drawings were at random. In fact, the newspaper made a small mistake here. If the goal is to calculate the probability of Adams' winning those two lotteries, the reporter should have taken into account the fact that Adams purchased multiple tickets; the newspaper's very low figure should thus have been somewhat higher. However, the sophisticated response is that this modest correction misses the point. For sophisticates, the relevant event to consider is not that Adams' won those two lotteries, but the fact that someone won two state lotteries at some time or other. Given the many millions of people who have purchased lottery tickets, this is "practically a sure thing" (Diaconis and Mosteller 1989, Myers 2002).

Another example of reasoning about coincidence in which the same dialectic unfolds begins with the fact that my birthday (06061948) occurs at the 16,769,633th position of the decimal expansion of π (not counting the initial "3").¹ The probability of this occurring is very small, if numbers appear at random in the decimal expansion. The naïve conclude that my birthday's occurring at that exact position cannot be a mere coincidence; perhaps my date of birth was so arranged that the number 16,769,633 would provide me with an encrypted message that points the way to my destiny. The sophisticated reply that the probability of my birthday's occurring somewhere in the first 100 million digits is actually very high – about 2/3. Given this, there is no reason to think that my birth date's showing up where it does is anything but a coincidence.

How the Naïve and the Sophisticated Reason

The naïve and the sophisticated² agree about one thing but disagree about another. Both rely on a rule of inference that I will call probabilistic *modus tollens*. This is the idea that you should reject a hypothesis if it tells you that what you observe is enormously improbable. The naïve think that the

¹ Go to <http://www.angio.net/pi/piquery> to see if your birthday appears in the first 100 million digits.

² The naïve and the sophisticated are characters in my story; I do not mean to suggest that all sophisticated thinkers in the real world reason exactly in the way I'll describe the sophisticated as reasoning.

hypothesis of Mere Coincidence strains our credulity too much. Since the hypothesis of Mere Coincidence says that the probability of Adams's double win is tiny, we should reject that hypothesis. Sophisticates grant the authority of probabilistic *modus tollens*, but contend that the hypothesis of Mere Coincidence should be evaluated by seeing what it says about the observation that someone or other wins two state lotteries at some time or other. Since this is very probable according to the hypothesis of Mere Coincidence, we should decline to reject that hypothesis. The naïve and the sophisticated thus seem to agree on the correctness of probabilistic *modus tollens*. Their disagreement concerns how the event to be explained should be described.

Sophisticates avoid rejecting the hypothesis of Mere Coincidence by replacing a logically stronger description of the observations with one that is logically weaker. The statement

(1) Evelyn Adams, having bought four tickets for each of two New Jersey lotteries, wins both.

is logically stronger than the statement

(2) Someone at sometime, having bought some number of tickets for two or more lotteries in one or more states, wins at least two lotteries in a single state.

It is a theorem in probability theory that logically weakening a statement can't lower its probability – the probability will either go up or stay the same. In the case at hand, the probability goes up -- *way up*.

Diaconis and Mosteller (1989, p. 859) say that the relevant principle to use when reasoning about coincidences is an idea they term the *Law of Truly Large Numbers*. This says that “with a large enough sample, any outrageous thing is likely to happen.” They cite Littlewood (1953) as having the same thought; with tongue in cheek, Littlewood defined a miracle as an event whose probability is less than 1 in a million. Using as an example the U.S. population of 250 million people, Diaconis and Mosteller observe that if a miracle “happens to one person in a million each day, then we expect 250 occurrences a day and close to 100,000 such occurrences a year.” If the human population of the earth is used as the reference class, miracles can be expected to be even more plentiful.

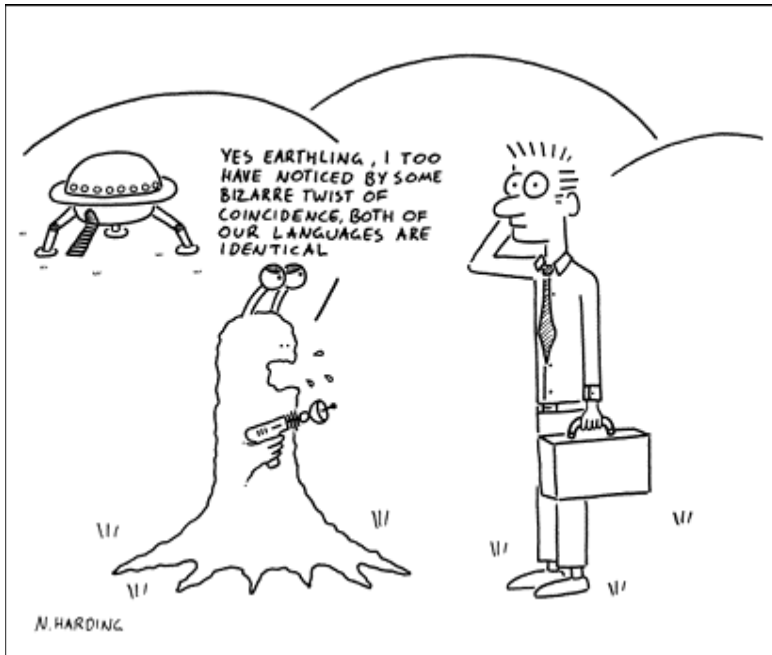
Two Problems for Sophisticates

Sophisticates bent on using probabilistic *modus tollens* should be wary about the strategy of replacing a logically stronger description of the observations with one that is logically weaker. The reason for wariness is that this strategy allows one to decline to reject hypotheses of Mere Coincidence no matter what they are and no matter what the data say. Even when there is compelling evidence that the observations should *not* be explained by this hypothesis, the hypothesis of Mere Coincidence can be defended by logically weakening the observations.

Consider, for example, Alfred Wegener's defense of the hypothesis of continental drift. Wegener noticed that the wiggles in the east coast of South America correspond rather exactly to the wiggles in the west coast of Africa. The pattern is as if a single sheet of paper were torn in two. He also noticed that the distribution of geological strata down one coast matches the distribution down the other. In addition, he observed that the distribution of organisms down the two coasts – both fossilized and

extant – shows the same detailed correlation. Wegener argued that this systematic matching should not be explained by the hypothesis of Mere Coincidence. His preferred alternative was that the continents had once been in contact and then had drifted apart.

Wegener encountered intense opposition from geophysicists, who didn't see how continents could plough through the ocean floor. I will return to this criticism later. My present point is that it would have been bizarre to counter Wegener's argument by weakening the data. A sophisticate bent on retaining the hypothesis of Mere Coincidence could point out that there are billions of planets in the universe that contain continents separated by wide oceans. If wiggles in coast lines and distributions of geological strata and of organisms are in each continent independently caused, there surely will exist at least one pair of continents on some planet or other that exhibits the kind of matching that Wegener found so interesting. With the data suitably weakened, probabilistic *modus tollens* tells you not to reject the hypothesis of Mere Coincidence.



A similar point is illustrated by the accompanying cartoon. If life forms from another planet turned out to speak English, the irresistible inference would be that we and they have had some sort of contact in the past. The idea that the detailed resemblance of the two languages is a Mere Coincidence strains our credulity too much. However, if we wish to hold fast to the belief that the resemblance is a Mere Coincidence, we can avoid having probabilistic *modus tollens* force us to reject that hypothesis merely by weakening our description of what the two languages have in common. Instead of focusing on the fact that the two languages match in a thousand specific ways, we can restrict our attention to the modest fact that both contain nouns. We then can reply that it isn't at all surprising that two languages

should both contain nouns if they developed independently; after all, nouns are useful.³ Notice that I just weakened the description of the data in a way that differs from the kind of weakening I considered in connection with Wegener. I didn't ask what the probability is that somewhere in the universe two languages would match even though they evolved independently (which is not to deny that that question might lead to the same conclusion). This brings out a further problem with the strategy of weakening the data at will. There are many ways to weaken the data. Which weakening should one employ? Why not simply replace the data with a tautology?

I began by noting that the naïve seem to think that *nothing* is a Mere Coincidence. Sophisticates who weaken their description of the data to avoid rejecting hypotheses of Mere Coincidence seem to think that *everything* is a Coincidence. These sophisticates are not just sophisticated – they are *jaded*. No correlation, no matter how elaborate and detailed, impresses them. In fact, none *can* impress them; their trick of weakening the data works against all comers. What we need is guidance on when the description of the data may be weakened, not the imperative to always do so, or the permission to do so whenever we please.

Statistics provides guidance on the question of when one's description of the data may be weakened. It is given in the theory of *sufficient statistics*. R.A. Fisher (1959) introduced this idea in the context of his theory of point estimation. Suppose you want to estimate a coin's probability θ of landing heads when tossed and you assume that the tosses are independent and identically distributed (i.i.d.) -- each toss has the same probability of landing heads and the results on some tosses don't influence the probabilities of others. To figure out which estimate is best, you toss the coin 1000 times, say, and obtain a particular sequence of heads and tails. Do you need to use this exact sequence as your description of the data, or can you just attend to the number of heads (which, let us suppose, was 503)? As it happens, this weaker description *suffices*; it is a sufficient statistic in the sense that it captures all the evidentially relevant information that the exact sequence contains. More specifically, the frequency of heads is a sufficient statistic in the context of using *maximum likelihood estimation* (MLE) as one's method for estimating θ because

$$(3) \quad \frac{\Pr(\text{the exact sequence} \mid \theta = p)}{\Pr(\text{the exact sequence} \mid \theta = q)} = \frac{\Pr(\text{the number of heads} \mid \theta = p)}{\Pr(\text{the number of heads} \mid \theta = q)}.$$

In all these conditional probabilities, I assume that the coin was tossed 1000 times. The reason (3) is true is that

$$(4) \quad \Pr(\text{the exact sequence} \mid \theta = x) = x^{503}(1-x)^{497}$$

and

$$(5) \quad \Pr(\text{number of heads} \mid \theta = x) = \binom{1000}{503} x^{503}(1-x)^{497}$$

³ Darwin (1859, ch. 13) argued that adaptive characters provide poor evidence of common ancestry and that it is useless characters that provide more powerful evidence. Darwin (1871, ch. 6) also noticed the parallel epistemological problems connecting historical linguistics and phylogenetic inference.

This is why the left-hand and right-hand ratios in (3) must have the same value. The maximum likelihood estimate of θ is the same whether you use the stronger or the weaker description of the data, and the likelihood ratio of that best estimate, compared to any inferior estimate, will be the same, again regardless of which description of the data you use. Notice that what counts as a sufficient statistic depends on the method of inference you use and on the range of possible hypotheses you want to consider.⁴ In the example just described, MLE is the method used and the assumption is that tosses are i.i.d. If MLE were used in the context of testing whether tosses are independent of each other, the number of heads would not be a sufficient statistic; information about the exact sequence would additionally be relevant.

With these ideas in mind, let's return to the example of Evelyn Marie Adams' double win in the New Jersey lottery. If we use probabilistic *modus tollens*, the weakened description of the data given in (2) is *not* endorsed by the idea of sufficient statistics. The point is that shifting from (1) to (2) makes a difference in the context of probabilistic *modus tollens*, whereas shifting from (4) to (5) does not matter from the point of view of MLE under the i.i.d. assumption. Shifting from a highly specific description of the data to one that is logically weaker is often permissible, but that is not enough to justify the sophisticate's pattern of reasoning about Adams. If a statistic is sufficient, you are *permitted* to shift to that weaker description of the data; you are not *obliged* to do so. The shift is permissible when and only when it doesn't change what you infer.

The problem of whether to weaken one's description of the evidence, and how to do so, is a problem for the sophisticate, not for the naïve. However, there is a second problem that both must face -- both rely on probabilistic *modus tollens*. This is a form of inference that no one should touch with a stick. The similarity between *modus tollens* and its probabilistic analog may suggest that the latter must be legitimate because the former is deductively valid; however, this is an illusion. *Modus tollens* says that if H entails O and O turns out to be false, that you may conclude that H is false. Probabilistic *modus tollens* says that if $\Pr(O \mid H)$ is very high and O turns out to be false, that you may likewise conclude that H is false. My beef with probabilistic *modus tollens* is not that the conclusion does not deductively follow from the premises. I've drawn a double line between premises and conclusion in Prob-MT below to acknowledge that this is so, but that isn't enough to rescue the principle. Rather, my objection is that the occurrence of an event that a hypothesis says is very improbable is often evidence *in favor* of the hypothesis, not evidence *against* it. What is evidence in favor of H cannot be a sufficient reason to reject H.

(MT) If H then O. not-O. ----- not-H	(Prob-MT) Pr(O H) is very high. not-O. =====	not-H
---	---	-------

Consider, for example, the use of DNA testing in forensic contexts. DNA evidence can be used to draw an inference about whether two individuals are related (for example, in paternity suits) or to draw an inference about whether a person suspected of a crime was at the crime scene. In both cases,

⁴ Notice also that the argument that appeals to (3) to show that the number of heads is a sufficient statistic depends on using the likelihood *ratio* as the relevant method for comparing the two estimates. If the *difference* in the likelihoods were used instead, the corresponding equality would not be true. How one measures weight of evidence matters; see Fitelson (1999) for further discussion.

you begin by determining whether two DNA samples match. This may seem to be a context in which probabilistic *modus tollens* is plausible. Suppose two individuals match at the loci examined, and that the probability of this match is only, say, 6.5×10^{-38} , if the two individuals are unrelated. This may seem to provide ample grounds for rejecting the hypothesis that the individuals are unrelated. However, what is missing from this exercise is any representation of how probable the data would be if the individuals *were* related. Crow *et al.* (2000, p. 66) discuss an example of this sort in which two individuals match at 13 loci for genes that happen to be rare. Crow calculated the above figure of 6.5×10^{-38} as the probability of the data under the hypothesis that the individuals are unrelated. However, it also is true that if the individuals were full sibs, the probability of the match would be 7.7×10^{-32} . Surely it would be absurd to apply probabilistic *modus tollens* twice over, first rejecting the hypothesis that the two individuals are unrelated and then rejecting the hypothesis that they are related. In fact, the data lend support to the hypothesis that the two individuals are sibs; it would be wrong to use the data to reject that hypothesis. The evidence *favors* the hypothesis that the two individuals are sibs over the hypothesis that they are unrelated because the observations are more probable under the first hypothesis than they are under the second. This is the Law of Likelihood (Hacking 1965, Edwards 1973, Royall 1997). It isn't the *absolute value* of the probability of the data under a single hypothesis that matters; rather, the relevant issue is how two such probabilities *compare*. The Law of Likelihood allows for the possibility that evidence may differentially support a hypothesis even though the hypothesis says that the evidence was very improbable. Notice also that the Law of Likelihood avoids an embarrassing question that defenders of probabilistic *modus tollens* must answer – how improbable is improbable enough for the hypothesis to be rejected? Defenders of *modus tollens* have had to admit that this question has only a conventional answer.

What I have dubbed probabilistic *modus tollens* is known in statistics as Fisher's test of significance. According to Fisher (1959, p. 39), you have two choices when a hypothesis says that your observations are very improbable. You can say that the hypothesis is false or that something very improbable has just occurred. Fisher was right about the disjunction. However, what does not follow is that the hypothesis is false; in fact, as just noted, it doesn't even follow that you have obtained evidence against the hypothesis (Hacking 1965, Edwards 1973, Royall 1997).

When the naïve and the sophisticated reasoned about whether Evelyn Marie Adams' double win was a Mere Coincidence, both helped themselves to probabilistic *modus tollens*. We need to understand this problem without appealing to that faulty rule of inference. The sophisticated also seemed to allow themselves to violate the Principle of Total Evidence. They were happy to substitute a weaker description of the data for a stronger one, even though that changed the conclusion that the rule of inference they use instructs one to draw. We need to explain why the naïve are wrong to think that nothing is a Mere Coincidence without violating that principle. This may seem to return us to square one, but it does not. There is something right about the sophisticate's demand that the data about Evelyn Adams be placed in a wider perspective. We need to consider not just her double win, but the track records that others have had and whether she bought tickets in other lotteries that did not turn out to be winners. However, moving to this wider data set does not involve *weakening* the initial description of the data, but *adding* to it; the key is to make the data stronger, not weaker.

Coinciding Observations, Coincidence Explanations, and Reichenbach's Principle of the Common Cause

Before I continue, some regimentation of vocabulary is in order. First of all, what is a coincidence? Diaconis and Mosteller (1989, p. 853) suggest a working definition: a coincidence is “a surprising concurrence of events, perceived as meaningfully related, with no apparent causal connection.” This is a good start, but it does entail that whether something is a coincidence is a subjective matter. There are two elements in this definition that we should separate. First, there is the idea of *coinciding observations*. When you and I meet on a street corner, our locations coincide. The same is true of the east coast of South America and the west coast of Africa – their wiggles, geological strata, and biogeography coincide. And perhaps it doesn't offend too much against the rules of English usage to say that the person who won the New Jersey lottery in one week “coincides” with the person who won it a few weeks later. Observations coincide when they are similar in some respect. There is no need to be precise about how much (or what kind of) similarity is required for two observations to coincide, since the main point is to distinguish the observations from a kind of hypothesis that might be offered to explain them. Here we need the idea of a *coincidence explanation*. A coincidence explanation asserts that the observations are not causally connected. By this I mean that neither causes the other and they do not have a common cause. Thus, to say that it is a coincidence that two events are similar is to suggest a certain kind of explanation; each event was produced via a separate and independent causal process. Saying that the similarity of the observations is a coincidence does *not* mean that the similarity is inexplicable or that there is no need to explain the similarity. Understood in this way, it is an objective matter whether a given coincidence explanation is true, assuming as I will that causation is an objective matter.

With coinciding observations distinguished from coincidence explanations, we can kick away the ladder and see that coinciding observations are not required for the question to arise of whether a hypothesis of Causal Connectedness is superior to a hypothesis of Mere Coincidence. We sometimes need to consider this choice when the observations exhibit a pattern of *dissimilarity*. Cartwright (1994, p. 117) suggests the following example. Suppose I go shopping each week at a grocery store with \$10 to spend. I spend some portion of the \$10 on meat and the rest on vegetables. When you observe my cash register receipts over the course of a year, you see that I rarely if ever spend exactly \$5 on the one and exactly \$5 on the other. The dollar amounts do not coincide. But the fact that they always sum to \$10 is not a coincidence. They are two effects of a common cause. So observations need not be similar for the question of coincidence to arise. If you and I always order different desserts when we dine together at a restaurant, the waiter may be right to suspect that this is not a coincidence.

The match, or mismatch, of two *token* events is a rather small data set. When there are many pairs of token events, a pattern between *kinds* of events may emerge. Based on the relative frequencies of kinds of events, one may infer that a correlation, either positive or negative, exists. Correlation is a probabilistic concept. Dichotomous event types A and B are positively correlated precisely when $\Pr(A\&B) > \Pr(A)\Pr(B)$. Cartwright's shopping example involves a negative correlation; let A = my spending more than \$5 on meat and B = my spending more than \$5 on vegetables. If you infer the probabilities of these two event types from their frequencies in the data set describing my 52 trips to the grocery store, you'll infer that $\Pr(A) \approx \Pr(B) \approx \frac{1}{2}$, but that $\Pr(A\&B) = 0$. Given a correlation (positive or negative), the question is whether the pattern of matching (or mismatching) of the token events

should be explained by saying that the correlates are causally connected or by saying that the correlation is a mere coincidence.

Reichenbach (1956) elevated our natural preference for hypotheses of causal connection to the status of a metaphysical principle.⁵ His *principle of the common cause* says that whenever two events are correlated, the explanation must be that the two correlates are causally connected. This principle is central to recent work on causal modeling and directed graphs (Spirtes, Glymour, and Shines 2001; Pearl 2000, Woodward 2003). I think it is better to treat Reichenbach's idea as an epistemological principle that should be evaluated in terms of the Law of Likelihood (Sober 1988a, 1988b, 2001). The question is whether a hypothesis of Causal Connection renders the observations more probable than does the hypothesis of Mere Coincidence. When this is so, the evidence favors the first hypothesis over the second; it does not guarantee that the Causal Connection hypothesis must be true.⁶

Reichenbach was able to show that the fact that two events are correlated deductively follows from a certain type of Common Cause model, one in which the postulated common cause raises the probability of each effect and renders them conditionally independent. Viewed from the point of view of the Law of Likelihood, Reichenbach's argument can be adapted to cases in which the *explanandum* is the coinciding of two token events, rather than the correlation of two event types (Sober 1988b). And the *mismatch* of two events sometimes points towards a common cause explanation and away from a separate cause explanation, depending again on the details of how the common cause and separate cause hypotheses are formulated. Thus, in a wide range of cases, the question of whether it is a mere coincidence that the two events E_1 and E_2 occurred can be addressed by comparing the likelihood of the hypothesis of Causal Connection with the likelihood of the hypothesis of Mere Coincidence.

The Limits of Likelihood

The Law of Likelihood is a useful tool in the project of reasoning about coincidences, but it doesn't provide the complete epistemology we need. The problem is that likelihood considerations favor hypotheses of causal connection in contexts in which this seems to be the wrong diagnosis of which of the competing hypothesis is better. Evelyn Adams won the lottery twice. Under the hypothesis that these events were causally unconnected and that each win was due to a random draw from the tickets purchased, the probability of the observations is very small. It is easy to construct hypotheses of Causal Connection that have much higher likelihoods. One of them says that her winning the first time was a random event, but that the occurrence of that first win guaranteed that she would win the next time. Another says that both lotteries were rigged so that she would win. This latter hypothesis has a likelihood than which none greater can be conceived; it has a likelihood of unity. The Law of Likelihood seems to endorse the naïve impulse to see conspiracies everywhere, to always think that a hypothesis of Causal Connection is better than the hypothesis of Mere Coincidence.

⁵ I do not use the term "metaphysical" here in the pejorative sense associated with logical positivism. Rather, I use the term in contrast with "epistemological." The former has to do with the way the world is, the latter with the beliefs we should form about the world.

⁶ One reason that Reichenbach's principle should not be formulated metaphysically is the fact that it is at least a defensible position to maintain that quantum mechanics describes event types that are lawfully correlated but not causally connected. Arthur Fine has pointed out to me that these correlations also show that my categories of Mere Coincidence and Causal Connection are not exhaustive.

Bayesianism provides a natural solution to this type of problem for a wide range of cases. If prior probabilities can be defended by appeal to evidence, and aren't merely reflections of someone's subjective degrees of belief, then perhaps the likelihood advantage that conspiracy theories have can be overcome. Do we know that most state lotteries are fair? If so, this frequency data allows us to estimate the value of the prior probability that a given lottery is fair. If the value of this defensible prior is high enough, we may be able to show that the conspiracy theory about Clarke's double win has a low posterior probability even if it has a high likelihood.

The Limits of Bayesianism

The problem with this Bayesian solution is that there are lots of cases in which it isn't possible to back up assignments of prior probability with evidence and yet we still feel that there is something fishy about conspiracy theories and other hypotheses of causal connection..

In discussing the example of Wegener and continental drift, I noted earlier that the hypothesis of Continental Drift has a much higher likelihood than the hypothesis of Continental Stasis: $\Pr(\text{Data} \mid \text{Drift}) \gg \Pr(\text{Data} \mid \text{Stasis})$. However, this doesn't settle the matter of which hypothesis has the higher posterior probability. To decide that question, we must say something about the values of the prior probabilities, $\Pr(\text{Drift})$ and $\Pr(\text{Stasis})$. Geophysicists rejected Wegener's theory on the grounds that it was impossible for the continents to plough through the ocean floor. Biologists and other friends of continental drift replied that this, or something like it, had to be possible, since the data are overwhelming. One aspect of the controversy that retarded the achievement of consensus was the way in which Wegener formulated his hypothesis. He could have restricted himself to the claim that the continents were once in contact, and not hazarded a guess about how they moved apart. He did not do this; as noted, he argued that the continents move across the ocean floor. He turned out to be right about the general claim, but wrong about the specifics. The continents don't move across the ocean floor. Rather, they and the ocean floor move together, riding on plates that slide across the viscous material deeper inside the earth.

A Bayesian will represent the disagreement between critics and defenders of the drift hypothesis by saying that they had different prior probabilities. Since the likelihoods overwhelmingly favor Drift over Stasis, the critics must have assigned to the drift hypothesis a prior probability that was incredibly small. Were they rational to do so? Or should they have assigned the hypothesis a somewhat larger prior, one that, though still small, allowed the data to give the drift hypothesis the higher posterior probability? It is hard to see how there can be an objective answer to that question. The prior probabilities were not estimated from frequency data. It's not as if a team of scientists visited a large number of planets, recording in each case whether the continents move, and then estimating from that data how probable it is that the continents move here on earth. Of course, there's another possible source of objective probabilities – ones that are derived from a well-confirmed theory. Did geophysicists have such a theory? If so, what probability did that theory entail for the hypothesis of continental drift? If the theory entails that continental drift is impossible, the Bayesian has a problem. The problem derives from the fact that a hypothesis assigned a prior probability of zero cannot have its probability increase, no matter what the evidence is. This is why Bayesians usually advise us to assign priors of zero only to truth-functional contradictions. Following this advice, we should decline to assign continental drift a prior of zero, even if our best confirmed theories say that drift is impossible. But what small prior should one then choose? If we choose a value that is extremely tiny, Drift will have a

lower posterior probability than Stasis, even though Drift has the higher likelihood. If the prior probability is assigned a value that is a bit bigger, though still very small, Drift will end up with the larger posterior probability. No wonder the two communities were so divided. It is hard to see how the Bayesian can help us decide what the correct assignment of prior probabilities is. Different groups of scientists had different degrees of belief; that appears to be all one can say.

Another scientific problem exhibits the same pattern. Consider the fact that the correlation of the phases of the moon and the tides were known for hundreds of years. It was not until Newton's theory of gravity that a systematic explanation of the correlation was developed. Newton's theory says that the two events are causally connected – the moon exerts a gravitational attraction on the earth's surface, with the result that there are tides. It is an objective matter that this hypothesis of causal connection has a higher likelihood than the hypothesis that says that it is a Mere Coincidence that the tides and the phases of the moon coincide: $\Pr(\text{data} \mid \text{Newtonian Theory}) \gg \Pr(\text{data} \mid \text{Mere Coincidence})$. But does that mean that Newtonian theory is more probable than the hypothesis that the moon and the tides are causally unconnected? That depends on one's choice of priors. If $\Pr(\text{Newtonian Theory})$ isn't enormously tiny, then $\Pr(\text{Newtonian Theory} \mid \text{data}) > \Pr(\text{Mere Coincidence} \mid \text{data})$. But if Newtonian theory is assigned a small enough prior, the theory will not be more probable than the hypothesis of Mere Coincidence. Unfortunately, there appears to be no objective basis for assigning priors in one way rather than the other.

Does a Bayesian analysis provide a convincing explanation of why Evelyn Adams' double win on the New Jersey lottery should be thought of as a Mere Coincidence? We need priors on the two hypotheses. Does any of us have frequency data on how often state lotteries, and the lottery in New Jersey specifically, are fixed? Surely if fixes occur, the parties will have every reason to prevent them from becoming public. How often they will succeed is another matter. My hunch is that the slogan "the truth will out" is an exaggeration, and how often the truth outs is more or less unknown. For this reason, we should be somewhat reluctant to interpret absence of evidence as evidence of absence.⁷ I do not say that there is no objective basis for assigning prior probabilities here. However, it would be nice if an analysis of this problem could be developed that did not require this.

Models for a Larger Data Set

Imagine that we have data on all the people who bought tickets in all the New Jersey lotteries that have ever occurred, as well as information on who won what. Evelyn Adams's double win is part of this large data set, but only a small part. I want to consider a variety of models that might be offered for these multiple lotteries. What I mean by a "model" will be clarified in due course. To simplify discussion, I'll assume that there is just one winner in each lottery.

The first model I'll consider says that each lottery is fair – each ticket in a lottery has the same probability of winning:

(FAIR) If ticket t is purchased in lottery i ($1 \leq i \leq r$), $\Pr(t \text{ wins} \mid t \text{ was purchased in lottery } i) = \alpha_i$.

The FAIR model is an r -fold conjunction:

⁷ There is an observation selection effect here; for discussion, see Sober (2004b).

$$\begin{aligned} \Pr(t \text{ wins} \mid t \text{ was purchased in lottery 1}) &= \alpha_1. \\ \Pr(t \text{ wins} \mid t \text{ was purchased in lottery 2}) &= \alpha_2. \\ &\dots \\ \Pr(t \text{ wins} \mid t \text{ was purchased in lottery } r) &= \alpha_r. \end{aligned}$$

By assigning a different parameter to each lottery, FAIR allows, *but does not require*, that the probability a given ticket has of winning one lottery differs from the probability another ticket has of winning another. Notice also that this model doesn't say what the probability is of a ticket's winning any lottery. Those probabilities must be estimated from the data. In each lottery i , there are n_i tickets sold and exactly one ticket was the winner. This means that the maximum likelihood estimate (the MLE) of α_i is $1/n_i$.

The second model I'll describe is more complicated than FAIR. It assigns a separate parameter to each player-lottery pair:

$$\begin{aligned} \text{(PL)} \quad &\text{If ticket } t \text{ is purchased in lottery } i \ (1 \leq i \leq r) \text{ by player } j \ (1 \leq j \leq s), \\ &\Pr(t \text{ wins} \mid t \text{ was purchased in lottery } i \text{ by player } j) = \beta_{ij}. \end{aligned}$$

This model is a conjunction that contains $r(s)$ conjuncts. It allows for the possibility that some or all the lotteries are unfair, but does not require this. The MLE of β_{ij} for player j on lottery i is 0 if the player lost, and $1/n_{ij}$ if the player won, where n_{ij} is the number of tickets the player purchased on that lottery.

The third model I'll consider is even more complicated. Like the one just described, it treats each player-lottery pair as a separate problem, but it introduces the possibility that different tickets purchased by the same player on the same lottery may have different probabilities of winning.

$$\begin{aligned} \text{(PLT)} \quad &\text{If ticket } t \text{ is the } k\text{th ticket purchased } (1 \leq k \leq n) \text{ in lottery } i \ (1 \leq i \leq r) \text{ by player } j \ (1 \leq j \leq s), \\ &\Pr(t \text{ wins} \mid t \text{ is the } k\text{th ticket purchased in lottery } i \text{ by player } j) = \gamma_{ijk}. \end{aligned}$$

This model is a conjunction with rsn conjuncts. Notice that FAIR has the smallest number of parameters of the models described so far, and that PL and PLT both say that each lottery might be unfair but need not be.

The fourth and last model I'll consider (not that there aren't many others), involves circling back to the beginning to find a model that is even simpler than FAIR. FAIR allows that tickets in different lotteries may have different probabilities of winning. This is why that model has r parameters in it, one for each lottery. If we constrain tickets in all lotteries to have the same probability of winning, we obtain the following one-parameter model:

$$\text{(ONE)} \quad \text{If ticket } t \text{ is purchased in any lottery, } \Pr(t \text{ wins} \mid t \text{ was purchased in a lottery}) = \delta.$$

In a sense, this model says the lotteries have a greater degree of "fairness" than FAIR itself asserts. According to FAIR, players who buy a ticket in one lottery might have better odds than players who buy a ticket in another. The ONE model stipulates that this isn't so – every ticket in every lottery is in the same boat.

These different conceptualizations of how the lotteries work are “models” in the sense of that term that is standard in statistics. Each contains adjustable parameters whose values can be estimated from the data. To clarify how these models are related to each other, let me describe two of their properties. First, notice that the models are nested; they are linked to each other by the relation of logical implication:

$$\text{ONE} \rightarrow \text{FAIR} \rightarrow \text{PL} \rightarrow \text{PLT}$$

Logically stronger models are special cases of models that are logically weaker. A stronger model can be obtained from a weaker one by stipulating that various parameters in the weaker model have equal values. Because of this, FAIR cannot be more probable than either PL or PLT, regardless of what the data are. Bayesians who want to argue that one of the simpler models has a higher prior or posterior probability than a model that is more complex might reply that the right way to set up models is to ensure that they are incompatible with each other; they should not be nested. This imperative requires that we compare ONE with FAIR*, PL*, and PLT*, where each of the starred models stipulates that different parameters must have different values. Now there is no logical barrier to stipulating that FAIR has a higher prior probability than either PL* or PLT*. However, it is questionable whether there is a convincing reason to think that this stipulation is true. Is it really more probable that all tickets have exactly the same probability of winning a lottery than that they differ, if only by a little? I myself think it is very improbable that lotteries are *exactly* fair; I think they are no better than so-called fair coins. I think coins in the real world have probabilities of landing heads that are *approximately* $\frac{1}{2}$, not *exactly* $\frac{1}{2}$. The other property of these models that I want to mention concerns the likelihoods they have when adjustable parameters are replaced by their maximum likelihood estimates. What I want to consider, for example, is not $\text{Pr}(\text{data} \mid \text{FAIR})$, but $\text{Pr}[\text{data} \mid \text{L}(\text{FAIR})]$, where L(FAIR) denotes the instance of FAIR obtained by assigning values to its parameters that make the data most probable. The point of interest here is that L(FAIR) can't have a higher likelihood than either L(PL) or L(PLT).⁸ Increasing the number of adjustable parameters allows the resulting, more complex, model to fit the data better. In fact, the two most complex models, PL and PLT, are so complex that L(PL) and L(PLT) both say that Evelyn Adams was certain to win the two lotteries she did win, and that the winners of the other lotteries also had probabilities of unity of winning theirs. L(PLT) goes even farther; it says, not just that Adams was certain to win each of those two lotteries, but that it was a certainty that the tickets that won the two lotteries for her would do so. L(PL) doesn't go that far; if Adams purchased multiple tickets on one of the lotteries she won, L(PL) says that those tickets had equal probabilities of winning.

Comparing these models leads to a point that I think is of the first importance in our quest to understand how we should reason about coincidences. The naïve think that nothing is a Mere Coincidence. And the explanations they suggest for coinciding observations often seem to be very simple. When the naïve propose to explain Adams' double win by saying that the two lotteries were fixed, it would seem perverse to complain that this is a complicated explanation. What's so complicated about it? However, if we view this explanation as deriving from a model whose parameters are estimated from the data, and if we require that model to address a data set that is considerably more inclusive than these two facts about Adams, it turns out that the model that the naïve are implicitly using is vastly complex. They seem to be using a model that, when fitted to the data, says that each event that

⁸ L(FAIR) can't have a higher likelihood than L(PL*) or L(PLT*), either.

occurred had to occur. The hypothesis that all state lotteries have been FAIR is much simpler. Understanding the epistemic relevance of simplicity would throw light on the problem at hand.

Simplicity and Model Selection

Not only do we need to consider a larger data set instead of focusing exclusively on Adams's double win; we also must adjust our conception of what the goals are in model evaluation. The point is not just to find a model that in some sense summarizes the data we have, but a model that will do a good job predicting data that we do not yet have. For example, suppose we were to use data on past New Jersey lotteries to compare models where our goal is to figure out which model will allow us to make the most accurate predictions about next year's lotteries. Of course, there's no getting around the Humean point that we have no assurance that future lotteries will play by the rules that governed past lotteries. But let us assume that this is true. How can we use the old data to estimate how well models will do in predicting new data?

Scientists who work on empirical problems by trying out multiple models inevitably learn that hugely complicated models often do a poor job predicting new data when fitted to old data. These models are able to *accommodate* the old data; as noted earlier, adding parameters to a model will allow it to fit the data better, and if M is sufficiently complex, $\Pr[\text{old data} \mid L(M)] = 1$. However, $\Pr[\text{new data} \mid L(M)]$ will often be very low, or, more precisely, the distance between the predicted values and the observed values in the new data will often be great. This doesn't lead scientists to think that they should use the simplest possible model to make predictions. Rather, some sort of trade-off is needed – the best model of the candidate models considered will embody the most nearly optimal trade-off between its fit to old data and its simplicity. How is that optimal balancing to be ascertained? Is it a matter of art, but not of science? Must young scientists simply work away at a given problem and gradually develop a feel for what works? Is this the “tacit dimension” that Polanyi (1966) discussed? Well, there's no substitute for practical experience. However, there is, in addition, a body of results in mathematical statistics that shows that it is not a mere coincidence that very complicated models often make very inaccurate predictions. One central result in this literature is a theorem due to H. Akaike (1973), which says that

An unbiased estimate of the predictive accuracy of model $M \approx \log[\Pr(\text{data} \mid L(M))] - k$,

where k is the number of adjustable parameters in M . Akaike's theorem shows how good fit-to-data, as measured by the log-likelihood, improves expected predictive accuracy, while complexity, as measured by the number of adjustable parameters, diminishes that expectation. It also specifies a precise rate-of-exchange between log-likelihood and simplicity. It tells you how much of an improvement in fit-to-data is needed for the shift from a simpler to a more complex model to embody a net improvement in expected predictive accuracy.

Akaike's theorem is the basis for the Akaike Information Criterion (AIC), which scores a model by computing $-2[\log[\Pr(\text{data} \mid L(M))] - k]$; the best model will have the lowest AIC value. There are other model selection criteria on the market. Most of them are intended to help one identify models that are predictively accurate, and most of them include a penalty for complexity;⁹ for discussion, see

⁹ Cross validation makes no explicit mention of simplicity, but shares with AIC the goal of finding models that will be predictively accurate. It is interesting that there is a form of cross-validation (“take-one-out” cross validation) that is

Burnham and Anderson (2002). There seems to be a broad consensus that different model selection criteria are appropriate for different inference problems.

If we use AIC to evaluate different models of the New Jersey lotteries, what will be the upshot? That will depend on the data, but not only on the data. L(FAIR) will have a lower log-likelihood than L(LP) and L(LPT), but that doesn't ensure that FAIR is the worst of the three. The reason is that FAIR is far simpler than LP and LPT. It would not be surprising if FAIR scored better than these two more complicated models, but I cannot assert that this is true, since I have not looked at the data. However, the relevant epistemological point is visible without us having to carry out this set of calculations. FAIR may be a better model of the New Jersey lotteries than models like LP and LPT, which say that one or all of the lotteries may have been rigged, even though L(FAIR) has a lower likelihood than L(LP) and L(LPT).

The model selection framework is not a magic bullet that will instantaneously convert the naïve into sophisticates. The naïve might reject the goal of predictive accuracy; they also may insist on focusing just on Adams' double win, and refuse to consider the other data that constitute the history of the New Jersey Lottery. If they do so, they will have built a mighty fortress. If you look just at the double win, and don't want anything besides a hypothesis of maximum likelihood, there is no denying that the hypothesis that the two lotteries were twice fixed to ensure that Adams would win beats the pants off the hypothesis that the two lotteries were fair.¹⁰ But if you are prepared to ask the data to help you decide among the models just described, it may turn out that the FAIR model is superior to the PL and the PLT models. It is interesting that you don't have to evaluate the prior probabilities of PL and PLT to see what is wrong with these models.¹¹

As noted before, it may be possible to provide an objective Bayesian treatment of Adams' double win. Even though the FIX hypothesis has a higher likelihood than the FAIR hypothesis, perhaps there is a way to justify an assignment of prior probabilities that has the consequence that the FAIR hypothesis has the higher posterior probability. I used the lottery example to illustrate how the model selection approach works because it is easy to describe a variety of models. The reader may wonder, however, whether this approach can be applied to problems in which there seems to be little hope of providing objective prior probabilities. This is important, if the virtue of that approach is that it takes over where Bayesianism leaves off. In fact, I think that the two examples I gave before of inference problems in which objective priors are not available can be fruitfully represented as problems of model selection. The examples, recall, were Wegener's explanation of the resemblance of the coasts of South America and Africa and Newton's explanation of the correlation of the phases of the moon and the tides. The key

asymptotically equivalent with AIC (Stone 1977).

¹⁰ It might be suggested that the hypothesis that the two lotteries were fixed to ensure that Adams would win is a hypothesis that would occur to you only *after* you observe Adams' double win, and that it is a rule of scientific inference that hypotheses must be formulated *before* the data are gathered to test them. This temporal requirement is a familiar idea in frequentist statistics. For discussion, see Hitchcock and Sober (2004). It is a point in favor of the model selection approach that one does not have to invoke this temporal requirement to explain what is wrong with the PL and the PLT models.

¹¹ Just as conspiracy hypotheses that contain no adjustable parameters can be invented that have maximum likelihood, it also is true that "conspiracy models" that contain adjustable parameters can be formulated that have very good AIC scores. Consider a model that stipulates that Adams was certain to win the two lotteries she entered but which is otherwise just like FAIR in what it says about *other* lotteries. This model will have a better AIC score than FAIR. We have here an example of what Forster and Sober (1994) call *the subfamily problem*.

to the analysis in both cases is the fact that unifying explanations, which postulate a Causal Connection between the correlated events, use models that are simpler than those that articulate the hypothesis that the correlation as a Mere Coincidence. The unified models include parameters that apply to the two correlated events, whereas the disunified models have different parameters for the two correlates (Forster and Sober 1994; Sober 2003). In these cases, both likelihood and simplicity may favor models that postulate Causal Connection over models that assert Mere Coincidence.¹²

Conclusion

Having come this far – from probabilistic *modus tollens* to the Law of Likelihood to Bayesianism and then to model selection – let’s return to an idea I mentioned towards the beginning. This is Diaconis and Mosteller’s (1989, p. 859) *Law of Truly Large Numbers*, which says that “with a large enough sample, any outrageous thing is likely to happen.” This principle implicitly assumes a certain type of model. As Diaconis and Mosteller are well aware, it isn’t true in a suitably arranged deterministic model that any outrageous thing is likely to happen with enough trials, and the same point applies to many models that are probabilistic. The heuristic value of their principle is that it recommends that we look at the world in a certain way – we should use models that say that coinciding events can and do occur as Mere Coincidences, and have very high probabilities of doing so when the sample size is very large. But what are the rules of inference that recommend such models above others? The *Law of Truly Large Numbers* is not intended to address this question.

The Law of Likelihood allows us to compare hypotheses of Mere Coincidence with hypotheses of Causal Connection, but seems unable to identify a respect in which the first type of hypothesis is superior to the second. This is especially clear when the Causal Connection Hypothesis is deterministic and the Mere Coincidence hypothesis is probabilistic. The Bayesian response to this problem is to assign prior probabilities. Sometimes these can be justified by appeal to evidence; at other times, they seem to be merely subjective. It is in the latter kind of case that model selection criteria seem like a breath of fresh air. However, it also is interesting to consider inference problems in which prior probabilities *can* be based on evidence, and the two approaches still seem to disagree.

Some years ago, cognitive psychologists discussed the phenomenon of “hot hands” in sports. Everyone with even the most superficial familiarity with professional basketball believes that players occasionally have “hot hands.” When players are hot, their chance of scoring improves, and team-mates try to feed the ball to them. However, a statistical analysis of scoring patterns in the NBA yielded the result that one cannot reject the null hypothesis that each player has a constant probability of scoring throughout the season (Gilovich *et al.* 1985). Statistically sophisticated scientists concluded that belief in hot hands is a “cognitive illusion.” A scoring streak is not due to the player’s getting hot, but is a Mere Coincidence. Basketball mavens reacted to this statistical pronouncement with total incredulity.

What would a Bayesian analysis of this problem look like? Surely we have lots of evidence that physical injury, influenza, upset stomach, lack of sleep, and migraine impair athletic performance. The idea that a player’s probability of scoring through the season is *absolutely constant* should therefore

¹² Phylogenetic inference provides another context in which model selection criteria seem to be promising, in part because prior probabilities apparently lack an objective interpretation; see Sober (2004a) for discussion.

be assigned a very low prior probability. For this reason, Bayesianism seems predestined to side with common sense on this issue. I do not see this as a defect in Bayesianism, nor do I have any sympathy with the argument that defends the null hypothesis by pointing out that the data do not sanction its rejection. Is this another case of probabilistic *modus tollens*' rearing its ugly head? In any event, the model selection framework provides a very different and useful perspective.

Recall that the goal in model selection is to find models that will be predictively accurate. It is an important philosophical fact about this framework that false models can sometimes be better predictors than true ones (Sober 2002). Bayesians are right to say that the null hypothesis has very low prior and posterior probabilities. The idea that players never waiver in their scoring probabilities, even a little, *is* preposterous. However, this doesn't settle which model will make the most accurate predictions. Presumably, the truth about basketball players is very complex. Their scoring probabilities change as subtle responses to a large number of interacting causes. Given this complexity, players and coaches may make better predictions by relying on simplified models. Hot hands may be a reality, but trying to predict when players have hot hands may be a fool's errand.

Acknowledgements

I thank Ellery Eells, Arthur Fine, Malcolm Forster, George Gale, Daniel Hausman, Stephen Leeds, Wouter Meijs, and David Myers and for their comments.

References

- Akaike, H. (1973): "Information Theory as an Extension of the Maximum Likelihood Principle." In B. Petrov and F. Csaki (eds.), *Second International Symposium on Information Theory*. Budapest: Akademiai Kiado, pp. 267-281.
- Burnham, K. and Anderson, D. (2002): *Model Selection and Inference – a Practical Information-Theoretic Approach*. New York: Springer, 2nd edition.
- Cartwright, N. (1994): *Nature's Capacities and their Measurement*. Oxford: Oxford University Press.
- Crow, J., Budowle, B, Erlich, H., Lederberg, J., Reeder, D., Schumm, J., Thompson, E., Walsh, P., Weir, B. (2000): "The Future of Forensic DNA Testing -- Predictions of the Research and Development Working Group." *National Institute of Justice*: NCJ 183697.
- Darwin, C. (1859): *The Origin of Species*. London: Murray.
- Darwin, C. (1871): *The Descent of Man and Selection in Relation to Sex*. London: Murray.
- Dawid, P. (2002): "Bayes's Theorem and Weighing Evidence by Juries." In R. Swinburne (ed.), *Bayes's Theorem*. Oxford: Oxford University Press, pp. 71-90.

- Diaconis, P. and Mosteller, F. (1989): "Methods of Studying Coincidences." *J. Amer. Statist. Assoc.* 84: 853-861.
- Edwards, A. (1972): *Likelihood*. Cambridge: Cambridge University Press.
- Fisher, R.A. (1959): *Statistical Methods and Scientific Inference*. New York: Hafner, 2nd edition.
- Fitelson, B. (1999): "The Plurality of Bayesian Measures of Confirmation and the Problem of Measure Sensitivity," *Philosophy of Science* 66: S362-S378.
- Forster, M. and Sober, E. (1994): "How to Tell When Simpler, More Unified, or Less *Ad Hoc* Theories Will Provide More Accurate Predictions." *British Journal for the Philosophy of Science* 45: 1-36.
- Gilovich, T., Valone, R., and Tversky, A. (1985): "The Hot Hand in Basketball – On the Misperception of Random Sequences." *Cognitive Psychology* 17: 295-314.
- Hacking, I. (1965): *The Logic of Statistical Inference*. Cambridge: Cambridge University Press.
- Hitchcock, C. and Sober, E. (2004): "Prediction versus Accommodation and the Risk of Overfitting," *British Journal for the Philosophy of Science* 55: 1-34.
- Kahnemann, D., Slovic, P. Tversky, A. (1982): *Judgment under Uncertainty – Heuristics and Biases*. Cambridge: Cambridge University Press.
- Klarreich, E. (2004): "Toss Out the Toss-Up – Bias in Heads-or-Tails." *Science News Online* 165: 131. Available at <http://www.sciencenews.org/articles/20040228/fob2.asp>.
- Littlewood, J. (1953): *A Mathematician's Miscellany*. London: Methuen.
- Myers, David G. (2002): *Intuition – Its Powers and Perils*. New Haven: Yale University Press.
- Pearl, J. (2000): *Causality – Models, Reasoning, Inference*. New York: Cambridge University Press.
- Polanyi, M. (1966): *The Tacit Dimension*. New York: Doubleday.
- Reichenbach, H. (1956): *The Direction of Time*. Berkeley: University of California Press.
- Royall, R. (1997): *Statistical Evidence -- a Likelihood Paradigm*. London: Chapman and Hall.
- Sober, E. (1988a): "The Principle of the Common Cause." In J. Fetzer (ed.) *Probability and Causation: Essays in Honor of Wesley Salmon*. Dordrecht: Reidel, 211-28; Reprinted in *From a Biological Point of View*. New York: Cambridge University Press, 1994.
- Sober, E. (1988b): *Reconstructing the Past – Parsimony, Evolution, and Inference*. Cambridge: MIT Press.

- Sober, E. (2001): "Venetian Sea Levels, British Bread Prices, and the Principle of the Common Cause." *British Journal for the Philosophy of Science* 52: 1-16.
- Sober, E. (2002): "Instrumentalism, Parsimony, and the Akaike Framework." *Philosophy of Science* 69: S112-S123.
- Sober, E. (2003): "Two Uses of Unification." In F. Stadler (ed.), *The Vienna Circle and Logical Empiricism -- Vienna Circle Institute Yearbook 2002*. Kluwer, pp. 205-216.
- Sober, E. (2004a): "The Contest Between Likelihood and Parsimony." *Systematic Zoology* xx: xxx-xxx.
- Sober, E. (2004b): "The Design Argument." In W. Mann (ed.), *Blackwell Guide to the Philosophy of Religion*. Oxford: Blackwell.
- Spirtes, P. Glymour, C. and Scheines, R. (2001): *Causality, Prediction, and Search*. Cambridge: MIT Press.
- Stone, M. (1977): "An Asymptotic Equivalence of Choice of Model by Cross-Validation and Akaike's Criterion." *Journal of the Royal Statistical Society B* 39: 44-47.
- Woodward, J. (2003): *Making Things Happen*. Oxford: Oxford University Press.