



UNIVERSITY OF GOTHENBURG  
SCHOOL OF BUSINESS, ECONOMICS AND LAW

**WORKING PAPERS IN ECONOMICS**

**No 594**

**Confirmation: What's in the evidence?**

**Mitesh Kataria**

**(Maj 2014)  
Revised June 2015**

**ISSN 1403-2473 (print)  
ISSN 1403-2465 (online)**

# Confirmation: What's in the evidence?

Mitesh Kataria<sup>†</sup>

## Abstract

The difference between accommodated evidence (i.e. when evidence is known first and a hypothesis is proposed to explain and fit the observations) and predicted evidence (i.e., when evidence verifies the prediction of a hypothesis formulated before observing the evidence) is investigated. According to Bayesian confirmation theory, accommodated and predicted evidence constitute equally strong confirmation. Using a survey experiment on a sample of students, however, it is shown that predicted evidence is perceived to constitute stronger confirmation than accommodated evidence and in line with the decision analytical framework that is presented we show that predictions work as a signal about the scientists' knowledge which in turn provides stronger confirmation. The existence of such an indirect relationship between hypothesis and evidence can be considered to impose undesirable subjectivity and arbitrariness on questions of evidential support. Evidential support is ideally a direct and impersonal relationship between hypothesis and evidence and not an indirect and personal relationship as it has shown to be in this paper.

*Keywords:* Subjective beliefs, Evidence, Prediction, Postdiction, Retrodiction

*JEL classification:* C11, C12, C80

---

<sup>†</sup> Department of Economics, University of Gothenburg, Box 640, 40530 Gothenburg, Sweden. Tel: +46 317861352. E-mail: [mitesh.kataria@economics.gu.se](mailto:mitesh.kataria@economics.gu.se)

## 1. Introduction

Is it relevant for scientific confirmation whether evidence is known first and a hypothesis is thereafter proposed to explain and fit it (henceforth called accommodation), or whether evidence verifies predictions from a hypothesis formulated before observing the evidence?<sup>1</sup> Bayesian confirmation theory implies that whether evidence is predicted or accommodated is irrelevant. A hypothesis is confirmed if the posterior is greater than the prior, and this occurs if evidence supports the hypothesis independent of the timing of the empirical claim, i.e. whether the hypothesis is stated before or after the data has been observed. This view, also known as the purely logical approach to confirmation, holds that the timing of the empirical claim is irrelevant for scientific confirmation.

Musgrave (1974) discusses the possibility of how to move away from the purely logical approach to confirmation by a detailed review of three different views of the historical approach. The historical approach takes into account the historical setting in which that the theory was proposed when judging the evidence. It holds that predicted evidence is more important than accommodated evidence unless special circumstances prevail. The first and most extreme, the strictly temporal view of background knowledge, holds that facts known before a hypothesis is proposed cannot confirm the hypothesis since it is already part of background knowledge. But for many this view is too conservative. A famous historical example is that Einstein showed that general relativity agrees closely with the observed amount of perihelion shift, which was not the case with Newtonian physics. Although the motion of the perihelion of Mercury was known long before Einstein proposed his theory, the evidence was considered to support the theory and to be a powerful argument motivating the adoption of general relativity. The second, heuristic view of background knowledge, claims that an old fact can confirm and be novel to a new theory, provided the theory has not been constructed to explain the fact but is still in the process of explaining it. Finally, the third view of background knowledge holds that a new theory is independently testable or novel where old facts can confirm the new theory if and only if its prediction is unique such that it cannot be explained by the old theory or contradict it.

While scientific confirmation has been debated heavily in philosophy for over 400 years, it is more recently it gained some well-deserved attention in economics. Kahn et al. (1996) developed a model that focuses on different scientific methods. In their decision theoretical framework it is shown that if the scientist has different abilities to propose truthful

---

<sup>1</sup> In the case of prediction the hypothesis is usually partly based on existing observations, however, a prediction requires the empirical claim to be verified by at least some observations that are made after the empirical claim Lipton (2005).

theories and can choose either to predict the evidence or construct a theory that accommodates it, an observer will have a stronger belief in the truthfulness of the theory, if the theory is proposed before the evidence has been considered. The observer, assumingly unaware of the scientist's abilities, updates the probability that the consistent theory is proposed by a scientist with greater ability to propose truthful theories if evidence supports the theory, thus providing stronger confirmation. If the scientist constructs a theory that accommodates evidence, however, nothing is learned about the scientist's type, and no updating takes place. While Bayesian epistemology traditionally avoids the relationship between evidence and personal or psychological attributes, which is considered to impose undesirable subjectivity and arbitrariness on questions of evidential support, Kahn et al. (1996) focus on this link.<sup>2</sup>

In this paper, we first test whether people take a purely logical stand on confirmation (i.e. holding the belief that timing of the empirical claim is irrelevant for scientific confirmation) or whether they believe in a research hypothesis that predicts evidence more than in one that accommodates it. Second, we test whether prediction constitutes stronger confirmation than accommodation because the observer infers that the scientist is more knowledgeable when the scientist provides a correct prediction (e.g., Kahn et al. 1996). Hence, we investigate whether it is the fit between the hypothesis and the evidence that effects belief concerning whether or not a hypothesis is true, or if people weight in beliefs considering the mental ability of the scientist when judging the evidence. As far as we know, there are no empirical studies that address these questions. Yet, it seems essential for our understanding of how scientific results are perceived. This is true not only for results produced in economics but in general for all empirical sciences.<sup>3</sup> To that extent economist wants to inform policymakers and improve evidence-based policymaking the understanding of how people actually judge evidence is crucial. In this paper we offer insights that could be useful to understand perception of evidence-based policymaking. Our main findings are that predicted evidence constitutes stronger confirmation than accommodated evidence and that prediction works as a convincing signal about the scientists' knowledge.

---

<sup>2</sup> A reason to avoid such a link is that norms in science value "universalism" which means that a person's attributes and social background is irrelevant to the scientific value of a person's ideas. I.e. scientific findings must be judged by impersonal criterias.

<sup>3</sup> Communication of scientific results to nonscientists is at times not only important but crucial. E.g., for the science of climate changes it is essential to understand how nonscientists judge scientific evidence. Though many climate science studies show evidence that a long-term change (to year 2100) in the average atmospheric temperature could occur, a non-negligible share of the public persists in distrusting the results. Whether short-term predictions can reduce the credibility gap between scientists and the public is a question that partly inspired this project.

The remainder of the paper is organized as follows. In section 2, simple models are discussed and based on the predictions of our models hypotheses are formulated. In section 3 the experimental design is explained, followed by the results in section 4. Section 5 concludes the paper.

## 2. A Direct and an Indirect Model of Confirmation

The belief a (Bayesian) agent assigns to some (binary) hypothesis ( $H$ ) given evidence ( $E$ ) is expressed by Bayes' theorem:

$$P(H|E) = \frac{P(E|H) \cdot P(H)}{P(E)} = \frac{P(E|H) \cdot P(H)}{P(E|H) \cdot P(H) + P(E|\neg H) \cdot P(\neg H)}$$

Dividing the numerator and denominator in the above equation with  $P(E|H)$  we have

$$P(H|E) = \frac{P(H)}{P(H) + \frac{P(E|\neg H)}{P(E|H)} \cdot P(\neg H)} \text{ stating that } P(H|E) = f(P(H), \frac{P(E|\neg H)}{P(E|H)}) \text{ where } f \text{ is an increasing}$$

function of the first argument which is the prior probability and decreasing function of the likelihood ratio  $\frac{P(E|\neg H)}{P(E|H)}$ . Hence, for a given value of the likelihood ratio, the posterior

probability  $P(H|E)$  increases with the prior. Furthermore, for a given value of the prior, the posterior probability of  $H$  is greater, the less probable  $E$  is relatively to  $\neg H$  compared to  $H$ .

Hence, the more likely it is that certain observation is made if (and only if) the hypothesis is true the stronger is the confirmation if the observation is actually made. For our purposes the important observation is that the Bayesian formula does not distinguish prediction from accommodation as there is nothing that suggests that the timing of the evidence is important relatively to the timing of the hypothesis. The reason is that the likelihood ratio is unaffected by the timing of the evidence. Stating it differently, Bayesian inference show by the means of conditional probabilities how a subjective degree of belief should rationally change to account for evidence, however, whether evidence is known before, or the hypothesis is formulated after knowing the evidence is irrelevant.

We now turn to a simple and direct model where accommodated evidence is not perceived differently from predicted evidence (see Kahn et al., 1996 for details). Assume that theories can be divided into four mutually exclusive sets. Constructing theories involves some randomness and is modeled as the selection of a ball from an urn, the balls representing possible theories. Set  $A$  consists of true theories and is drawn with probability  $p$ , set  $B$  consists of false theories that could be falsified in future but are consistent with current observations and is drawn with probability  $q$ . Set  $C$  consists of theories that are consistent

with all past observations but not with the latest observation(s) and is drawn with probability  $1 - p - q$ . Finally, set  $D$  consist of inconsistent and false theories. We will assume that the scientists avoid drawing from set  $D$  so that the urn contains no type  $D$  balls. A scientist who predicts evidence by constructing a theory before learning that the theory is consistent with past and new observations, has constructed a theory which is true with probability  $P(A|A \& B) = \frac{p}{(p+q)}$ . A scientist who considers the evidence before constructing theory avoids drawing from set  $C$  and set  $D$ , and probability that the theory is true is again  $P(A|A \& B) = \frac{p}{(p+q)}$ . Hence, the first main result is:

*“...the probability that a theory is true, conditional on its being consistent with all (old and new) observations is independent of the research strategy.”* (Kahn et al., 1996).

Note that in this simplistic model we do not specify a decision framework of when to reject a hypothesis from data that is subject to random variation in terms of sample variability. In a frequentist decision framework for data that is subjected to random variation, a hypothesis is a conjecture about the distribution of one or more random variables. A statistical hypothesis defines the rule of when to reject the conjecture. In a single hypothesis test, an acceptable maximum probability of committing a Type 1 error (i.e. rejecting the true null hypothesis) is defined which is known as the testwise alpha. This is in practice often compared to a p-value which states the probability under the null to observe the sampled or more extreme data. If the probability is sufficiently small (i.e.  $p\text{-value} < \alpha$ ) the data is considered to be too extreme to be explained as an outcome of chance and therefore viewed as evidence against the null hypothesis.<sup>4</sup> The first main result can in the frequentist decision framework be stated as the probability to observe the sampled or more extreme data under the null is independent of the scientists' behavior such as if the scientist considers the data or not before proposing the hypothesis. Since the probability is independent of the research strategy, the decision whether to reject the hypothesis or not is also independent of the research strategy i.e. whether the evidence is accommodated or predicted does not matter for whether the hypothesis is rejected or not as long as we are analyzing the same amount of hypotheses.<sup>5</sup>

---

<sup>4</sup> Bayesian inference is based on the probability that a hypothesis is true given data i.e.  $P(H|E)$ , while the frequentist inference is based on the inverse probability which is  $P(E|H)$ . Notably they are not the same e.g., assuming that the probability to die after falling from a high building is 1, does not mean that a person who's dead has fallen from a high building.

<sup>5</sup> Note that the focus of this paper is not on multiple testing but on the difference between accommodated and predicted evidence, which differs in the timing of the evidence. Accommodating data, could but does not

Let us continue to consider an indirect model of confirmation presupposing a link between evidence and personal attributes. The aim of this extension is to identify a situation when accommodated evidence is perceived differently from predicted evidence. Kahn et al. (1996) shows that if there are different types of researchers, where one type is more likely to construct true and consistent theories independent of the research strategy chosen, then the outcome of the prediction conveys information to an observer about the type of researcher. Instead of the selection of a ball from one urn, this could be thought of as a draw from one of two types of urns, one containing more true and consistent theories and the other containing more inconsistent theories. Expressing this formally, assume the following for the two types  $i$ , and  $j = 1 - i$ :  $P_i(A) + P_i(B) = p + q > P_j(A) + P_j(B) = r + s$ . As before we assume that the inconsistent theories are drawn with probability  $1 - P_i(A) - P_i(B)$  for scientists type  $i$  and  $1 - P_j(A) - P_j(B)$  for scientists type  $j$ . The posterior probability that the theory is suggested by the high ability type  $i$  if the scientist theorizes first and the theory is consistent is found using Bayes theorem:

$$P(i|A\&B) = i' = \frac{P(i)P(A\&B|i)}{P(A\&B)} = i \frac{p + q}{i(p + q) + j(r + s)} > i$$

Assume that the scientist who is more likely to construct true and consistent theories is also more likely to suggest true theories given they are consistent i.e.  $\frac{p}{p+q} > \frac{r}{r+s}$ . From an observer's perspective, the probability that a certain theory is true depends on beliefs about how likely the observer finds it that the theory is suggested by the two type of scientists. The more likely it is found that the theory is suggested by type  $i$  scientist the higher probability the observer will attach to proposition that the theory is true. Moreover, the probability  $\gamma'$  that the theory is true if the scientist theorize first and the theory survives the test is higher than the probability  $\gamma$  that the theory is true if the scientist constructs the theory after having considered the data since in this case no updating takes place.

$$\gamma' = i' \frac{p}{p + q} + j \frac{r}{r + s} > \gamma = i \frac{p}{p + q} + j \frac{r}{r + s}$$

We can now state the second main result:

---

necessary imply multiple testing. Similarly, predicted evidence could but does not necessary imply multiple testing. If several hypothesis are tested a familywise error rate should possibly taken into account insted of he traditional testwise alpha since multiple testing increases the probability to make a Type 1 error. This is true for accommodated as well as predicted evidence.

*“Assume that the scientist's type is unknown. Then the probability that the theory is true, conditional on its being consistent with all (old and new) observations, is higher if the scientist theorized first than if he looked first.”* (Kahn et al., 1996).

The second main result suggests that prediction constitutes stronger confirmation than accommodation because prediction allows the observer to update the belief of whether the theory is suggested by the low or high ability type. Absurdly it also implies that if there are two scientists who propose identical theories, the only difference being that one scientist is more careful in observing a full sample of events while the other, less observant scientist neglects a subsample of the events which instead are predicted, the more observant scientist will be penalized for being knowledgeable about the occurred events and mistrusted (see, e.g., Musgrave 1974). The mistrust toward the scientist who observes more could be rooted in the belief that more observations increase the risk of over-fitting the data (e.g., Hitchcock and Sober 2004; Lipton 2005), i.e., explaining patterns that have appeared by chance..<sup>6</sup>

Our survey experiment is designed to test two hypotheses related to the models of confirmation we discussed. The first null hypothesis is:

*H0: People believe that the probability that a theory is true is independent of whether evidence is predicted or accommodated, i.e., the research strategy of the scientist.*

In case the first null hypothesis is rejected, we want to test if people trust scientists who theorize first to be more knowledgeable. The second null hypothesis is:

*H0: People do not believe that a scientist who makes a correct prediction is more knowledgeable than a scientist who accommodates evidence.*

---

<sup>6</sup> One can argue that the literature takes an asymmetric stand on this issue as it seldom discusses the risk of under-fitting the data (i.e. the practice of missing robust structural relationships by insisting on a-priorism) while the phenomena of over-fitting have gained much more attention. An exemption is found in a famous book chapter on how to write empirical papers. Bem (2003) provides the following advice to students in psychology: “There are two possible articles you can write: (a) the article you planned to write when you designed your study or (b) the article that makes the most sense now that you have seen the results. They are rarely the same, and the correct answer is (b).”. Bem is more concerned with the problems of under-fitting the data and missing the chance to discover than with the problems of over-fitting the data.



### 3. Experimental Survey Design

The data for this study was collected using an online experimental survey with three treatments conducted in Jena, Germany, in May 2013. The subjects were students with various educational backgrounds and were recruited using the Online Recruitment System for Economic Experiments (ORSEE). Participation was incentivized by a lottery procedure with a 10 percent probability to win 25 euros. In total, 243 (=81x3) subjects were recruited using between-subject design (i.e. each subject only participates in one of the treatments) . The treatments are summarized in Table 1, followed by the scenarios in T1 and T2. Unique features of the scenario in T1 are marked with curly brackets { } while unique features of T2 are marked with square brackets [ ]. The scenario in T3 is presented in the Appendix. Common for all treatments are that they contain a scenario about El Niño and questions about the predictability of El Niño. El Niño refers to the extensive warming of the central and eastern tropical Pacific that leads to extreme climate change patterns across the Pacific.<sup>7</sup>

In the baseline treatment, T1, subjects are presented with a scenario where a research team accommodates data before constructing the model that aims at predicting future El Niño events. In the second treatment, T2, subjects are presented with a scenario where a research team accommodates data before constructing the model, which results in five predictions that are subsequently shown to be correct. The model of the scenario in the baseline treatment will henceforth be referred as the *A-model* and the model of the scenario in T2 will henceforth be referred as the *P-model*. Both the A-model and the P-model share the feature that they are partly based on existing observations. The main difference is that the P-model is verified by at least some observations that are made after constructing the model. In both treatments, subjects are asked about how likely they think it is that the model they are presented with will be correct in future predictions. In the third treatment, T3, subjects are presented with a scenario where two teams of scientists each employ one of the two research strategies to “accommodate” or “predict” data. The main aim of this treatment was to ask subjects whether they believe that either of the two teams is more knowledgeable. The subjects’ were also asked about how likely they think it is that the predictions of the two teams will be correct in future trails.

---

<sup>7</sup> For more information about El Niño see e.g. Cane et al. (1997).

**Table 1: Treatments**

| Treatment | No. Obs. | Description                                 |
|-----------|----------|---|
| T1        | 81       | A scenario with no predictions              |
| T2        | 81       | A scenario with five successful predictions |
| T3        | 81       | A joint scenario                            |

**Scenario in {T1} and [T2]**

El Niño is characterized of abnormal warm ocean water temperatures that occasionally develops off the western coast of South America and can cause climatic changes across the Pacific Ocean. El Niño events happen irregularly.

The El Niño phenomenon drastically affects the weather in many parts of the world. Developing countries which are dependent on agriculture and fishing, particularly those bordering the Pacific Ocean, are the most affected. If forecasts could provide warnings before an El Niño episode, human suffering and economic losses could be reduced. Please consider the partly fictional scenario below and answer the question.

A team of scientists have recently constructed a new hybrid model, where an ocean model is coupled to a statistical atmospheric model that accommodates {25} [20] of the known El Niño events of the twentieth century (i.e., 1901 – 2000). The model is constructed to fit observations that have already been made. Using old data (i.e., the known El Niño events) the model is rigorously tested and able to detect El Niño events 12 months before it starts in {20} [15] of the {25} [20] cases without causing any false alarms. Without any prior knowledge, the chances to detect El Niño before it starts are only 5 percent. A model is considered to be good, if it detects El Niño events 12 months before they start with a probability of 80 percent. {The model has never been tested regarding how well it predicts future El Niño events but has rigorously been tested using old data.} [The model has recently also been tested on how well it predicts future El Niño events. More specifically, after the model was developed, the El Niño event has occurred 5 times, and the model successfully predicted these events 12 months before they started in all of these 5 cases without causing any false alarms. In total that would imply 15 correct predictions using old data and 5 correct predictions using new data, i.e., a total of 20 correct predictions out of 25.]

At this stage, a few remarks are necessary about the scenarios before we discuss the results. First of all, note that the performance of the two models in the two scenarios is the same in terms of the ability to detect the El Niño events with a probability of 80 percent. Also note that the total number of events is kept the same in the two scenarios to ensure comparability. The main difference between the two scenarios is the amount of evidence that is accommodated and predicted. In T1, the scenario consist more confirmation of accommodated evidence while in T2 the scenario consist more of predicted evidence.<sup>8</sup> The task of the subjects in the experiment is to state their beliefs that the model(s) in the scenarios will make correct prediction(s), given the evidence they have about the historical performance of the model(s) and how the model(s) were developed. In particular, we are interested if their beliefs differ depending on whether evidence is accommodated or predicted. More specifically, to interpret the responses to the scenarios in a Bayesian framework it is assumed

<sup>8</sup> In our scenarios a scientific hypothesis is represented by a model. While a model can be used to represent a scientific hypothesis, these terms should in general not be used interchangeably. Furthermore, while a hypothesis can be true or false, this terminology seems less appropriate for discussing models, and we therefore talk about models being good or bad.

that the subjects have prior beliefs whether a certain model is good or bad. Given the parameterization of our scenarios and assuming the existence of two types of models, a good model is defined to make 80 percent correct predictions and a bad model to make 5 percent correct predictions, subjects are presented with historical evidence of whether the model produces correct predictions. Based on the evidence, subjects could update their posterior beliefs about whether the model is good (positively if confirming evidence) or bad (positively if disconfirming evidence) using Bayes' theorem of binary hypotheses. Subjects are then asked how likely they think that the model they been presented with will be successful in predicting future El Niño events, which in turn depends on whether they believe that the model is good or bad and should be sufficient to address our research question whether accommodated evidence is treated differently than predicted evidence. Alternatively, we could have asked the subjects more directly about their posterior beliefs that the model is good. The reason we instead asked them indirectly about their beliefs about the predictability of the model(s) was that it facilitated a classical frequentist interpretation of the scenarios. The aim was to present the subjects with a short but content-rich and meaningful scenario without nudging them toward applying either the classical frequentist or Bayesian framework.

Let us now turn to the question of how a frequentist might react to the scenarios. For a frequentist the probability for the El Niño event to be detected in the future equals the relative frequency of detection to occurrence of the event in the past. Hence, given that subjects weight accommodated and predicted evidence equally, the probability that the prediction of the two models will be correct in future outcomes should be the same in the two scenarios.

## 4. Results

In this section, the three main results of the paper will be presented. The first result utilizes the single scenario data to test the difference between accommodated and predicted evidence, the second result compares the result of the single scenario data with the joint scenarios data, and the third result focuses on the joint scenarios data to test for subjective (psychological) links between the evidence and the proposer of the evidence.

**Result 1:** *A model that predicts evidence is assessed to be more correct than a model that accommodates evidence.*

The probability that the P-model will be correct in future predictions is appreciated to 77 percent while the A-model is appreciated to make correct predictions with a probability of

65 percent. The difference is statistically significant for any conventional significance level using a two-sided Mann-Whitney-Wilcoxon test.

**Result 2:** *Trust in a model that accommodates evidence increases when the model is compared (side-by-side) to a model that predicts evidence.*

The probability that the P-model will be correct in future predictions is appreciated to 76 percent in the joint scenario T3, which is not significantly different to 77 percent in T2, the single scenario data, using a two-sided Mann-Whitney-Wilcoxon test. The A-model in T3, however, is appreciated to make correct predictions in the future with a probability of 72 percent. The increase from the 65 percent in T1 is statistically significant using a two-sided Mann-Whitney-Wilcoxon test (p-value: 0,030). Hence people take a more logical approach (i.e. trust accommodated and predicted evidence more equally) when both models are compared side-by-side.<sup>9</sup> The increase notwithstanding, a model that predicts evidence is still assessed to be more correct than a model that accommodates evidence. The difference is statistically significant for any conventional significance level using a two-sided signed-rank test. This means that the effect that prediction constitutes stronger confirmation is robust to the framing of the problem.

**Result 3:** *Trust in a model's capacity to predict seems to be intrinsically related to trust in the scientists' level of knowledge.*

Note that most subjects (60%) do not believe that a scientist who predicts evidence is more knowledgeable than a scientist who accommodates evidence or the other way around. However, a substantial share (32%) of subjects believes that a scientist who predicts evidence is more knowledgeable. Only a small share of subjects (8%) believes that a scientist who accommodates evidence is more knowledgeable.

Subjects that believe that a scientist who accommodates evidence is equally knowledgeable to a scientist who predicts evidence, trust the P-model (which uses predicted evidence) to be correct in future predictions with a probability of 77 percent and the A-model (which uses accommodated evidence) to be correct with a probability of 74 percent. This

---

<sup>9</sup> Forty-nine (49)% of the subjects stated that the A-model and P-model will be equally correct in future prediction(s), while 37% stated a higher probability for the P-model and 14% stated a higher probability for the A-model.

small difference is not statistically significant for any conventional significance level using the Wilcoxon signed-rank test.

Subjects that believe that a scientist who predicts evidence is more knowledgeable than a scientist who accommodates evidence, trust the P-model to be correct in future predictions with a probability of 77 percent and the A-model to be correct with a probability of 65 percent. This difference is statistically significant for any conventional significance level using the Wilcoxon signed-rank test.

The third result offers interesting insights. In the philosophical literature many attempts has been done to revise the purely logical approach to confirmation and to show that predictions constitute stronger confirmation than accommodation. The aim has been to formulate a theory without introducing undesirable subjectivity in the relationship between hypothesis and evidence. However, no theory has to date been widely accepted. Interestingly, our results show that the intuition that predicted evidence constitute stronger confirmation is driven by subjects' that not only judge the relationship between the hypothesis and the evidence but that uses the evidence to infer something about the abilities of the scientist which in turn provides stronger confirmation.

## **5. Conclusions**

In this paper, we turn to folk intuitions to empirically investigate whether predicted evidence constitutes stronger confirmation than accommodated evidence. Two key findings have emerged. First, we find that predicted evidence constitutes stronger confirmation than accommodated evidence, and the effect is robust to the framing of the problem. While most subjects (49%) believe that a model that accommodates evidence is as good as one that has been shown to predict evidence, a substantial share (37%) of subjects believe that a model that has been shown to predict evidence is more likely to perform better in future trials. Only a small share of subjects (14%) had a higher trust in the model that accommodated evidence. Hence it seems that people perceive the risk of over-fitting the data as more severe than under-fitting the data. Second, our results suggests that trusting a model to predict correctly is intrinsically related to trust in the proposers' (i.e., the scientists') level of knowledge, and relatively more subjects are persuaded by proposers' abilities to utilize this knowledge to predict in the future if they, the proposers, in the past shown to be successful in predicting rather than accommodating evidence. We find that prediction works as a convincing signal about the scientists' knowledge. This confirms the conjecture in Kahn et al. (1996) and Lipton (2005) that evidential support is not a direct relationship between a hypothesis and evidence

but indirectly linked through the proposer of the model. Notably, this link can be argued to be epistemologically unwarranted.

Accommodating the data to find the model that is considered to best explain the data might make some readers to think about the econometric literature of data mining (Leamer, 1983) and data snooping (White, 2000) which deals with the problems involved in conducting an extensive specification search of a model. Therefore this literature deserves a comment especially on how it relates to the scope of the study presented here. One way to articulate the problem of data mining is that in any specification search there will be a multiple amount of hypotheses tested, while the tools that econometricians traditionally apply are developed for testing a single hypothesis. White (2000) notes that a method controlling for multiple hypotheses: "... permits data snooping/mining to be undertaken with some degree of confidence that one will not mistake results that could have been generated by chance for genuinely good results".<sup>10</sup> What is important to remember, however, is that although accommodating evidence is related to multiple testing, it is not the distinguish feature between accommodated and predicted evidence. The distinguish feature between accommodated and predicted evidence is the timing of the evidence, which also been the focus of this study, which means that we are interested in explaining the differences between accommodated and predicted evidence everything else (including number of hypotheses tested) equal.

The most interesting result of our paper is that we showed the existence of an indirect relationship between hypothesis and evidence. This specific relationship can in turn be argued to be epistemologically unwarranted as it violates important scientific norms such universalism. On the other hand, such indirect relationship is in line with the academic environment where hierarchical ranking structure amongst scholars is very clear and could carry information about researchers' level of knowledge, which in turn can affect trust in the claims that a researcher makes.

---

<sup>10</sup> But far from everyone seems to agree. Westfall et al. (1997) notes: "Multiple testing is difficult and controversial on either side of the Bayes'/frequentist fence, with arguments over whether and how multiplicity adjustment should be performed." They also discuss when and how to adjust prior assessments to account for multiplicity and specify conditions for which the resulting posterior probabilities roughly correspond to Bonferroni adjusted p-values. Berry and Hochberg (1999) also discuss Bayesian attitudes and methods for adjusting inferences for multiplicities. One objection is that multiplicity adjustments are too subjective; whether the correction should be made depending on the number of tests per article or for all tests considered in the scientists' lifetime, or whether each field should correct the number of tests or whether any corrections should be made at all – all of them questions that have puzzled many practicing statisticians. Hoover and Perez (2000) confirm the view that econometricians have highly different attitudes toward data mining, which range from 'it is to be avoided' to 'it is inevitable' or even to 'it is essential'.

## References

- Berry, D. A., & Hochberg, Y. (1999). Bayesian perspectives on multiple comparisons. *Journal of Statistical Planning and Inference*, 82(1), 215-227.
- Bem, D. J. (2003). Writing the empirical journal article. In J. M. Darley, M. P. Zanna, & H. L. Roediger III (Eds.), *The complete academic: A career guide* (pp. 171–201). Washington, DC: American Psychological Association.
- Cane, M. A., Zebiak, S. E. and Dolan, S. C. (1986). Experimental forecasts of El Niño. *Nature*, **321**, 827–832.
- Hitchcock, C and E. Sober (2004). Prediction versus accommodation and the risk of overfitting. *Brit. J. Philos. Sci.*, 55.
- Hoover, K. D and Perez, S. J. (2000). Three attitudes towards data mining. *Journal of Economic Methodology*, 7:2, 195-210.
- Kahn, J.A., Landsburg, S.E., and Stockman, A.C. (1996). The Positive Economics of Methodology. *Journal of Economic Theory*, 68(1), 64–76.
- Leamer, E. (1983). Let's Take the Con Out of Econometrics. *American Economic Review*, 73: 31-43.
- Lipton, P. (2005). Testing hypotheses: prediction and prejudice. *Science*, 307, 219-221.
- Maddala G.S. (1992). Introduction to Econometrics. 2nd ed., Macmillan.
- Musgrave, A. (1974). Logical versus historical theories of confirmation. *Brit. J. Philos. Sci.*, 22.
- Westfall, P.H., Johnson, W.O., Utts, J.M. (1997). A Bayesian perspective on the Bonferroni adjustment. *Biometrika*, 84, 419-427.
- White, H. (2000). A reality check for data snooping. *Econometrica*, 68, 1067–1084.

### Appendix: Scenario in Treatment 3

A team of scientists, henceforth called team A, have recently constructed a new hybrid model, where an ocean model is coupled to a statistical atmospheric model that accommodates 20 of the known El Niño events of the twentieth century (i.e., 1901 – 2000). The model is constructed to fit observations that have already been made. Using old data (i.e., the known El Niño events), the model is rigorously tested and able to detect El Niño events 12 months before they start in 15 of the 20 cases without causing any false alarms. Without any knowledge, the chance to detect El Niño before it starts is only 5 percent. A model is considered to be good if it detects El Niño events 12 months before they start with a probability of 80 percent. The model has recently also been tested on how well it predicts future El Niño events. More specifically, after the model was developed, the El Niño event has occurred 5 times, and the model successfully predicted that event 12 months before it started in all of these 5 cases without causing any false alarms. In total, that would imply 20 correct predictions using old data and 5 correct predictions using new data, i.e., in total 20 correct predictions out of 25.

Independent of team A's work, but knowing that the latest five El Niño events occurred, another group of scientists, henceforth called team B, developed a different hybrid model (i.e., where an ocean model is coupled to a statistical atmospheric model) which is constructed to accommodate the occurred events. The model is constructed to fit observations that have already been made. Testing the model on old data (including the 5 latest El Niño events), team B's model is able to detect El Niño events 12 months before they start in 20 of the 25 cases without causing any false alarms. Team B's model has never been tested on how well it predicts future El Niño events but has been rigorously tested using old data.

The two teams disagree about which model is more correct and will have greater predictive power in the future. Please state a) which of the two teams you consider to be more knowledgeable based on the information you have been given and b) how likely you believe it is that the models will be correct in future predictions.