

## Beyond subjective and objective in statistics

Andrew Gelman

*Columbia University, New York, USA*

and Christian Hennig

*University College London, UK*

[Read before The Royal Statistical Society on Wednesday, April 12th, 2017, Professor P. J. Diggle in the Chair]

**Summary.** Decisions in statistical data analysis are often justified, criticized, or avoided by using concepts of objectivity and subjectivity. We argue that the words ‘objective’ and ‘subjective’ in statistics discourse are used in a mostly unhelpful way, and we propose to replace each of them with broader collections of attributes, with objectivity replaced by *transparency*, *consensus*, *impartiality* and *correspondence to observable reality*, and subjectivity replaced by awareness of *multiple perspectives* and *context dependence*. Together with *stability*, these make up a collection of virtues that we think is helpful in discussions of statistical foundations and practice. The advantage of these reformulations is that the replacement terms do not oppose each other and that they give more specific guidance about what statistical science strives to achieve. Instead of debating over whether a given statistical method is subjective or objective (or normatively debating the relative merits of subjectivity and objectivity in statistical practice), we can recognize desirable attributes such as transparency and acknowledgement of multiple perspectives as complementary goals. We demonstrate the implications of our proposal with recent applied examples from pharmacology, election polling and socio-economic stratification. The aim of the paper is to push users and developers of statistical methods towards more effective use of diverse sources of information and more open acknowledgement of assumptions and goals.

*Keywords:*

### 1. Introduction

We cannot do statistics without data, and as statisticians much of our efforts revolve around modelling the links between data and substantive constructs of interest. We might analyse national survey data on purchasing decisions as a way of estimating consumers’ responses to economic conditions, or gather blood samples over time on a sample of patients with the goal of estimating the metabolism of a drug, with the ultimate goal of coming up with a more effective dosing schedule; or we might be performing a more exploratory analysis, seeking clusters in a multivariate data set with the aim of discovering patterns that are not apparent in simple averages of raw data.

As applied researchers we are continually reminded of the value of integrating new data into an analysis, and the balance between data quality and quantity. In some settings it is possible to answer questions of interest by using a single clean data set, but increasingly we are finding that this simple textbook approach does not work.

*Address for correspondence:* Andrew Gelman, Department of Statistics, Columbia University, Room 1016, 1255 Amsterdam Avenue, New York, NY 10027, USA.  
E-mail: [gelman@stat.columbia.edu](mailto:gelman@stat.columbia.edu)

1 External information can come in many forms, including

- 2 (a) recommendations on what variables to adjust for non-representativeness of a survey or
- 3 imbalance in an experiment or observational study,
- 4 (b) the extent to which outliers should be treated as regular, or erroneous, or as indicating
- 5 something that is meaningful but essentially different from the main body of observations,
- 6 (c) issues of measurement, confounding, and substantively meaningful effect sizes,
- 7 (d) population distributions that are used in poststratification, age adjustment and other
- 8 procedures that attempt to align inferences to a common population of interest,
- 9 (e) restrictions such as smoothness or sparsity that serve to regularize estimates in high
- 10 dimensional settings,
- 11 (f) the choice of the functional form in a regression model (which in economics might be
- 12 chosen to work with a particular utility function, or in public health might be motivated
- 13 on the basis success in similar studies in the literature) and
- 14 (g) numerical information about particular parameters in a model.

15  
16 Of all these, only the final item is traditionally given the name ‘prior information’ in a statistical  
17 analysis, but all can be useful in serious applied work. Other relevant information concerns not  
18 the data-generating process but rather how the data and results of an analysis are to be used or  
19 interpreted.

20 We were motivated to write the present paper because we felt that our applied work, and that  
21 of others, was impeded because of the conventional framing of certain statistical analyses as  
22 subjective. It seemed to us that, rather than being in opposition, subjectivity and objectivity both  
23 had virtues that were relevant in making decisions about statistical analyses. We have earlier  
24 noted (Gelman and O’Rourke, 2015) that statisticians typically choose their procedures on the  
25 basis of non-statistical criteria, and philosophical traditions and even the labels attached to  
26 particular concepts can affect real world practice.

27 In Section 2 we show how the discussions of objectivity and subjectivity affect statistical  
28 practice and statistical thinking, followed by an outline of our own philosophical attitude to  
29 objectivity and subjectivity in science; Appendix A provides an overview of what the philosophy  
30 of science has to say on the matter. In Section 3 we present our proposal, exploding objectiv-  
31 ity and subjectivity into several specific virtues to guide statistical practice. In Section 4 we  
32 demonstrate the relevance of these ideas for three of our active applied research projects: a  
33 hierarchical population model in pharmacology, a procedure for adjustment of opt-in surveys  
34 and a cluster analysis of data on socio-economic stratification. In Section 5 we revisit funda-  
35 mental approaches to statistics by using the proposals of Section 3, demonstrating how they  
36 can elaborate advantages and disadvantages of the various approaches in a more helpful way  
37 than the traditional labelling of ‘objective’ and ‘subjective.’ Section 6 contains a final discussion,  
38 including a list of issues in scientific publications that could be addressed by using the virtues  
39 that are proposed here.

## 40 **2. Objectivity and subjectivity**

### 41 *2.1. Objectivity, subjectivity and decision making in statistics*

42 Concepts of objectivity and subjectivity are often used to justify, criticize, avoid or hide the deci-  
43 sions that are made in data analysis. Typically, the underlying idea is that science should be ob-  
44 jective, understood as something like ‘independence of personal biases’, without referring to any  
45 clear definition. Concepts of objectivity can be implicitly invoked when making choices, so that  
46 certain decisions are avoided or hidden in order not to open an analysis to charges of subjectivity.  
47  
48

1 For many statistical methods tuning constants need to be decided such as the proportion  
2 of trimmed observations when computing a trimmed mean or bandwidths for smoothing in  
3 density estimation or non-parametric regression; one could also interpret the conventional  
4 use of the 0.05 level of significance as a kind of tuning parameter. In the statistical literature,  
5 methods are advertised by stating that they do not require any tuning decisions by the user.  
6 Often these choices are hidden so that users of statistics (particularly those without specific  
7 background knowledge in statistics) expect that for their data analysis task there is a unique  
8 correct statistical method. This expectation is exploited by the marketing strategies for some  
9 data analysis software packages such as Beyondcore that suggest that a default analysis is only  
10 one click away. Although such an approach obviously tempts the user by its simplicity, it also  
11 appeals on the level of avoiding individual impact or subjectivity. Decisions that need to be made  
12 are taken out of the hand of the user and are made by the algorithm, removing an opportunity  
13 for manipulation but ignoring valuable information about the data and their background. This  
14 is in stark contrast with the typical trial-and-error way of building one or more statistical models  
15 with plenty of subjective decisions starting from data preprocessing via data exploration and  
16 choice of method onto the selection of how to present which results. Realistically, even one-  
17 click methods require user choices on data coding and data exclusion, and these inputs can  
18 have big influences on end results such as  $p$ -values and confidence intervals (Steege *et al.*,  
19 2006).

20 More mathematically oriented statisticians design and choose methods that optimize criteria  
21 such as unbiased minimum variance, often relying on restrictive assumptions. This can allow for  
22 elegant theoretical results with narrow scope. When methods are to be compared in simulation  
23 studies, typically there is a huge variety of choices including data-generating processes (distrib-  
24 tions, parameters, dimensionality, etc.), performance measures and tuning of competitors.  
25 This can easily discourage researchers from running such studies at all or at least beyond one  
26 or two illustrative toy set-ups but, despite their subjective flavour and the difficulty in finding a  
27 path through a potentially confusing jungle of results, such studies can be informative and raise  
28 important issues. Already in 1962, Tukey criticized an obsession with mathematical optimiza-  
29 tion problems concentrating on one-dimensional criteria and simple parametric models in the  
30 name of objectivity, and stated that ‘in data analysis we must look to a very heavy emphasis on  
31 judgment’ (Tukey, 1962).

32 Researchers often rely on the seeming objectivity of the  $p < 0.05$  criterion without realizing  
33 that theory behind the  $p$ -value is invalidated when analysis is contingent on data (Simmons  
34 *et al.*, 2011, Gelman and Loken, 2014). Significance testing can be part of a misguided ideology  
35 that leads researchers to hide, even from themselves, the iterative searching process by which a  
36 scientific theory is mapped into a statistical model or choice of data analysis (Box, 1983). More  
37 generally, overreaction to concerns about subjectivity can lead researchers to avoid incorporat-  
38 ing relevant and available information in their analyses and not to adapt analyses appropriately  
39 to their research questions and potential uses of their results.

40 Personal decision making cannot be avoided in statistical data analysis and, for want of  
41 approaches to justify such decisions, the pursuit of objectivity degenerates easily to a pursuit  
42 to merely *appear* objective. Scientists whose methods are branded as subjective have the awkward  
43 choice of either saying, ‘No, we are really objective’, or else embracing the subjective label  
44 and turning it into a principle, and the temptation is high to avoid this by hiding researcher  
45 degrees of freedom from the public unless they can be made to appear ‘objective’. Such attitudes  
46 about objectivity and subjectivity can be an obstacle to good practice in data analysis and its  
47 communication, and we believe that researchers can be guided in a better way by a list of more  
48 specific scientific virtues when choosing and justifying their approaches.

1 The continuing interest in and discussion of objectivity and subjectivity in statistics are,  
2 we believe, a necessary product of a fundamental tension in science: on one hand, scientific  
3 claims should be impersonal in the sense that a scientific argument should be understandable  
4 by anyone with the necessary training, not just by the person promulgating it, and it should  
5 be possible for scientific claims to be evaluated and tested by outsiders. On the other hand, a  
6 reality that is assumed to be objective in the sense of being independent of its observers is only  
7 accessible through observations that are made by observers and dependent on their perspectives;  
8 communication about the observations and the process of observation and measurement relies  
9 on language constructs. Thus objective and subjective elements arise in the practice of science,  
10 and similar considerations hold in statistics.

11 Within statistics, though, discourse on objectivity and subjectivity is at an impasse. Ideally  
12 these concepts would be part of a consideration of the role of different sorts of information and  
13 assumptions in statistical analysis, but instead they often seemed to be used in restrictive and  
14 misleading ways.

15 One problem is that the terms ‘objective’ and ‘subjective’ are loaded with so many associations  
16 and are often used in a mixed descriptive–normative way. For example, a statistical method that  
17 does not require the specification of any tuning parameters is objective in a descriptive sense (it  
18 does not require decisions by the individual scientist). Often this is presented as an advantage of  
19 the method without further discussion, implying objectivity as a norm, but the lack of flexibility  
20 that is caused by the impossibility of tuning can actually be a disadvantage (and indeed can  
21 lead to subjectivity at a different point in the analysis, when the analyst must make the decision  
22 of whether to use an autotuned approach in a setting where its inferences do not appear to  
23 make sense). The frequentist interpretation of probability is objective in the sense that it locates  
24 probabilities in an objective world that exists independently of the observer, but the definition  
25 of these probabilities requires a subjective definition of a reference set. Some proponents of  
26 frequentism consider its objectivity (in the sense of impersonality, conditional on the definition  
27 of the reference set) as a virtue, but this property is ultimately only descriptive; it does not  
28 imply on its own that such probabilities indeed exist in the objective world, nor that they are a  
29 worthwhile target for scientific inquiry.

30 In discussions of the foundations of statistics, objectivity and subjectivity are seen as opposites.  
31 Objectivity is typically seen as a good thing; many see it as a major requirement for good  
32 science. Bayesian statistics is often presented as being subjective because of the choice of a  
33 prior distribution. Some Bayesians (notably Jaynes (2003) and Berger (2006)) have advocated  
34 an objective approach, whereas others (notably de Finetti (1974)) have embraced subjectivity.  
35 It has been argued that the subjective–objective distinction is meaningless because all statistical  
36 methods, Bayesian or otherwise, require subjective choices, but the choice of prior distribution  
37 is sometimes held to be particularly subjective because, unlike the data model, it cannot be  
38 determined even in the asymptotic limit. In practice, subjective prior distributions often have  
39 well-known empirical problems such as overconfidence (Alpert and Raiffa, 1984; Erev *et al.*,  
40 1994), which motivates efforts to check and calibrate Bayesian models (Rubin, 1984; Little, 2012)  
41 and to situate Bayesian inference within an error statistical philosophy (Mayo, 1996; Gelman  
42 and Shalizi, 2013).

43 de Finetti can be credited with acknowledging honestly that subjective decisions cannot be  
44 avoided in statistics, but it is misleading to think that the required subjectivity always takes  
45 the form of prior belief. The confusion arises from two directions: first, prior distributions are  
46 not necessarily any more subjective than other aspects of a statistical model; indeed, in many  
47 applications priors can and are estimated from data frequencies (see chapter 1 of Gelman *et al.*  
48 (2013) for several examples). Second, somewhat arbitrary choices come into many aspects of

1 statistical models, Bayesian and otherwise, and therefore we think that it is a mistake to consider  
2 the prior distribution as the exclusive gate at which subjectivity enters a statistical procedure.

3 On one hand, statistics is sometimes said to be the science of defaults: most applications  
4 of statistics are performed by non-statisticians who adapt existing general methods to their  
5 particular problems, and much of the research within the field of statistics involves devising,  
6 evaluating and improving such generally applicable procedures (Gelman, 2014a). It is then seen  
7 as desirable that any required data analytic decisions or tuning are performed in an objective  
8 manner, either determined somehow from the data or justified by some kind of optimality  
9 argument.

10 On the other hand, practitioners must apply their subjective judgement in the choice of what  
11 method to use, what assumptions to invoke and what data to include in their analyses. Even  
12 using 'no need for tuning' as a criterion for method selection or prioritizing bias, for example,  
13 or mean squared error, is a subjective decision. Settings that appear completely mechanical  
14 involve choice: for example, if a researcher has a checklist saying to apply linear regression for  
15 continuous data, logistic regression for binary data, and Poisson regression for count data, he  
16 or she still has the option of coding a response as continuous or to use a threshold to define  
17 a binary classification. And such choices can be far from trivial; for example, when modelling  
18 elections or sports outcomes, one can simply predict the winner or instead predict the numerical  
19 point differential or vote margin. Modelling the binary outcome can be simpler to explain but in  
20 general will throw away information, and subjective judgement arises in deciding what to do in  
21 this sort of problem (Gelman, 2014b). And, in both classical and Bayesian statistics, subjective  
22 choices arise in defining the sample space and considering what information to condition on.  
23  
24

## 25 2.2. *Objectivity, subjectivity and quantification in scientific measurement*

26 Another issue that is connected to objectivity and subjectivity relevant to statisticians has to do  
27 with where the data to be analysed come from. There is an ideology that is widespread in many  
28 areas of science that sees quantification and numbers and their statistical analysis as key tools  
29 for objectivity. An important function of quantitative scientific measurement is the production  
30 of observations that are thought of as independent of individual points of view. But, even apart  
31 from the generally difficult issue of measurement validity, the focus on what can be quantified  
32 can narrow down what can be observed and may not necessarily do the measured entities justice.

33 The social sciences have seen endless arguments over the relative importance of objective  
34 conditions and what Keynes (1936) called 'animal spirits'. In macroeconomics, for example,  
35 the debate has been between the monetarists who tend to characterize recessions as necessary  
36 consequences of underlying economic conditions (as measured, for example, by current account  
37 balances, business investment and productivity), and the Keynesians who focus on more subjective  
38 factors such as stock market bubbles and firms' investment decisions. These disagreements  
39 also turn methodological, with much dispute, for example, over the virtues and defects of various  
40 attempts to measure objectively the supply and velocity of money, or consumer confidence or  
41 various other inputs to economic models. In psychology, there is a big effort to measure personality  
42 traits and subjective states scientifically. For example, Kahneman (1999) defined 'objective  
43 happiness' as 'the average of utility over a period of time'. Whether or not this definition makes  
44 much sense, it illustrates a movement in the social and behavioural sciences to measure, in supposedly  
45 objective manners, what might previously have been considered unmeasurable. Another  
46 example is the use of quantitative indicators for human rights in different countries; although it  
47 has been argued that it is of major importance that such indicators should be objective to have  
48 appropriate impact on political decision making (Candler *et al.*, 2011), many aspects of their

1 definition and methodology are subject to controversy and reflect specific political interests and  
2 views (Merry, 2011), and we think that it will help the debate to communicate such indicators  
3 transparently together with their limitations and the decisions involved rather than to sell them  
4 as objective and unquestionable.

5 Connected to quantification as a means of objectification is an attitude to statistics of many  
6 researchers in various areas who use standard routines in statistical software without much  
7 understanding of how the methods' assumptions and motivation relate to their specific research  
8 problem, in the expectation that the software can condense their research into a single summary  
9 (most often a  $p$ -value) that 'objectifies' their results. This idea of objectivity is in stark contrast  
10 with the realization by many of these researchers at some point that depending on individual  
11 inventiveness there are many ways to arrive at such a number.

12 See Porter (1996), Desrosieres (2002) and Douglas (2009) for more discussion of the connec-  
13 tion between quantification and objectivity. As with choices in statistical modelling and analysis,  
14 we believe that when considering measurement the objective–subjective antagonism is less help-  
15 ful than a more detailed discussion of what quantification can achieve and what its limitations  
16 are.

### 19 *2.3. Our attitude towards objectivity and subjectivity in science*

20 Many users of the terms 'objective' and 'subjective' in discussions concerning statistics do not  
21 acknowledge that these terms are quite controversial in the philosophy of science (as is 'realism')  
22 and that they are used with a variety of meanings and are therefore prone to misunderstandings.  
23 An overview is given in Appendix A.

24 The attitude that is taken in the present paper is based on Hennig (2010). According to  
25 this perspective, human inquiry starts from observations that are made by personal observers  
26 ('personal reality'). Through communication, people share observations and generate 'social  
27 realities' that go beyond a personal point of view. These shared realities include for example  
28 measurement procedures that standardize observations, and mathematical models that connect  
29 observations to an abstract formal system that is meant to create a thought system that is  
30 cleaned from individually different points of view. Nevertheless, human beings only have access  
31 to 'observer-independent reality' through personal observations and how these are brought  
32 together in social reality.

33 Science aims at arriving at a view of reality that is stable and reliable and can be agreed  
34 freely by general observers and is therefore as observer independent as possible. In this sense we  
35 see objectivity as a scientific ideal. But at the same time we acknowledge what gave rise to the  
36 criticism of objectivity: the existence of different individual perspectives and also of perspectives  
37 that differ between social systems, and therefore the ultimate inaccessibility of a reality that is  
38 truly independent of observers is a basic human condition. Objectivity can only be attributed  
39 by observers and, if observers disagree about what is objective, there is no privileged position  
40 from which this can be decided. Ideal objectivity can never be achieved.

41 How to resolve scientific disputes by scientific means without throwing up our hands and  
42 giving up on the possibility of scientific consensus is a key problem, and science should be  
43 guided by principles that at the same time aim at stable and reliable consensus as usually asso-  
44 ciated with 'objectivity' while remaining open to a variety of perspectives, often associated with  
45 'subjectivity', exchange between which is needed to build a stable and reliable scientific world  
46 view.

47 Although there is no objective access to observer-independent reality, we acknowledge that  
48 there is an almost universal human experience of a reality perceived as located outside the

1 observer and as not controllable by the observer. We see this reality as a target of science, which  
 2 makes observed reality a main guiding light for science. We are therefore ‘active scientific realists’  
 3 in the sense of Chang (2012), who wrote

4 ‘I take reality as whatever is not subject to one’s will, and knowledge as an ability to act without being  
 5 frustrated by resistance from reality. This perspective allows an optimistic rendition of the pessimistic  
 6 induction, which celebrates the fact that we can be successful in science without even knowing the truth.  
 7 The standard realist argument from success to truth is shown to be ill-defined and flawed.’

8 Or, more informally, ‘Reality is that which, when you stop believing in it, doesn’t go away’  
 9 (Dick, 1981). Active scientific realism implies that finding out the truth about objective reality is  
 10 not the ultimate aim of science, but that science rather aims at supporting human actions. This  
 11 means that scientific methodology must be assessed relatively to the specific aims and actions  
 12 that is connected to its use.

13 Because science aims at agreement, communication is central to science, as are transparency  
 14 and techniques for supporting the clarity of communication. Among these techniques are formal  
 15 and mathematical language, standardized measurement procedures and scientific models. Such  
 16 techniques provide a basis for scientific discussion and consensus, but at the same time the  
 17 scientific consensus should not be based on authority and it always needs to be open to new  
 18 points of view that challenge an established consensus. Therefore, in science there is always a  
 19 tension between the ideal of general agreement and the reality of heterogeneous perspectives,  
 20 and the virtues that are listed in Section 3 are meant to help statisticians navigating this tension.  
 21

### 22 3. Our proposal

23 To move the conversation towards principles of good science, we propose to replace, wherever  
 24 possible, the words ‘objectivity’ and ‘subjectivity’ with broader collections of attributes, namely  
 25 by *transparency*, *consensus*, *impartiality* and *correspondence to observable reality*, all related to  
 26 objectivity, awareness of *multiple perspectives* and *context dependence*, related to subjectivity,  
 27 and *investigation of stability*, related to both.  
 28

29 The advantage of this reformulation is that the replacement terms do not oppose each other.  
 30 Instead of debating over whether a given statistical method is subjective or objective (or nor-  
 31 matively debating the relative merits of subjectivity and objectivity in statistical practice), we  
 32 can recognize attributes such as transparency and acknowledgement of multiple perspectives as  
 33 complementary.  
 34

#### 35 3.1. ‘Transparency’, ‘consensus’, ‘impartiality’ and ‘correspondence to observable reality’, 36 instead of ‘objectivity’

37 Science is practised by human beings, who have access to the real world only through inter-  
 38 pretation of their perceptions. Taking objectivity seriously as an ideal, scientists need to make  
 39 the sharing of their perceptions and interpretations possible. When applied to statistics, the  
 40 implication is that choices in data analysis (including the prior distribution, if any, but also the  
 41 model for the data, methodology and the choice of what information to include in the first place)  
 42 should be motivated on the basis of factual, externally verifiable information and transparent  
 43 criteria. This is similar to the idea of the concept of ‘institutional decision analysis’ (section  
 44 9.5 of Gelman *et al.*, 2013), under which the mathematics of formal decision theory can be  
 45 used to ensure that decisions can be justified on the basis of clearly stated criteria. Different  
 46 stakeholders will disagree on decision criteria, and different scientists will differ on statistical  
 47 modelling decisions, so, in general, there is no unique ‘objective’ analysis, but we can aim at  
 48

1 communicating and justifying analyses in ways that support scrutiny and eventually consensus.  
2 Similar thoughts have motivated the slogan ‘transparency is the new objectivity’ in journalism  
3 (Weinberger, 2009).

4 In the context of statistical analysis, a key aspect of objectivity is therefore a process of *trans-*  
5 *parency*, in which the choices that are involved are justified on the basis of external, potentially  
6 verifiable sources or at least transparent considerations (ideally accompanied by sensitivity anal-  
7 yses if such considerations leave alternative options open), a sort of ‘paper trail’ leading from  
8 external information, through modelling assumptions and decisions about statistical analysis,  
9 all the way to inferences and decision recommendations. The current push of some journals to  
10 share data and computer code and the advent of tools to organize code and projects such as  
11 Github and version control better go in this direction. Transparency also comprises spelling out  
12 explicit and implicit assumptions about the data production, some of which may be unverifiable.

13 But transparency is not enough. Science aims at stable *consensus* in potentially free exchange  
14 (see Section 2.3), which is one reason that the current crisis of non-replication is taken so seriously  
15 in psychology (Yong, 2012). Transparency contributes to this building of consensus by allowing  
16 scholars to trace the sources and information that are used in statistical reasoning (Gelman and  
17 Basbøll, 2013). Furthermore, scientific consensus, as far as it deserves to be called ‘objective’,  
18 requires rationales, clear arguments and motivation, along with elucidation of how this relates  
19 to already existing knowledge. Following generally accepted rules and procedures counters the  
20 dependence of results on the personalities of individual researchers, although there is always a  
21 danger that such generally accepted rules and procedures are inappropriate or suboptimal for  
22 the specific situation at hand. For such reasons, one might question the inclusion of consensus  
23 as a virtue. Its importance stems from the impossibility to access observer-independent reality  
24 which means that exchange between observers is necessary to find out about what can be taken  
25 as real and stable. Consensus cannot be enforced; as a virtue it refers to behaviour that facilitates  
26 consensus.

27 In any case, consensus can only be achieved if researchers attempt to be *impartial* by taking  
28 into account competing perspectives, avoiding to favouring pre-chosen hypotheses, and being  
29 open to criticism. In the context of epidemiology, Greenland (2012) proposed transparency and  
30 neutrality as replacements for objectivity.

31 Going on, the world outside the observer’s mind plays a key role in usual concepts of objec-  
32 tivity, and as explained in Section 2.3 we see it as a major target of science. We acknowledge  
33 that the ‘real world’ is accessible to human beings through observation only, and that scientific  
34 observation and measurement cannot be independent of human preconceptions and theories.  
35 As statisticians we are concerned with making general statements based on systematized obser-  
36 vations, and this makes *correspondence to observed reality* a core concern regarding objectivity.  
37 This is not meant to imply that empirical statements about observations are the only meaning-  
38 ful statements that can be made about reality; we think that scientific theories that cannot be  
39 verified (but can potentially be falsified) by observations are meaningful thought constructs,  
40 particularly because observations are never ‘pure’ and truly independent of thought constructs.  
41 Certainly in some cases the measurements, i.e. the observations that the statistician deals with,  
42 require critical scrutiny before discussing any statistical analysis of them; see Section 2.2.

43 Formal statistical methods contribute to objectivity as far as they contribute to the fulfilment  
44 of these *desiderata*, particularly by making procedures and their implied rationales transparent  
45 and unambiguous.

46 For example, Bayesian statistics is commonly characterized as ‘subjective’ by Bayesians and  
47 non-Bayesians alike. But, depending on how exactly prior distributions are interpreted and used  
48 (see Sections 5.3–5.5), they fulfil or aid some or all of the virtues that were listed above. Priors

1 make the researchers' prior point of view transparent; different approaches of interpreting them  
2 provide different rationales for consensus; 'objective Bayesians' (see Section 5.4) try to make  
3 them impartial; and if suitably interpreted (see Section 5.5) they can be properly grounded in  
4 observations.

### 6 3.2. 'Multiple perspectives' and 'context dependence', instead of 'subjectivity'

7 Science is normally seen as striving for objectivity, and therefore acknowledging subjectivity  
8 can be awkward. But, as noted above already, reality and the facts are accessible only through  
9 individual personal experiences. Different people bring different information and different view-  
10 points, and they will use scientific results in different ways. To enable clear communication and  
11 consensus, differing perspectives need to be acknowledged, which contributes to transparency  
12 and thus to objectivity. Therefore, subjectivity is important to the scientific process. Subjectivity  
13 is valuable in statistics in that it represents a way to incorporate the information coming from  
14 differing perspectives, which are the building blocks of scientific consensus.

15 We propose *awareness of multiple perspectives* and *context dependence* as key virtues making  
16 explicit the value of subjectivity. To the extent that subjectivity in statistics is a good thing, it is  
17 because information truly is dispersed, and, for any particular problem, different stakeholders  
18 have different goals. A counterproductive implication of the idea that science should be 'ob-  
19 jective' is that there is a tendency in the communication of statistical analyses either to avoid  
20 or hide decisions that cannot be made in an automatic, seemingly 'objective' fashion by the  
21 available data. Given that all observations of reality depend on the perspective of an observer,  
22 interpreting science as striving for a unique ('objective') perspective is illusory. Multiple per-  
23 spectives are a reality to be reckoned with and should not be hidden. Furthermore, by avoiding  
24 personal decisions, researchers often waste opportunities to adapt their analyses appropriately  
25 to the context, the specific background and their specific research aims, and to communicate  
26 their perspective more clearly. Therefore we see the acknowledgement of multiple perspectives  
27 and context dependence as virtues, making clearer in which sense subjectivity can be productive  
28 and helpful.

29 The term 'subjective' is often used to characterize aspects of certain statistical procedures  
30 that cannot be derived automatically from the data to be analysed, such as Bayesian prior  
31 distributions, tuning parameters (e.g. the proportion of trimmed observations in trimmed means,  
32 or the threshold in wavelet smoothing), or interpretations of data visualization. Such decisions  
33 are entry points for multiple perspectives and context dependence. The first decisions of this  
34 kind are typically the choice of data to be analysed and the family of statistical models to be  
35 fitted.

36 To connect with the other part of our proposal, the recognition of different perspectives  
37 should be done in a transparent way. We should not say that we set a tuning parameter to 2.5  
38 (say) just because that is our belief. Rather, we should justify the choice explaining clearly how it  
39 supports the research aims. This could be by embedding the choice in a statistical model that can  
40 ultimately be linked back to observable reality and empirical data, or by reference to desirable  
41 characteristics (or avoidance of undesirable artefacts) of the methodology given the use of the  
42 chosen parameter; actually, many tuning parameters are related to such characteristics and  
43 aims of the analysis rather than to some assumed underlying 'belief' (see Section 4.3). In many  
44 cases, such a justification may be imprecise, for example because background knowledge may  
45 be only qualitative and not quantitative or not sufficiently precise to tell possible alternative  
46 choices apart, but often it can be argued that even then conscious tuning or specification of a  
47 prior distribution comes with benefits compared with using default methods of which the main  
48 attraction often is that seemingly 'subjective' decisions can be avoided.

1 To consider an important example, regularization requires such decisions. Default priors on  
2 regression coefficients are used to express the belief that coefficients are typically close to 0, and,  
3 from a non-Bayesian perspective, lasso shrinkage can be interpreted as encoding an external  
4 assumption of sparsity. Sparsity assumptions can be connected to an implicit or explicit model  
5 in which problems are in some sense being sampled from some distribution or probability  
6 measure of possible situations; see Section 5.5. This general perspective (which can be seen  
7 as Bayesian with an implicit prior on states of nature, or classical with an implicit reference  
8 set for the evaluation of statistical procedures) provides a potential basis to connect choices  
9 to experience; at least it makes transparent what kind of view of reality is encoded in the  
10 choices.

11 Tibshirani (2014) wrote that enforcing sparsity is not primarily motivated by beliefs about  
12 the world, but rather by benefits such as computability and interpretability, hinting at the fact  
13 that considerations other than being ‘close to the real world’ often play an important role in  
14 statistics and more generally in science. Even in areas such as social science where no underlying  
15 truly sparse structure exists, imposing sparsity can have advantages such as supporting stability  
16 (Gelman, 2013).

17 In a wider sense, if one is performing a linear or logistic regression, for example, and consid-  
18 ering options of maximum likelihood, the lasso or hierarchical Bayes with a particular structure  
19 of priors, all of these choices are ‘subjective’ in the sense of encoding aims regarding possi-  
20 ble outputs and assumptions, and all are ‘objective’ as far as these aims and assumptions are  
21 made transparent and the assumptions can be justified on the basis of past data and ultimately  
22 be checked given enough future data. So the conventional labelling of Bayesian analyses or  
23 regularized estimates as ‘subjective’ misses the point.

24 For another example, the binomial data confidence interval based on  $(y + 2)/(n + 4)$  gives  
25 better coverage than the classical interval based on  $y/n$  (Agresti and Coull, 1998). Whereas the  
26 latter has a straightforward justification, the former is based on trading interval width against  
27 conservatism and involves some approximation and simplification, which Agresti and Coull  
28 (1998) justified by the fact that the resulting formula can be presented in elementary courses.  
29 Debating whether this is more subjective than the classical approach, and whether this is a  
30 problem, is not helpful. Similarly, when comparing Bayesian estimates of public opinion by using  
31 multilevel regression and poststratification to taking raw survey means (which indeed correspond  
32 to Bayesian analyses under unreasonable flat priors), it is irrelevant which is considered more  
33 subjective.

34 Tuning parameters can be set or estimated on the basis of past data, and also on the basis  
35 of understanding of the effect of the choice on results and a clear explanation why a certain  
36 impact is desired or not. In robust statistics, for example, the breakdown point of some methods  
37 can be tuned and may be chosen lower than the optimal 50% because, if there is too large a  
38 percentage of data deviating strongly from the majority, one may rather want the method to  
39 deliver a compromise between all observations but, if the percentage of outliers is quite low,  
40 one may rather want them to be disregarded, with borderline percentages depending on the  
41 application (particularly on to what extent outliers are interpreted as erroneous observations  
42 rather than as somewhat special but still relevant cases).

43 Here is an example in which awareness of multiple perspectives can help with a problem with  
44 impartiality. Simulation studies for comparing statistical methods are often run by the designers  
45 of one of the competing approaches and, even if this is not so, the person running the study  
46 may have prior opinions about the competitors that may affect the study. There is simply no  
47 ‘objective’ way how this can be avoided; taking into account multiple perspectives by for example  
48 asking designers of all competing methods to provide simulation setups might help here.

### 3.3. Stability

As outlined in Section 2.3, we believe that science aims at a stable and reliable view of reality. Human beings do not have direct access to observer-independent reality, but phenomena that remain stable when perceived through different channels, at different times, and that are confirmed as stable by different observers, are the best contenders to be attributed ‘objective existence.’

The term *stability* is difficult to find in philosophical accounts of objectivity, but it seems that Mayo’s (1996) view of the growth of experimental knowledge through piecemeal testing of aspects of scientific theories and learning from error (which to her is a key feature of objectivity) implicitly aims at probing the stability of these theories. Stability is also connected to subjectivity in the sense that in the best case stability persists under inquiry from as many perspectives and in as many contexts as possible.

The accommodation and analysis of variability is something that statistical modelling brought to science, and in this sense statisticians investigate stability (of observations as well as of the statistics and estimators computed from them) all the time. An investigation of stability just based on variability assuming a parametric model is quite narrow, though, and there are many further sources of potential instabilities. Stability can refer to reproducibility of conclusions on new data, or to alternative analyses of the same data making different choices regarding for example tuning constants, Bayesian priors, transformations, resampling, removing outliers or even completely different methodology as far as this aims at investigating the same issue (alternative analyses that can be interpreted as doing something essentially different cannot be expected to deliver a similar result). On the most basic (but not always trivial) level, the same analysis on the same data should be replicable by different researchers. In statistical theory, basic variability assuming a parametric model can be augmented by robustness against various violations of model assumptions and Bayesian sensitivity analysis.

There are many aspects of stability that can be investigated, and only so much can be expected from a single study or publication; the generation of reliable scientific knowledge generally requires investigation of phenomena from more points of view than that of a single researcher or team.

### 3.4. A list of specific virtues

To summarize the above discussion, we give a more detailed list of the virtues that were discussed above, which we think will improve on discussions in which approaches, analyses and arguments are branded ‘subjective’ or ‘objective’ (Table 1).

In the subsequent discussion we refer to the item numbers in the list in Table 1 starting by V for virtue, such as V4(b) for ‘clear conditions for reproduction, testing, and falsification’.

We are aware that in some situations some of these virtues may oppose each other; for example ‘consensus’ can contradict ‘awareness of multiple perspectives’, and indeed dissent is essential to scientific progress. This tension between impersonal consensus and creative debate is an unavoidable aspect of science. Sometimes the consensus can only be that there are different legitimate points of view. Furthermore, the virtues listed are not all fully autonomous; clear reference to observations may be both a main rationale for consensus and a key contribution to transparency; and the subjective virtues contribute to both transparency and openness to criticism and exchange.

Not all items on the list apply to all situations. For example, in Section 5 we apply the list to the foundations of statistics, but some virtues (such as full communication of procedures) rather apply to specific studies.

Table 1. Virtues

<p><i>V1. Transparency</i></p> <ul style="list-style-type: none"> <li>(a) Clear and unambiguous definitions of concepts</li> <li>(b) Open planning and following agreed protocols</li> <li>(c) Full communication of reasoning, procedures, spelling out of (potentially unverifiable) assumptions and potential limitations</li> </ul> <p><i>V2. Consensus</i></p> <ul style="list-style-type: none"> <li>(a) Accounting for relevant knowledge and existing related work</li> <li>(b) Following generally accepted rules where possible and reasonable</li> <li>(c) Provision of rationales for consensus and unification</li> </ul> <p><i>V3. Impartiality</i></p> <ul style="list-style-type: none"> <li>(a) Thorough consideration of relevant and potentially competing theories and points of view</li> <li>(b) Thorough consideration and if possible removal of potential biases: factors that may jeopardize consensus and the intended interpretation of results</li> <li>(c) Openness to criticism and exchange</li> </ul> <p><i>V4. Correspondence to observable reality</i></p> <ul style="list-style-type: none"> <li>(a) Clear connection of concepts and models to observables</li> <li>(b) Clear conditions for reproduction, testing and falsification</li> </ul> <p><i>V5. Awareness of multiple perspectives</i></p> <p><i>V6. Awareness of context dependence</i></p> <ul style="list-style-type: none"> <li>(a) Recognition of dependence on specific contexts and aims</li> <li>(b) Honest acknowledgement of the researcher's position, goals, experiences, and subjective point of view</li> </ul> <p><i>V7. Investigation of stability</i></p> <ul style="list-style-type: none"> <li>(a) Consequences of alternative decisions and assumptions that could have been made in the analysis</li> <li>(b) Variability and reproducibility of conclusions on new data</li> </ul>
---

#### 4. Applied examples

In conventional statistics, assumptions are commonly minimized. Classical statistics and econometrics are often framed in terms of robustness, with the goal being methods that work with minimal assumptions. But the decisions about what information to include and how to frame the model—these are typically buried, not stated formally as assumptions but just baldly stated: ‘Here is the analysis we did. . .,’ sometimes with the statement or implication that these have a theoretical basis but typically with little clear connection between subject matter theory and details of measurements. From the other perspective, Bayesian analyses are often boldly assumption based but with the implication that these assumptions, being subjective, need no justification and cannot be checked from data.

We would like statistical practice, Bayesian and otherwise, to move towards more transparency regarding the steps linking theory and data to models, and recognition of multiple perspectives in the information that is included in this paper trail and this model. In this section we show how we are trying to move in this direction in some of our recent research projects. We present these examples not as any sort of ideals but rather to demonstrate how we are grappling with these ideas and, in particular, the ways in which active awareness of the concepts of transparency, consensus, impartiality, correspondence to observable reality, multiple perspectives and context dependence is changing our applied work.

##### 4.1. A hierarchical Bayesian model in pharmacology

Statistical inference in pharmacokinetics–pharmacodynamics involves many challenges: data are indirect and often noisy; the mathematical models are non-linear and computationally

expensive, requiring the solution of differential equations; parameters vary by person but often with only a small amount of data on each experimental subject. Hierarchical models and Bayesian inference are often used to get a handle on the many levels of variation and uncertainty (see, for example, Sheiner (1984) and Gelman *et al.* (1996)).

One of us is currently working on a project in drug development involving a Bayesian model that was difficult to fit, even when using advanced statistical algorithms and software. Following the so-called folk theorem of statistical computing (Gelman, 2008), we suspected that the problems with computing could be attributed to a problem with our statistical model. In this case, the issue did not seem to be a lack of fit, or a missing interaction, or unmodelled measurement error—problems we had seen in other settings of this sort. Rather, the fit appeared to be insufficiently constrained, with the Bayesian fitting algorithm being stuck going through remote regions of parameter space that corresponded to implausible or unphysical parameter values.

In short, the model as written was only weakly identified, and the given data and priors were consistent with all sorts of parameter values that did not make scientific sense. Our iterative Bayesian computation had poor convergence—i.e. the algorithm was having difficulty approximating the posterior distribution—and the simulations were going through zones of parameter space that were not consistent with the scientific understanding of our pharmacology colleagues.

To put it another way, our research team had access to prior information that had not been included in the model. So we took the time to specify a more informative prior. The initial model thus played the role of a placeholder or default which could be elaborated as needed, following the iterative prescription of falsificationist Bayesianism (Box (1980) and Gelman *et al.* (2013), section 5.5).

In our experience, informative priors are not so common in applied Bayesian inference and, when they are used, they often seem to be presented without clear justification. In this instance, though, we decided to follow the principle of transparency and to write a note explaining the genesis of each prior distribution. To give a sense of what we are talking about, we present a subset of these notes here:

- $\gamma_1$ : mean of population distribution of  $\log(\text{BVA}_j^{\text{latent}}/50)$ , centered at 0 because the mean of the BVA values in the population should indeed be near 50. We set the prior sd to 0.2 which is close to  $\log(60/50) = 0.18$  to indicate that we're pretty sure the mean is between 40 and 60.
- $\gamma_2$ : mean of pop dist of  $\log(k_j^{\text{in}}/k_j^{\text{out}})$ , centered at 3.7 because we started with  $-2.1$  for  $k^{\text{in}}$  and  $-5.9$  for  $k^{\text{out}}$ , specified from the literature about the disease. We use a sd of 0.5 to represent a certain amount of ignorance: we're saying that our prior guess for the population mean of  $k^{\text{in}}/k^{\text{out}}$  could easily be off by a factor of  $\exp(0.5) = 1.6$ .
- $\gamma_3$ : mean of pop dist of  $\log k_j^{\text{out}}$ , centered at  $-5.8$  with a sd of 0.8, which is the prior that we were given before, from the time scale of the natural disease progression.
- $\gamma_4$ :  $\log E_{\text{max}}^0$ , centered at 0 with sd 2.0 because that's what we were given earlier.'

The  $\gamma$ s here already represent a transformation of the original parameters, BVA (baseline visual acuity; this is a drug for treating vision problems),  $k_{\text{in}}$  and  $k_{\text{out}}$  (rate constants for differential equations that model the diffusion of the drug) and  $E_{\text{max}}^0$ , a saturation parameter in the model. One goal in this sort of work is to reparameterize to unbounded scales (so that normal distributions are more reasonable, and we can specify parameters based on location and scale) and to aim for approximate independence in the prior distribution because of the practical difficulties of eliciting prior correlations. The 'literature about the disease' comes from previously published trials of other drugs for this disease; these trials also include control arms which give us information on the natural progression of visual acuity in the absence of any treatment.

1 We see this sort of painfully honest justification as a template for future Bayesian data analyses.  
 2 The above snippet certainly does not represent an exemplar of best practices, but we see it as  
 3 a ‘good enough’ effort that presents our modelling decisions in the context in which they were  
 4 made.

5 To label this prior specification as ‘objective’ or ‘subjective’ would miss the point. Rather,  
 6 we see it as having some of the virtues of objectivity and subjectivity—notably, transparency  
 7 (virtue V1) and some aspects of consensus (virtue V2) and awareness of multiple perspectives  
 8 (virtue V5)—while recognizing its clear imperfections and incompleteness. Other desirable fea-  
 9 tures would derive from other aspects of the statistical analysis—for example, we use external  
 10 validation to approach correspondence to observable reality (virtue V4), and our awareness of  
 11 context dependence (virtue V6) comes from the placement of our analysis within the larger goal,  
 12 which is to model dosing options for a particular drug.

13 One concern about our analysis which we have not yet thoroughly addressed is sensitivity  
 14 to model assumptions. We have established that the prior distribution makes a difference but  
 15 it is possible that different reasonable priors yield posteriors with greatly differing real world  
 16 implications, which would raise concern about consensus (virtue V2) and impartiality (virtue  
 17 V3). Our response to such concerns, if this sensitivity is indeed a problem, would be to document  
 18 our choice of prior more carefully, thus doubling down on the principle of transparency (virtue  
 19 V1) and to compare with other possible prior distributions supported by other information,  
 20 thus supporting impartiality (virtue V3) and awareness of multiple perspectives (virtue V5).

21 The point is not that our particular choices of prior distributions are ‘correct’ (whatever that  
 22 means), or optimal, or even good, but rather that they are transparent, and in a transparent  
 23 way connected to knowledge. Subsequent researchers—whether supportive, critical or neutral  
 24 regarding our methods and substantive findings—should be able to interpret our priors (and, by  
 25 implication, our posterior inferences) as the result of some systematic process, a process which  
 26 is sufficiently open that it can be criticized and improved as appropriate.  
 27  
 28

#### 29 4.2. *Adjustments for pre-election polls*

30 Wang *et al.* (2014) described another of our recent applied Bayesian research projects, in this  
 31 case a statistical analysis that allows highly stable estimates of public opinion by adjustment of  
 32 data from non-random samples. The particular example that was used was an analysis of data  
 33 from an opt-in survey conducted on the Microsoft Xbox video game platform, a technique that  
 34 allowed the research team, effectively to interview respondents in their living rooms, without  
 35 ever needing to call or enter their houses.

36 The Xbox survey was performed during the two months before the 2012 US presidential  
 37 election. In addition to offering the potential practical benefits of performing a national survey  
 38 using inexpensive data, this particular project made use of its large sample size and panel  
 39 structure (repeated responses on many thousands of Americans) to learn something new about  
 40 US politics: we found that certain swings in the polls, which had been generally interpreted  
 41 as representing large swings in public opinion, actually could be attributed to differential non-  
 42 response, with Democrats and Republicans in turn being more or less likely to respond during  
 43 periods where there was good or bad news about their candidate. This finding was consistent  
 44 with some of the literature in political science (see Erikson *et al.* (2004)), but the Xbox study  
 45 represented an important empirical confirmation (Gelman *et al.*, 2016).

46 Having established the potential importance of the work, we next consider its controversial  
 47 aspects. For many decades, the gold standard in public opinion research has been proba-  
 48 bility sampling, in which the people being surveyed are selected at random from a list or

1 lists (for example, selecting households at random from a list of addresses or telephone num-  
 2 bers and then selecting a person within each sampled household from a list of the adults  
 3 who live there). From this standpoint, opt-in sampling of the sort that was employed in the  
 4 Xbox survey lacks a theoretical foundation, and the estimates and standard errors thus ob-  
 5 tained (and which we reported in our research papers) do not have a clear statistical interpre-  
 6 tation.

7 This criticism—that inferences from opt-in surveys lack a theoretical foundation—is interest-  
 8 ing to us here because it is *not* framed in terms of objectivity or subjectivity. We do use Bayesian  
 9 methods for our survey adjustment but the criticism from certain survey practitioners is not  
 10 about adjustment but rather about the data collection: they take the position that no good  
 11 adjustment is possible for data that are collected from a non-probability sample.

12 As a practical matter, our response to this criticism is that non-response rates in national  
 13 random-digit-dialled telephone polls are currently in the range of 90%, which implies that  
 14 realworld surveys of this sort are essentially opt-in samples in any case: if there is no theoretical  
 15 justification for non-random samples then we are all dead, which leaves us all with the choice  
 16 either to abandon statistical inference entirely when dealing with survey data, or to accept that  
 17 our inferences are model based and to do our best (Gelman, 2014c).

18 Our Bayesian adjustment model (Wang *et al.*, 2014) used prior information in two ways. First,  
 19 population distributions of demographics, state of residence and party identification were im-  
 20 puted by using exit poll data from the previous election; from the survey sampling perspective  
 21 this was a poststratification step, and from the political science perspective this represents an as-  
 22 sumption of stability in the electorate from 2008 to 2012. The second aspect of prior information  
 23 was encoded in our hierarchical logistic regression model, in which varying intercepts for states  
 24 and for different demographic factors were modelled as exchangeable batches of parameters  
 25 drawn from normal distributions. These assumptions are necessarily approximate and are thus  
 26 ultimately justified on pragmatic grounds.

27 We shall now express this discussion by using the criteria from Section 4. Probability sampling  
 28 has the clear advantage of transparency (virtue V1) in that the population and sampling mecha-  
 29 nism can be clearly defined and accessible to outsiders, in a way that an opt-in survey such as the  
 30 Xbox survey is not. In addition, the probability sampling has the benefits of consensus (virtue  
 31 V2), at least in the USA, where such surveys have a long history and are accepted in marketing  
 32 and opinion research. Impartiality (virtue V3) and correspondence to observable reality (virtue  
 33 V4) are less clearly present because of the concern with non-response, just noted. We would  
 34 argue that the large sample size and repeated measurements of the Xbox data, coupled with our  
 35 sophisticated hierarchical Bayesian adjustment scheme, put us well on the road to impartiality  
 36 (through the use of multiple sources of information, including past election outcomes, used  
 37 to correct for biases in the form of known differences between sample and observation) and  
 38 correspondence to observable reality (in that the method can be used to estimate population  
 39 quantities that could be validated from other sources).

40 Regarding the virtues that are associated with subjectivity, the various adjustment schemes  
 41 represent awareness of context dependence (virtue V6) in that the choice of variables to match in  
 42 the population depend on the context of political polling, both in the sense of which aspects of the  
 43 population are particularly relevant for this purpose, and in respecting the awareness of survey  
 44 practitioners of what variables are predictive of non-response. The researcher's subjective point  
 45 of view is involved in the choice of exactly what information to include in weighting adjustments  
 46 and exactly what statistical model to fit in regression-based adjustment. Users of probability  
 47 sampling on grounds of 'objectivity' may shrink from using such judgements and may therefore  
 48 ignore valuable information from the context.

#### 4.3. Transformation of variables in cluster analysis for socio-economic stratification

Cluster analysis aims at grouping similar objects and separating dissimilar objects, and as such is based, explicitly or implicitly, on some measure of dissimilarity. Defining such a measure, for example by using some set of variables characterizing the objects to be clustered, can involve many decisions. Here we consider an example of Hennig and Liao (2013), where we clustered data from the 2007 US Consumer Finances Survey, comprising variables on income, savings, housing, education, occupation, number of checking and savings accounts, and life insurance with the aim of data-based exploration of socio-economic stratification. The choice of variables and the decisions of how they are selected, transformed, standardized and weighted has a strong effect on the results of the cluster analysis. This effect depends to some extent on the clustering technique that is afterwards applied to the resulting dissimilarities but will typically be considerable, even for cluster analysis techniques that are not directly based on dissimilarities. One of the various issues that were discussed by Hennig and Liao (2013) was the transformation of the variables treated as continuous (namely income and savings amount), with the view of basing a cluster analysis on a Euclidean distance after transformation, standardization and weighting of variables.

There is some literature on choosing transformations, but the usual aims of transformation, namely achieving approximate additivity, linearity, equal variances or normality, are often not relevant for cluster analysis, where such assumptions apply only to model-based clustering, and only within the clusters, which are not known before transformation.

The rationale for transformation when setting up a dissimilarity measure for clustering is of a different kind. The measure needs to formalize appropriately which objects are to be treated as ‘similar’ or ‘dissimilar’ by the clustering methods and should therefore be put into the same or different clusters respectively. In other words, the formal dissimilarity between objects should match what could be called the ‘interpretative dissimilarity’ between objects. This is an issue involving subject matter knowledge that cannot be decided by the data alone.

Hennig and Liao (2013) argued that the interpretative dissimilarity between different savings amounts is governed rather by ratios than by differences, so that \$2 million of savings is seen as about as dissimilar from \$1 million, as \$2000 is dissimilar from \$1000. This implies a logarithmic transformation. We do not argue that there is a precise argument that privileges the log-transformation over other transformations that achieve something similar, and one might argue from intuition that even taking logarithms may not be sufficiently strong. We therefore recognize that any choice of transformation is a provisional device and only an approximation to an ideal ‘interpretative dissimilarity’, even if such an ideal exists.

In the data set, there are no negative savings values as there is no information on debts, but there are many people who report zero savings, and it is conventional to kludge the logarithmic transformation to become  $x \mapsto \log(x + c)$  with some  $c > 0$ . Hennig and Liao (2013) then pointed out that, in this example, the choice of  $c$  has a considerable effect on clustering. The number of people with very low but non-zero savings in the data set is quite small. Setting  $c = 1$ , for example, the transformation creates a substantial gap between the zero-savings group and people with fairly low (but non-zero) amounts of savings, and of course this choice is also sensitive to scaling (for example, savings might be coded in dollars, or in thousands of dollars). The subsequent cluster analysis (done by ‘partitioning around medoids’; Kaufman and Rousseeuw (1990)) would therefore separate the zero-savings group strictly; no person with zero savings would appear together in a cluster with a person with non-zero savings. For larger values for  $c$ , the dissimilarity between the zero-savings group and people with a low savings amount becomes effectively sufficiently small that people with zero savings could appear in clusters together with other people, as long as values on other variables are sufficiently similar.

1 We do not believe that there is a true value of  $c$ . Rather, clusterings arising from different  
 2 choices of  $c$  are legitimate but imply different interpretations. The clustering for  $c = 1$  is based  
 3 on treating the zero-savings group as special, whereas the clustering for  $c = 200$ , say, implies that  
 4 a difference in savings between \$0 and \$100 is taken as not so important (although it is more  
 5 important in any case than the difference between \$100 and \$200). Similar considerations hold  
 6 for issues such as selecting and weighting variables and coding ordinal variables.

7 It can be frustrating to the novice in cluster analysis that such decisions for which there do  
 8 not seem to be an objective basis can make such a difference, and there is apparently a strong  
 9 temptation to ignore the issue and just to choose  $c = 1$ , which may look natural in the sense that  
 10 it maps zero onto zero, or even to avoid transformation at all to avoid the discussion, so that no  
 11 obvious lack of objectivity strikes the reader. Having the aim of socio-economic stratification  
 12 in mind, though, it is easy to argue that clusterings that result from ignoring the issue are less  
 13 desirable and useful than a clustering obtained from making a however imprecisely grounded  
 14 decision choosing  $c > 1$ , therefore avoiding either separation of the zero-savings group as a  
 15 clustering artefact or an undue domination of the clustering by people with large savings in case  
 16 of not applying any transformation at all.

17 We believe that this kind of tuning problem that cannot be interpreted as estimating an  
 18 unknown true constant (and does therefore not lend itself naturally to an approach through a  
 19 Bayesian prior) is not exclusive to cluster analysis and is often hidden in presentations of data  
 20 analyses.

21 Hennig and Liao (2013) pointed out the issue and did some sensitivity analysis about the  
 22 strength of the effect of the choice of  $c$  (virtue V7a). The way that we picked  $c$  in Hennig and  
 23 Liao (2013) made clear reference to the context dependence, while being honest that the subject  
 24 matter knowledge in this case provided only weak guidelines for making this decision (virtue  
 25 V6). We were also clear that alternative choices would amount to alternative perspectives rather  
 26 than being just wrong (virtues V5 and V3).

27 The issue how to foster consensus and to make a connection to observable reality (virtues V2  
 28 and V4) is of interest but is not treated here.

29 But it is problematic to establish rationales for consensus that are based on ignoring con-  
 30 text and potentially multiple perspectives. There is a tendency in the cluster analysis literature  
 31 to seek formal arguments for making such decisions automatically (see, for example, Everitt  
 32 *et al.*, (2011), section 3.7, on variable weighting; it is difficult to find anything systematic in the  
 33 clustering literature on transformations), trying to optimize ‘clusterability’ of the data set, or  
 34 preferring methods that are less sensitive to such decisions, because this amounts to making  
 35 the decisions implicitly without giving the researchers access to them. In other words, the data  
 36 are given the authority to determine not only which objects are similar (which is what we want  
 37 them to do), but also what similarity should mean. The latter should be left to the researcher,  
 38 although we acknowledge that the data can have a certain influence: for example the idea that  
 39 dissimilarity of savings amounts is governed by ratios rather than differences is connected to  
 40 (but not determined by) the fact that the distribution of savings amounts is skewed, with large  
 41 savings amounts sparsely distributed.

#### 44 4.4. *Testing for homogeneity against clustering*

45 An issue in Hennig and Liao (2013) was whether there is any meaningful clustering to be found  
 46 in the data. Some sociologists suspect that, in many modern democratic societies, stratification  
 47 may represent no more than a drawing of arbitrary borders through a continuum of socio-  
 48 economic conditions. We were interested in what the data have to say on this issue, and we

1 chose to address this by running a test of a homogeneity null hypothesis against a clustering  
2 alternative (knowing that there is some distance to go between the result of such an analysis and  
3 the ‘desired’ sociological interpretation).

4 If we had been concerned primarily with appearing objective and the ease to achieve a sig-  
5 nificant result, probably we would have chosen a likelihood ratio test of the null hypothesis of  
6 a standard homogeneity model (in the specific situation this could have been a Gaussian dis-  
7 tribution for the continuous variables, an uncorrelated adjacent category ordinal logit model  
8 for ordinal variables and a locally independent multinomial model for categorical data) for a  
9 single mixture component in the framework of mixture models as provided, for example, in the  
10 LatentGOLD software package (Vermunt and Magidson, 2016).

11 But even in the absence of meaningful clusters, real data do not follow such clean distributional  
12 shapes and therefore sufficiently large data sets (including ours, with  $n > 17000$ ) will almost  
13 always reject a simple homogeneity model. We therefore set out to build a null model that  
14 captured the features of the data set such as the dependence between variables and marginal  
15 distributions of the categorical variables as well as possible, without involving anything that  
16 could be interpreted as clustering structure. As opposed to the categorical variables, the marginal  
17 distributions of the ‘continuous’ variables such as the transformed savings amount were treated  
18 as potentially indicating clustering, and therefore the null model used non-parametric unimodal  
19 distributions for them. Data from this null model involving several characteristics estimated from  
20 the data could be simulated by using the parametric bootstrap.

21 As test statistic we used a cluster validity statistic of the clustering computed on the data,  
22 which was not model based but dissimilarity based. The idea behind this was that we wanted  
23 a test statistic which would measure the degree of clustering, so that we could find out how  
24 much ‘clustering’ we could expect to see even if no meaningful clustering was present (under  
25 the null model). Actually we computed clusterings for various numbers of clusters. Rather than  
26 somehow to define a single  $p$ -value from aggregating all these clusterings (or selecting the ‘best’  
27 one), we decided to show a plot of the values of the validity statistic for the different numbers of  
28 clusters for the real data set together with the corresponding results for many data sets simulated  
29 from the null model. The result of this showed clearly that a higher level of clustering was found  
30 in the real data set.

31 In doing this, we deviated from classical significance test logic in several ways, by not using  
32 a test statistic that was optimal against any specific alternative, by not arguing from a single  
33  $p$ -value and by using a null model that relied heavily on the data to try as hard as we can to  
34 model the data without clustering. Still, in case that the validity statistic values for the real data  
35 do not look clearly different from those of the bootstrapped data set, this can be interpreted as  
36 no evidence in the data for real clustering, whereas the interpretation of a clear (‘significant’)  
37 difference depends on whether we can argue convincingly that the null model is as good as  
38 possible at trying to model the data without clustering structure. Setting up a straw man null  
39 model for homogeneity and rejecting it would have been easy and not informative. The general  
40 principle is discussed in more detail in Hennig and Lin (2015), including real data examples  
41 where such a null model could not be rejected, as opposed to a straw man model.

42 The essence here is that we made quite a number of decisions that opened our analysis more  
43 clearly to the charge of ‘not being objective’ than following a standard approach, for the sake  
44 of adapting the analysis better to the specific data in hand, and of giving the null hypothesis the  
45 best possible chance (the non-rejection of it would have been a non-discovery here; the role of  
46 it was not to be ‘accepted’ as ‘true’ anyway).

47 We tried to do good science, though, by checking as impartially and transparently as we  
48 could (virtues V1 and V3), whether the data support the idea of a real clustering (virtue V4).

This involved context-dependent judgement (virtue V6) and the transparent choice of a specific perspective (the chosen validity index) among a potential variety (virtue V5), because we were after more qualitative statements than degrees of belief in certain models.

## 5. Decomposing subjectivity and objectivity in the foundations of statistics

In this section, we use the above list of virtues to revisit aspects of the discussion on fundamental approaches to statistics, for which the terms ‘subjective’ and ‘objective’ typically play a dominant role. We discuss what we perceive to be the major streams of the foundations of statistics, but within each of these streams there are several approaches, which we cannot cover completely in such a paper; rather we sketch the streams somewhat roughly and refer to only a single or a few leading references for details where needed.

Here, we distinguish between interpretations of probability, and approaches for statistical inference. For example, ‘frequentism’ as an interpretation of probability does not necessarily imply that Fisherian or Neyman–Pearson tests are preferred to Bayesian methods, despite the fact that frequentism is more often associated with the former than with the latter.

We shall go through several philosophies of statistical inference, for each laying out the connections that we see to the virtues of objectivity and subjectivity outlined in Section 3.4.

Exercising awareness of multiple perspectives, we emphasize that we do not believe that one of these philosophies is the correct or best; nor do we claim that reducing the different approaches to a single approach would be desirable. What is lacking here is not unification, but rather, often, transparency about which interpretation of probabilistic outcomes is intended when applying statistical modelling to specific problems. Particularly, we think that, depending on the situation, both ‘aleatory’ or ‘epistemic’ approaches to modelling uncertainty are legitimate and worthwhile, referring to data-generating processes in observer-independent reality on one hand and rational degrees of belief on the other.

We focus on approaches that are conventionally labelled as either Bayesian or frequentist, but we acknowledge that there are important perspectives on statistics that lie outside this traditional divide. Discussing them in detail would be worthwhile but is beyond our focus, and we hope that discussants of our paper will pick up these threads. Examples of other perspectives include *machine learning*, where the focus is on prediction rather than parameter estimation; thus there is more emphasis on correspondence to observable reality (virtue V4) compared with other virtues, *alternative models of uncertainty* such as belief functions, imprecise probabilities, and fuzzy logic that aim to circumvent some of the limitations of probability theory (most notoriously, the difficulty of distinguishing between ‘known unknowns’ and ‘unknown unknowns,’ or risk and uncertainty in the terminology of Knight, 1921) and *exploratory data analysis* (Tukey, 1977), which is sensitive to multiple perspectives (virtue V5) and context dependence (virtue V6), and tries to be more directly connected to the data than if it was mediated by probability models (virtue V4(a)). Whether avoidance of probability modelling contributes to transparency (virtue V1(a)) is rather problematic because implicit assumptions that may not be spelled out (virtue V1(c)) can be controversial.

### 5.1. Frequentist probabilities

‘Frequentism’ as an interpretation of probability refers, in a narrow sense, to the identification of the probability of an event in a certain experiment with a limiting relative frequency of occurrences if the experiment were to be carried out infinitely often in some kind of independent manner. Frequentist statistics is based on evaluating procedures based on a long-term average over a ‘reference set’ of hypothetical replicated data sets. Different choices of reference sets were

1 for example used by Fisher (1955) and Pearson (1955) when discussing permutation or  $\chi^2$ -tests  
2 for  $2 \times 2$  tables.

3 In the wider sense, we call probabilities ‘frequentist’ when they formalize observer-independent  
4 tendencies or propensities of experiments to yield certain outcomes (see, for example, Gillies  
5 (2000)), which are thought of as replicable and yielding a behaviour under infinite replication  
6 as suggested by what is assumed to be the ‘true’ probability model.

7 The frequentist mindset locates probabilities in the observer-independent world, so they are  
8 in this sense objective (and often called ‘objective’ in the literature, e.g. Kendall (1949)). This,  
9 however, does not guarantee that frequentist probabilities really exist; an infinite number of  
10 replicates cannot exist, and even a finite amount of real replicates will neither be perfectly  
11 identical nor perfectly independent. Ultimately the ideally infinite populations of replicates are  
12 constructed by the ‘statistician’s imagination’ (Fisher, 1955).

13 The decision to adopt the frequentist interpretation of probability regarding a certain phe-  
14 nomenon therefore requires idealization. It cannot be enforced by observation; nor is there  
15 generally enough consensus that this interpretation applies to any specific setup, although it  
16 is well discussed and supported in some physical settings such as radioactive decay (virtues  
17 V2 and, V4). Once a frequentist model has been adopted, however, it makes predictions about  
18 observations that can be checked, so the reference to the observable reality (virtue V4) is clear.

19 There is some disagreement about whether the frequentist definition of probability is clear and  
20 unambiguous (virtue V1(a)). On one hand, the idea of a tendency of an experiment to produce  
21 certain outcomes as manifested in observed and expected relative frequencies seems quite clear.  
22 On the other hand, it is difficult to avoid the circularity that would result from referring to inde-  
23 pendent and identical replicates when defining frequentist probabilities, because the standard  
24 definition of the terms ‘independent’ and ‘identical’ assumes a definition of probability that is  
25 already in place (see von Mises (1957) for a prominent attempt to solve this, and Fine (1973)  
26 for a criticism).

27 Frequentism implies that, in the observer-independent reality, true probabilities are unique,  
28 but there is considerable room for multiple perspectives (virtue V5) regarding the definition of  
29 replicable experiments, collectives or reference sets. The idea of replication is often constructed  
30 in a rather creative way. For example, in time series modelling the frequentist interpretation  
31 implies an underlying true distribution for every single time point, but there is no way to repeat  
32 observations independently at the same time point. This actually means that the effective sample  
33 size for time series data would be 1, if replication were not implicitly constructed in the statistical  
34 model, e.g. by assuming independent innovations in auto-regressive moving average type models.  
35 Such models, or, more precisely, certain aspects of such models, can be checked against the data  
36 but, even if such a check does not fail, it is still clear that there is no such thing in observable  
37 reality, even approximately, as a marginal ‘true’ frequentist distribution of the value of the time  
38 series  $x_t$  at fixed  $t$ , as implied by the model, because  $x_t$  is strictly not replicable.

39 The issue that useful statistical models require a construction of replication (or exchangeabil-  
40 ity) on some level by the statistician is, as we discuss below, not confined to frequentist models.  
41 To provide a rationale for the essential statistical task of pooling information from many ob-  
42 servations to make inference relevant for future observations, all these observations need to be  
43 assumed to represent the same process somehow.

44 The appropriateness of such assumptions in a specific situation can often only be tested in  
45 quite a limited way by observations. All kinds of informal arguments can apply about why it is  
46 a good or bad idea to consider a certain set of observations (or unobservable implied entities  
47 such as error terms and latent variables) as independent and identically distributed frequentist  
48 replicates.

1 Unfortunately, although such an openness to multiple perspectives and potential context  
 2 dependence (virtue V6(a)) can be seen as positive from our perspective, those issues involved in  
 3 the choices of a frequentist reference set are often not clearly communicated and discussed. The  
 4 existence of a true model with implied reference set is typically taken for granted by frequentists,  
 5 motivated at least in part by the desire for objectivity.

## 6 5.2. Frequentist inference

7 This section is about inference from data about characteristics of an assumed true frequen-  
 8 tist probability model. Traditionally, this comprises hypothesis tests, confidence intervals and  
 9 parameter estimators but is not limited to them; see below.

10 According to Mayo and Spanos (2010) and Cox and Mayo (2010), a fundamental feature of  
 11 frequentist inference is the evaluation of error probabilities, i.e. probabilities of wrong decisions.  
 12 Traditionally these would be the type I and type II errors of Neyman–Pearson hypothesis testing,  
 13 but the the error statistical perspective could also apply to other constructs such as errors of  
 14 sign and magnitude (‘type S’ and ‘type M’ errors; Gelman and Carlin (2014)).

15 Mayo and co-workers see the ability to learn from error and to test models severely (in such  
 16 a way that it would be difficult for a model to pass a test if it was wrong regarding the specific  
 17 aspect that is assessed by a test) against data as a major feature of objectivity, which is made  
 18 possible by the frequentist interpretation of probability measures as ‘data generators’. In our  
 19 list of virtues, this feature is captured in virtue V4(b) (reference to observations: reproduction;  
 20 testing; falsification). The underlying idea, with which we agree, is that learning from error is  
 21 a main driving force in science: a lifetime contract between the mode of statistical investiga-  
 22 tion and its object. This corresponds to Chang’s active scientific realism that was mentioned  
 23 above.

24 The error probability characteristics of methods for frequentist inference rely, in general, on  
 25 model assumptions. In principle, these assumptions can be tested, also, and are therefore, ac-  
 26 cording to Mayo and co-workers, no threat to the objectivity of the account. But this comes with  
 27 two problems. Firstly, derivations of statistical inference based on error probabilities typically  
 28 assume the model as fixed and do not account for prior model selection based on the data.  
 29 This issue has recently attracted some research (e.g. Berk *et al.*, (2013)), but this still requires  
 30 a transparent listing of all the possible modelling decisions that could be made virtue (V1(b)),  
 31 which often is missing, and which may not even be desirable as long as the methods are used  
 32 in an exploratory fashion (Gelman and Loken, 2014). Secondly, any data set can be consistent  
 33 with many models, which can lead to divergent inferences. Davies (2014) illustrates this with  
 34 the analysis of a data set on amounts of copper in drinking water, which can be fitted well by a  
 35 Gaussian, a double-exponential and a comb distribution, but yields vastly different confidence  
 36 intervals for the centre of symmetry (which is assumed to be the target of inference) under these  
 37 three models.

38 Davies (2014) suggested that it is misleading to hypothesize models or parameters to be ‘true’.  
 39 According to Davies, statistical modelling is about approximating the data in the sense that  
 40 ‘adequate’ models are not rejected by tests based on characteristics of the data that the statisti-  
 41 cian is interested in (allowing for multiple perspectives and context dependence virtues V5 and  
 42 V6), i.e. they generate data that ‘look like’ the observed data with respect to the chosen char-  
 43 acteristics. Regarding these characteristics, according to Davies, there is no essential difference  
 44 between parameter values and distributional shapes or structural assumptions, and therefore no  
 45 conceptual separation as in traditional frequentist inference between checking model assump-  
 46 tions and inference about parameters assuming a parametric model to be true. Such an approach  
 47 is tied to the observations in a more direct way without making metaphysical assumptions about  
 48

unobservable features of observer-independent reality (virtues V1(a) and V4). It is frequentist inference in the sense that the probability models are interpreted as ‘data generators’.

Two further streams in frequentist inference are concerned about the restrictiveness of parametric model assumptions. Robust statistics explores the stability (virtue V7) of inferences in case the ‘true’ model is not equal to the nominal model, but rather in some neighbourhood, and strives to develop methods that are stable in this respect. There are various ways to define such neighbourhoods and to measure robustness, so robustness considerations can bring in multiple perspectives (virtue V5) but may cause problems with reaching consensus (virtue V2).

Non-parametric statistics allows us to remove bias (virtue V3(c)) by minimizing assumptions regarding, for example, distributional shapes (structural assumptions such as independence are still required). In some cases, particularly with small data sets, this must be afforded by decreased stability (virtue V7).

Overall, there is no shortage of entry points for multiple perspectives (virtue V5) in frequentist inference. This could be seen as something positive, but it runs counter to some extent to the way that the approach is advertised as objective by some of its proponents. Many frequentist analyses could in our opinion benefit from acknowledging honestly their flexibility and the researcher’s choices made, many of which cannot be determined by data alone.

### 5.3. *Subjectivist Bayesianism*

We call ‘subjectivist epistemic’ the interpretation of probabilities as quantifications of strengths of belief of an individual, where probabilities can be interpreted as derived from, or implementable through, bets that are coherent in that no opponent can cause sure losses by setting up some combinations of bets. From this requirement of coherence, the usual probability axioms follow (virtue V2(c)). Allowing conditional bets implies Bayes’s theorem, and therefore, as far as inference concerns learning from observations about not (yet) observed hypotheses, Bayesian methodology is used for subjectivist epistemic probabilities: hence the term ‘subjectivist Bayesianism’.

A major proponent of subjectivist Bayesianism was de Finetti (1974). de Finetti was not against objectivity in general. He viewed observed facts as objective, as well as mathematics and logic and certain formal conditions of random experiments such as the set of possible outcomes. But he viewed uncertainty as something subjective and he held that objective (frequentist) probabilities do not exist. He claimed that his subjectivist Bayesianism appropriately takes into account both the objective (see above) and subjective (opinions about unknown facts based on known evidence) components for probability evaluation.

In de Finetti’s work the term ‘prior’ refers to all probability assignments using information that is external to the data at hand, with no fundamental distinction between the ‘parameter prior’ assigned to parameters in a model, and the form of the ‘sampling distribution’ given a fixed parameter, in contrast with common Bayesian practice today, in which the term ‘prior’ is used to refer only to the parameter prior. In the following discussion we shall use the term ‘priors’ in de Finetti’s general sense.

Regarding the list of virtues in Table 1 in Section 3.4, de Finetti provided a clear definition of probability (virtue V1(a)) based on principles that he sought to establish as generally acceptable (virtue V2(c)). Unlike objectivist Bayesians, subjectivist Bayesians do not attempt to enforce agreement regarding prior distributions, not even given the same evidence; still, de Finetti (1974) and other subjectivist Bayesians proposed rational principles for assigning prior probabilities. There is also some work on (partial) intersubjective agreement on prior specifications, e.g. Dawid (1982a), providing a rationale for consensus (virtue V2(c)). The difference between the objectivist and subjectivist Bayesian point of view is rooted in the general tension in science that

1 was explained above; the subjectivist approach can be criticized for not supporting agreement  
 2 sufficiently—conclusions based on one prior may be seen as irrelevant for somebody who holds  
 3 another (virtue V2(c))—but can be defended for honestly acknowledging that prior information  
 4 often does not come in ways that allow a unique formalization (virtue V6(b)). In any case it  
 5 is vital that subjectivist Bayesians explain transparently how they arrive at their priors, so that  
 6 other researchers can decide to what extent they can support the conclusions (virtue V1(c)).

7 In de Finetti's conception, probability assessments, prior and posterior, can ultimately only  
 8 concern observable events, because bets can only be evaluated if the experiment on which a bet  
 9 is placed has an observable outcome, and so there is a clear connection to observables (virtue  
 10 V4(a)).

11 However, priors in the subjectivist Bayesian conception are not open to falsification (virtue  
 12 V4(b)), because by definition they must be fixed before observation. Adjusting the prior after  
 13 having observed the data to be analysed violates coherence. The Bayesian system as derived  
 14 from axioms such as coherence (as well as those used by objectivist Bayesians; see Section 5.4)  
 15 is designed to cover all aspects of learning from data, including model selection and rejection,  
 16 but this requires that all potential later decisions are already incorporated in the prior, which  
 17 itself is not interpreted as a testable statement about yet unknown observations. In particular  
 18 this means that, once a coherent subjectivist Bayesian has assessed a set-up as exchangeable *a*  
 19 *priori*, he or she cannot drop this assumption later, whatever the data are (think of observing 20  
 20 0s, then 20 1s, and then 10 further 0s in a binary experiment). This is a major problem, because  
 21 subjectivist Bayesians use de Finetti's theorem to justify working with parameter priors and  
 22 sampling models under the assumption of exchangeability, which is commonplace in Bayesian  
 23 statistics. Dawid (1982b) discussed calibration (quality of match between predictive probabilities  
 24 and the frequency of predicted events to happen) of subjectivist Bayesians inferences, and he  
 25 suggested that badly calibrated Bayesians could do well to adjust their future priors if this is  
 26 needed to improve calibration, even at the cost of violating coherence.

27 Subjectivist Bayesianism scores well on virtues V5 and V6(b). But it is a limitation that the  
 28 prior distribution exclusively formalizes belief; context and aims of the analysis do not enter  
 29 unless they have implications about belief. In practice, an exhaustive elicitation of beliefs is  
 30 rarely feasible, and mathematical and computational convenience often plays a role in setting  
 31 up subjective priors, despite de Finetti's having famously accused frequentists of 'ad hoceries  
 32 for mathematical convenience'. Furthermore, the assumption of exchangeability will hardly  
 33 ever precisely match an individual's beliefs in any situation—even if there is no specific reason  
 34 against exchangeability in a specific set-up, the implicit commitment to stick to it whatever  
 35 will be observed seems too strong—but some kind of exchangeability assumption is required by  
 36 Bayesians for the same reason for which frequentists need to rely on independence assumptions:  
 37 some internal replication in the model is needed to allow generalization or extrapolation to future  
 38 observations; see Section 5.1.

39 Summarizing, we view much of de Finetti's criticism of frequentism as legitimate, and objec-  
 40 tivist Bayesianism comes with a commendable honesty about the effect of subjective decisions  
 41 and allows for flexibility accommodating multiple perspectives. But checking and falsification  
 42 of the prior are not built into the approach, and this can obstruct agreement between observers.

#### 44 5.4. Objectivist Bayesianism

45 Given the way that objectivity is often advertised as a key scientific virtue (often without speci-  
 46 fying what exactly it means), it is not surprising that de Finetti's emphasis on subjectivity is not  
 47 shared by all Bayesians, and that there have been many attempts to specify prior distributions  
 48 in a more objective way. Currently the approach of Jaynes (2003) seems to be among the most

1 popular. As with many of his predecessors such as Jeffreys and Carnap, Jaynes saw probability  
 2 as a generalization of binary logic to uncertain propositions. Cox (1961) proved that, given  
 3 a certain list of supposedly commonsense *desiderata* for a ‘plausibility’ measurement, all such  
 4 measurements are equivalent, after suitable scaling, to probability measures. This theorem is the  
 5 basis of Jaynes’s objectivist Bayesianism, and the claim to objectivity comes from postulating  
 6 that, given the same information, everybody should come to the same conclusions regarding  
 7 plausibilities: prior and posterior probabilities (virtue V2(c)), a statement with which subjectivist  
 8 Bayesians disagree.

9 In practice, this objectivist ideal seems to be difficult to achieve, and Jaynes (2003) admitted  
 10 that setting up objective priors including all information is an unsolved problem. One may  
 11 wonder whether his ideal is achievable at all. For example, in chapter 21, he gave a full Bayesian  
 12 ‘solution’ to the problem of dealing with and identifying outliers, which assumes that prior  
 13 models must be specified for both ‘good’ and ‘bad’ data (between which therefore there must be  
 14 a proper distinction), including parameter priors for both models, as well as a prior probability  
 15 for any number of observations to be ‘bad’. It is difficult to see, and no information about this  
 16 was provided by Jaynes himself, how it can be possible to translate the unspecific information of  
 17 knowing of some outliers in many kinds of situations, some of which are more or less related, but  
 18 none identical (say) to the problem at hand, into precise quantitative specifications as needed  
 19 for Jaynes’s approach in an objective way, all before seeing the data.

20 Setting aside the difficulties of working with informally specified prior information, a key  
 21 issue of objectivist Bayesianism is the specification of an objective prior distribution formalizing  
 22 the absence of information. Various principles for doing this have been proposed (maximum  
 23 entropy, Jaynes (2003); maximum missing Shannon information, Berger *et al.* (2009); and a set  
 24 of desirable properties, Bayarri *et al.*, 2012). Such principles have their difficulties and disagree  
 25 in many cases (Kass and Wasserman, 1996). Objectivity seems to be an ambition rather than  
 26 a description of what indeed can be achieved by setting up objectivist Bayesian priors. More  
 27 modestly, therefore, Berger *et al.* (2009) used the term ‘reference priors’, avoiding the term  
 28 ‘objective’, and emphasizing that it would be desirable to have a convention for such cases  
 29 (virtue V2(b)), but admitting that it may not be possible to prove any general approach for  
 30 arriving at such a convention uniquely correct or optimal in any rational sense. However, the  
 31 proposal and discussion of such principles certainly served transparency (virtues V1 (a) and  
 32 I(c)) and provided rationales for consensus (virtue V2(c)).

33 Apart from the issue of the objectivity of the specification of the prior, by and large the  
 34 objectivist Bayesian approach has similar advantages and disadvantages regarding our list of  
 35 virtues as its subjectivist cousin. Particularly it comes with the same difficulties regarding the  
 36 issue of falsifiability from observations. Prior probabilities are connected to logical analysis of  
 37 the situation rather than to betting rates for future observations as in de Finetti’s subjectivist  
 38 approach, which makes the connection of objectivist Bayesian prior probabilities to observations  
 39 even weaker than in the subjectivist Bayesian approach (probabilistic logic has applications other  
 40 than statistical data analysis, for which this may not be a problem).

41 The merit of objectivist Bayesianism is that the approach comes with a much stronger drive  
 42 to justify prior distributions in a transparent way using principles that are as clear and general  
 43 as possible.

#### 44 45 46 5.5. *Falsificationist Bayesianism, and frequentist probabilities in Bayesian statistics*

47 For both subjectivist and objectivist Bayesians, probability models including both parameter  
 48 priors and sampling models do not model the data-generating process, but rather represent

1 plausibility or belief from a certain point of view. Plausibility and belief models can be modified  
 2 by data in ways that are specified *a priori*, but they cannot be falsified by data.

3 In much applied Bayesian work, in contrast, the sampling model is interpreted, explicitly  
 4 or implicitly, as representing the data-generating process in a frequentist or similar way, and  
 5 parameter priors and posteriors are interpreted as giving information about what is known  
 6 about the ‘true’ parameter values. It has been argued that such work does not directly run  
 7 counter to the subjectivist or objectivist philosophy, because the ‘true parameter values’ can  
 8 often be interpreted as expected large sample functions given the prior model (Bernardo and  
 9 Smith, 1994), but the way in which classical subjectivist or objectivist statistical data analysis is  
 10 determined by the untestable prior assignments is seen as unsatisfactory by many statisticians.

11 In any case, the frequentist interpretation of a probability distribution as ‘data generator’  
 12 is regularly used to investigate how Bayesian analyses perform under such assumptions, theo-  
 13 retically, often by analysis of asymptotic properties or by simulation. Wasserman (2006) called  
 14 Bayesian methods with good frequentist properties ‘objective’, referring to the ‘representing  
 15 things in the observer-independent world’ sense of objectivity, but also providing a connection  
 16 of Bayesian models to observables (virtue V4(a)). Rubin (1984) discussed frequentist approaches  
 17 for studying the characteristics of Bayesian methods under misspecified models, i.e. stability  
 18 (virtue V7).

19 The suggestion of testing aspects of the prior distribution by observations using error statisti-  
 20 cal techniques has been around for some time (Box, 1980). Gelman and Shalizi (2013) incor-  
 21 porated this in an outline of what we refer to here as ‘falsificationist Bayesianism’, a philosophy  
 22 that openly deviates from both objectivist and subjectivist Bayesianism, integrating Bayesian  
 23 methodology with an interpretation of probability that can be seen as frequentist in a wide sense  
 24 and with an error statistical approach to testing assumptions in a bid to satisfy virtue V4(b).

25 Falsificationist Bayesianism follows the frequentist interpretation of the probabilities that is  
 26 formalized by the sampling model given a true parameter, so that these models can be tested  
 27 by using frequentist inference (with the limitations that such techniques have, as discussed in  
 28 Section 5.2). Gelman and Shalizi (2013), argued as some frequentists do, that such models  
 29 are idealizations and should not be believed to be literally true, but that the scientific process  
 30 proceeds from simplified models through test and potential falsification by improving the models  
 31 where they are found to be deficient.

32 To put it another way, it is desirable for Bayesian intervals to have close to nominal coverage  
 33 both conditionally on any observables and unconditionally; the desire for this coverage leads  
 34 naturally to calibration checks, which in turn motivates the modification or even rejection  
 35 of models that are not well calibrated empirically. This process serves the correspondence to  
 36 observable reality (virtue V4) while putting more of a burden on transparency (virtue V1) and  
 37 stability (virtue V7) in that the ultimate choice of model can depend on the decision of what  
 38 aspects of the fitted model will be checked.

39 A key issue regarding transparency of falsificationist Bayesianism is how to interpret the  
 40 parameter prior, which does not usually (if occasionally) refer to a real mechanism that produces  
 41 frequencies. Major options are firstly to interpret the parameter prior in a frequentist way, as  
 42 formalizing a more or less idealized data-generating process generating parameter values. A  
 43 bold idealization would be to view ‘all kinds of potential studies with the (statistically) same  
 44 parameter’ as the relevant population, even if the studies are about different topics with different  
 45 variables, in which case more realizations exist, but it is difficult to view a specific study of interest  
 46 as a ‘random draw’ from such a population.

47 Alternatively, the parameter prior may be seen as a purely technical device, serving aims such  
 48 as regularization, without making any even idealized assumption it corresponds to anything

1 that ‘exists’ in the real world. In this case the posterior distribution does not have a proper  
 2 direct interpretation, but statistics such as the posterior mean or uncertainty intervals could be  
 3 interpreted on the basis of their frequentist properties.

4 Overall, falsificationist Bayesianism combines the virtue of error statistical falsifiability with  
 5 virtues V5 and V6 connected to subjectivity. However, the flexibility of the falsificationist  
 6 Bayesian approach—its openly iterative and tentative nature—creates problems regarding clar-  
 7 ity and unification.

## 8 9 10 **6. Discussion**

### 11 *6.1. Implications for statistical theory and practice*

12 At the level of discourse, we would like to move beyond a subjective *versus* objective shouting  
 13 match. But our goals are larger than this. Gelman and Shalizi (2013) on the philosophy of  
 14 Bayesian statistics sought not just to clear the air but also to provide philosophical and rhetorical  
 15 space for Bayesians to feel free to check their models and for applied statisticians who were  
 16 concerned about model fit to feel comfortable with a Bayesian approach. In the present paper,  
 17 our goals are for scientists and statisticians to achieve more of the specific positive qualities  
 18 into which we decompose objectivity and subjectivity in Section 3.4. At the present time, we  
 19 feel that concerns about objectivity are obstructing researchers trying out different ideas and  
 20 considering different sources of inputs to their model, whereas an ideology of subjectivity is  
 21 limiting the degree to which researchers are justifying and understanding their model.

22 There is a tendency for hard-core believers in objectivity needlessly to avoid the use of valuable  
 23 external information in their analyses, and for subjectivists, but also for statisticians who want to  
 24 make their results seem strong and uncontroversial, to leave their assumptions unexamined. We  
 25 hope that our new framing of transparency, consensus, avoidance of bias, reference to observable  
 26 reality, multiple perspectives, dependence on context and aims, investigation of stability and  
 27 honesty about the researcher’s position and decisions will give researchers of all stripes the  
 28 impetus and, indeed, permission, to integrate different sources of information in their analyses,  
 29 to state their assumptions more clearly and to trace these assumptions backwards to past data  
 30 that justify them and forwards to future data that can be used to validate them.

31 Also, we believe that the pressure to appear objective has led to confusion and even dishonesty  
 32 regarding data coding and analysis decisions which cannot be motivated in supposedly objective  
 33 ways; see van Loo and Romeijn (2015) for a discussion of this point in the context of psychiatric  
 34 diagnosis. We prefer to encourage a culture in which it is acceptable to be open about the reasons  
 35 for which decisions are made, which may at times be a mathematical convenience, or the aim  
 36 of the study, rather than strong theory or hard data. It should be recognized openly that the  
 37 aim of statistical modelling is not always to make the model as close as possible to observer-  
 38 independent reality (which always requires idealization anyway), and that some decisions are  
 39 made, for example, to make outcomes more easily interpretable for specific target audiences.

40 Our key points are as follows:

- 41 (a) multiple perspectives correspond to multiple lines of reasoning, not merely to mindless  
 42 and unjustified guesses and
- 43 (b) what is needed is not just a prior distribution or a tuning parameter, but a statistical  
 44 approach in which these choices can be grounded, either empirically or by connecting  
 45 them in a transparent way to the context and aim of the analysis.  
 46

47 For these reasons, *we do not think it at all accurate to limit Bayesian inference to ‘the analysis*  
 48 *of subjective beliefs’*. Yes, Bayesian analysis can be expressed in terms of subjective beliefs, but

1 it can also be applied to other settings that have nothing to do with beliefs (except to the extent  
 2 that all scientific inquiries are ultimately about what is believed about the world).

3 Similarly, *we would not limit classical statistical inference to 'the analysis of simple random*  
 4 *samples'*. Classical methods of hypothesis testing, estimation, and data reduction can be applied  
 5 to all sorts of problems that do not involve random sampling. There is no need to limit the  
 6 applications of these methods to a narrow set of sampling or randomization problems; rather,  
 7 it is important to clarify the foundation for using the mathematical models for a larger class of  
 8 problems.

## 10 6.2. Beyond 'objective' and 'subjective'

11 The list in Table 1 in Section 3.4 is the core of the paper. The list may not be complete, and such a  
 12 list may also be systematized in different ways. Particularly, we developed the list having partic-  
 13 ularly applied statistics in mind, and we may have missed aspects of objectivity and subjectivity  
 14 that are not connected in some sense to statistics. In any case, we believe that the given list can be  
 15 helpful in practice for researchers, for justifying and explaining their choices, and for recipients  
 16 of research work, for checking to what extent the virtues listed are practised in scientific work. A  
 17 key issue here is transparency, which is required for checking all the other virtues. Another key  
 18 issue is that subjectivity in science is not something to be avoided at any cost, but that multiple  
 19 perspectives and context dependence are actually basic conditions of scientific inquiry, which  
 20 should be explicitly acknowledged and taken into account by researchers. We think that this is  
 21 much more constructive than the simple objective–subjective duality.

22 We do not think that this advice represents empty truisms of the 'mom and apple pie' variety.  
 23 In fact, we repeatedly encounter publications in top scientific journals that fall foul of these  
 24 virtues, which indicates to us that the underlying principles are subtle.

25 Instead of pointing at specific bad examples, here is a list of some common problems (dis-  
 26 cussed, for example, in Gelman (2015) and Gelman and Zelizer (2015)), where we believe that  
 27 exercising one or more of our listed virtues would improve matters:

- 29 (a) presenting analyses that are contingent on data without explaining the exploration and  
 30 selection process and without even acknowledging that it took place;
- 31 (b) justifying decisions by reference to specific literature without acknowledging that what  
 32 was cited may be controversial, not applicable in the given situation or without proper  
 33 justification in the cited literature as well (or not justifying the decisions at all);
- 34 (c) failure to reflect on whether model assumptions are reasonable in the given situation, what  
 35 effect it would have if they were violated or whether alternative models and approaches  
 36 could be reasonable as well;
- 37 (d) choosing methods because they do not require tuning or are automatic and therefore  
 38 seem 'objective' without discussing whether the methods chosen can handle the data  
 39 more appropriately in the given situation than alternative methods with tuning;
- 40 (e) choosing methods for the main reason that they 'do not require assumptions' without  
 41 realizing that every method is based on implicit assumptions about how to treat the data  
 42 appropriately, regardless of whether these are stated in terms of statistical models;
- 43 (f) choosing Bayesian priors without justification or explanation of what they mean and  
 44 imply;
- 45 (g) using non-standard methodology without justifying the deviation from standard ap-  
 46 proaches (where they exist);
- 47 (h) using standard approaches without discussion of their appropriateness in a specific  
 48 context.

1 Most of these are concerned with the unwillingness to admit to having made decisions, to justify  
2 them, and to take into account alternative possible views that may be equally reasonable. In  
3 some sense perhaps this can be justified on the basis of a sociological model of the scientific  
4 process in which each paper presents just one view, and then the different perspectives battle  
5 it out. But we think that this idea ignores the importance of communication and facilitating  
6 consensus for science. Scientists normally believe that each analysis aims at the truth and, if  
7 different analyses give different results, this is not because there are different conflicting truths  
8 but rather because different analysts have different aims, perspectives and access to different  
9 information. Letting the issue aside of whether it makes sense to talk of the existence of different  
10 truths or not, we see aiming at general agreement in free exchange as essential to science and,  
11 the more perspectives that are taken into account, the more the scientific process is supported.

12 We see the listed virtues as ideals which in practice cannot generally be fully achieved in any  
13 real project. For example, tracing all assumptions to observations and making them checkable  
14 by observable data is impossible because one can always ask whether and why results from the  
15 specific observations that are used should generalize to other times and other situations. As  
16 mentioned in Section 5.1, ultimately a rationale for treating different situations as ‘identical and  
17 independent’ or ‘exchangeable’ needs to be constructed by human thought (people may appeal  
18 to historical successes for justifying such idealizations, but this does not help much regarding  
19 specific applications). At some point—but, we hope, not too early—researchers must resort to  
20 somewhat arbitrary choices that can be justified only by logic or convention, if that.

21 And it is likewise unrealistic to suppose that we can capture all the relevant perspectives on any  
22 scientific problem. Nonetheless, we believe that it is useful to set these as goals which, in contrast  
23 with the inherently opposed concepts of ‘objectivity’ and ‘subjectivity’, can be approached  
24 together.

## 25 **Acknowledgements**

26 We thank the US National Science Foundation and Office of Naval Research for partial support  
27 of this work, and Sebastian Weber, Jay Kadane, Arthur Dempster, Michael Betancourt, Michael  
28 Zyphur, E. J. Wagenmakers, Deborah Mayo, James Berger, Prasanta Bandyopadhyay, Laurie  
29 Paul, Jan-Willem Romeijn, Gianluca Baio, Keith O’Rourke, Laurie Davies and the reviewers  
30 for helpful comments.  
31  
32  
33

## 34 **Appendix A: Objectivity in the philosophy of science**

35 Megill (1994) listed four basic senses of objectivity: ‘absolute objectivity’ in the sense of ‘representing the  
36 things as they really are’ (independently of an observer), ‘disciplinary objectivity’ referring to a consensus  
37 among experts within a discipline and highlighting the role of communication and negotiation, ‘procedural  
38 objectivity’ in the sense of following rules that are independent of the individual researcher, and ‘dialectical  
39 objectivity,’ referring to active human ‘objectification’ required to make phenomena communicable and  
40 measurable so that they can then be treated in an objective way so that different subjects can understand  
41 them in the same way. These ideas appear under various names in many places in the literature. Porter (1996)  
42 listed the ideal of impartiality of observers as another sense of objectivity. Douglas (2004) distinguished  
43 three modes of objectivity: human interaction with the world, individual thought processes and processes  
44 to reach an agreement. Daston and Galison (2007) called the ideal of scientific images that attempt to  
45 capture reality in an unmanipulated way ‘mechanical objectivity’ as opposed to ‘structural objectivity’,  
46 which refers to mathematical and logical structures. The latter emerged from the insight of scientists and  
47 philosophers such as Helmholtz and Poincare that observation of reality cannot exclude the observer and  
48 will never be as reliable and pure as ‘mechanical objectivists’ would hope.

More generally, virtually all senses of objectivity have been criticized at some point in history for being  
unachievable, which often prompted the postulation of new scientific virtues and new senses of objectivity.

For example, the realist ideal of ‘absolute objectivity’ has been branded as metaphysical, meaningless and illusory by positivists including Pearson (1911), and more contemporarily by empiricists such as van Fraassen (1980). The latter took observability and the ability of theory to account for observed facts as objective from an antirealist perspective.

Some writers even criticize the idea that objectivity is a generally desirable virtue in science, e.g. for its implication of a denial of the specific conditions of an observer’s point of view (Feyerabend, 1978; MacKinnon, 1987; Maturana, 1988) and its use as a rhetorical device or tool of power (see Fuchs (1997) for a critical overview of such ideas).

The core benefit of such controversies around objectivity and subjectivity for statisticians is the elaboration of aspects of good science, which should inform statistical data analysis and decision making. Hacking (2015) wrote a paper called ‘Let’s not talk about objectivity’, and with him we believe that, for discussing the practice of statistics (or more generally science), the objectivity *versus* subjectivity discourse should be replaced by looking at more specific virtues of scientific work, the awareness of which could have a more direct influence on the work of scientists. The virtues that we have listed in Section 3 are all connected either to senses of objectivity as summarized above, or to reasons for criticizing certain concepts of objectivity.

## References

- Agresti, A. and Coull, B. A. (1998) Approximate is better than exact for interval estimation of binomial proportions. *Am. Statistn*, **52**, 119–126.
- Alpert, M. and Raiffa, H. (1984) A progress report on the training of probability assessors. In *Judgment Under Uncertainty: Heuristics and Biases* (eds D. Kahneman, P. Slovic and A. Tversky), pp. 294–305. New York: Cambridge University Press.
- Bayarri, M. J., Berger, J. O., Forte, A. and Garcia-Donato, G. (2012) Criteria for Bayesian model choice with application to variable selection. *Ann. Statist.*, **40**, 1550–1577.
- Berger, J. O. (2006) The case for objective Bayesian analysis. *Bayesn Anal.*, **1**, 385–402.
- Berger, J. O., Bernardo, J. M. and Sun, D. (2009) The formal definition of reference priors. *Ann. Statist.*, **37**, 905–938.
- Bernardo, J. M. and Smith, A. F. M. (1994) *Bayesian Theory*. Chichester: Wiley.
- Berk, R., Brown, L., Buja, A., Zhang, K. and Zhao, L. (2013) Valid post-selection inference. *Ann. Statist.*, **41**, 802–837.
- Box, G. E. P. (1980) Sampling and Bayes’ inference in scientific modelling and robustness (with discussion). *J. R. Statist. Soc. A*, **143**, 383–430.
- Box, G. E. P. (1983) An apology for ecumenism in statistics. In *Scientific Inference, Data Analysis, and Robustness* (eds G. E. P. Box, T. Leonard and C. F. Wu), pp. 51–84. New York: Academic Press.
- Candler, J., Holder, H., Hosali, S., Payne, A. M., Tsang, T. and Vizard, P. (2011) Human rights measurement framework: prototype panels, indicator set and evidence base. *Research Report 81*. Equality and Human Rights Commission, Manchester.
- Chang, H. (2012) *Is Water H<sub>2</sub>O?: Evidence, Realism and Pluralism*. Dordrecht: Springer.
- Cox, R. T. (1961) *The Algebra of Probable Inference*. Baltimore: Johns Hopkins University Press.
- Cox, D. and Mayo, D. G. (2010) Objectivity and conditionality in frequentist inference. In *Error and Inference* (eds D. G. Mayo and A. Spanos), pp. 276–304. Cambridge: Cambridge University Press.
- Daston, L. and Galison, P. (2007) *Objectivity*. New York: Zone Books.
- Davies, P. L. (2014) *Data Analysis and Approximate Models*. Boca Raton: CRC Press.
- Dawid, A. P. (1982a) Intersubjective statistical models. In *Exchangeability in Probability and Statistics* (eds G. Koch and F. Spizichino), pp. 217–232. Amsterdam: North-Holland.
- Dawid, A. P. (1982b) The well-calibrated Bayesian. *J. Am. Statist. Ass.*, **77**, 605–610.
- Desrosieres, A. (2002) *The Politics of Large Numbers*. Boston: Harvard University Press.
- Dick, P. K. (1981) *VALIS*. New York: Bantam Books.
- Douglas, H. (2004) The irreducible complexity of objectivity. *Synthese*, **138**, 453–473.
- Douglas, H. (2009) *Science, Policy and the Value-free Ideal*. Pittsburgh: University of Pittsburgh Press.
- Erev, I., Wallsten, T. S. and Budescu, D. V. (1994) Simultaneous over- and underconfidence: the role of error in judgment processes. *Psychol. Rev.*, **101**, 519–527.
- Erikson, R. S., Panagopoulos, C. and Wlezien, C. (2004) Likely (and unlikely) voters and the assessment of campaign dynamics. *Publ. Opin. Q.*, **68**, 588–601.
- Everitt, B. S., Landau, S., Leese, M. and Stahl, D. (2011) *Cluster Analysis*, 5th edn. Chichester: Wiley.
- Feyerabend, P. (1978) *Science in a Free Society*. London: New Left Books.
- Fine, T. L. (1973) *Theories of Probability*. Waltham: Academic Press.
- de Finetti, B. (1974) *Theory of Probability*. New York: Wiley.
- Fisher, R. (1955) Statistical methods and scientific induction. *J. R. Statist. Soc. B*, **17**, 69–78.

- van Fraassen, B. (1980) *The Scientific Image*. Oxford: Oxford University Press.
- Fuchs, S. (1997) A sociological theory of objectivity. *Sci. Stud.*, **11**, 4–26.
- Gelman, A. (2003) A Bayesian formulation of exploratory data analysis and goodness-of-fit testing. *Int. Statist. Rev.*, **71**, 369–382.
- Gelman, A. (2008) The folk theorem of statistical computing. *Statistical modeling, causal inference, and social Science blog*, May 13th. (Available from <http://andrewgelman.com/2008/05/13/the.folk.theore/>.)
- Gelman, A. (2013) Whither the “bet on sparsity principle” in a nonsparse world? *Statistical modeling, causal inference, and social science blog*, Feb. 25th (Available from <http://andrewgelman.com/2013/12/16/whither-the-bet-on-sparsity-principle-in-a-nonsparse-world/>.)
- Gelman, A. (2014a) How do we choose our default methods? In *Past, Present, and Future of Statistical Science* (eds X. Lin, C. Genest, D. L. Banks, G. Molenberghs, D. W. Scott and J. L. Wang), pp. 293–301. London: Chapman and Hall.
- Gelman, A. (2014b) Basketball stats: don’t model the probability of win, model the expected score differential. *Statistical modeling, causal inference, and social science blog*, Feb 25th. (Available from <http://andrewgelman.com/2014/02/25/basketball-stats-dont-model-probability-win-model-expected-score-differential/>.)
- Gelman, A. (2014c) President of American Association of Buggy-Whip Manufacturers takes a strong stand against internal combustion engine, argues that the so-called “automobile” has “little grounding in theory” and that “results can vary widely based on the particular fuel that is used”. *Statistical modeling, causal inference, and social science blog*. (Available from <http://andrewgelman.com/2014/08/06/president-american-association-buggy-whip-manufacturers-takes-strong-stand-internal-combustion-engine-argues-called-automobile-little-grounding-theory/>.)
- Gelman, A. (2015) The connection between varying treatment effects and the crisis of unreplicable research: a Bayesian perspective. *J. Mangmnt*, **41**, 632–643.
- Gelman, A. and Basbøll, T. (2013) To throw away data: plagiarism as a statistical crime. *Am. Scient.*, **101**, 168–171.
- Gelman, A., Bois, F. Y. and Jiang, J. (1996) Physiological pharmacokinetic analysis using population modeling and informative prior distributions. *J. Am. Statist. Ass.*, **91**, 1400–1412.
- Gelman, A., and Carlin, J. B. (2014) Beyond power calculations: assessing Type S (sign) and Type M (magnitude) errors. *Perspect. Psychol. Sci.*, **9**, 641–651.
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D., Vehtari, A. and Rubin, D. B. (2013) *Bayesian Data Analysis*, 3rd edn. London: Chapman and Hall.
- Gelman, A., Goel, S., Rivers, D. and Rothschild, D. (2016) The mythical swing voter. *Q. J. Polit. Sci.*, **11**, 103–130.
- Gelman, A. and Loken, E. (2014) The statistical crisis in science. *Am. Scient.*, **102**, 460–465.
- Gelman, A. and O’Rourke, K. (2015) Convincing evidence. In *Roles, Trust, and Reputation in Social Media Knowledge Markets* (eds S. Matei and E. Bertino). New York: Springer.
- Gelman, A. and Shalizi, C. (2013). Philosophy and the practice of Bayesian statistics (with discussion). *Br. J. Math. Statist. Psychol.*, **66**, 8–80.
- Gelman, A. and Zelizer, A. (2015) Evidence on the deleterious impact of sustained use of polynomial regression on causal inference. *Res. Polit.*, **2**, 1–7.
- Gillies, D. (2000) *Philosophical Theories of Probability*. London: Routledge.
- Greenland, S. (2012) Transparency and disclosure, neutrality and balance: shared values or just shared words? *J. Epidem. Commty Hlth*, **66**, 967–970.
- Hacking, I. (2015) Let’s not talk about objectivity. In *Objectivity in Science* (eds F. Padovani *et al.*)
- Hennig, C. (2010). Mathematical models and reality: a constructivist perspective. *Foundns Sci.*, **15**, 29–48.
- Hennig, C. and Liao, T. F. (2013) How to find an appropriate clustering for mixed-type variables with application to socio-economic stratification (with discussion). *Appl. Statist.*, **62**, 309–369.
- Hennig, C. and Lin, C.-J. (2015) Flexible parametric bootstrap for testing homogeneity against clustering and assessing the number of clusters. *Statist. Comput.*, **25**, 821–833.
- Huber, P. J. and Ronchetti, E. M. (2009) *Robust Statistics*, 2nd edn. New York: Wiley.
- Jaynes, E. T. (2003) *Probability Theory: the Logic of Science*. Cambridge: Cambridge University Press.
- Kahneman, D. (1999) Objective happiness. In *Well-being: Foundations of Hedonic Psychology*, pp. 3–25. New York: Russell Sage Foundation Press.
- Kass, R. E. and Wasserman, L. (1996) The selection of prior distributions by formal rules. *J. Am. Statist. Ass.*, **91**, 1343–1370.
- Kendall, M. G. (1949) On the reconciliation of theories of probability. *Biometrika*, **36**, 101–116.
- Keynes, J. M. (1936) *The General Theory of Employment, Interest and Money*. London: Macmillan.
- Knight, F. H. (1921) *Risk, Uncertainty, and Profit*. Boston: Hart, Schaffner and Marx.
- Lewis, D. (1980) A subjectivist’s guide to objective chance. In *Studies in Inductive Logic and Probability*, vol. II (ed. R. C. Jeffrey), pp. 263–293. Berkeley: University of California Press.
- Linstone, H. A. (1989) Multiple perspectives: concept, applications, and user guidelines. *Syst. Pract.*, **2**, 307–331.
- Little, R. J. (2012) Calibrated Bayes, an alternative inferential paradigm for official statistics. *J. Off. Statist.*, **28**, 309–334.
- van Loo, H. M. and Romeijn, J. W. (2015) Psychiatric comorbidity: fact or artifact? *Theoret. Med. Bioeth.*, **36**, 41–60.

- 1 MacKinnon, C. (1987) *Feminism Unmodified*. Boston: Harvard University Press.
- 2 Maturana, H. R. (1988) Reality: the search for objectivity or the quest for a compelling argument. *Ir. J. Psychol.*,  
3 9, 25–82.
- 4 Mayo, D. G. (1996) *Error and the Growth of Experimental Knowledge*. Chicago: University of Chicago Press.
- 5 Mayo, D. G. and Spanos, A. (2010) The error-statistical philosophy. In *Error and Inference*, (eds D. G. Mayo and  
6 A. Spanos), pp. 15–27. New York: Cambridge University Press.
- 7 Megill, A. (1994) Four senses of objectivity. In *Rethinking Objectivity*, (ed. A., Megill), pp. 1–20. Durham: Duke  
8 University Press.
- 9 Merry, S. E. (2011) Measuring the world: indicators, human rights, and global governance. *Curr. Anthropol.*, **52**,  
10 S3, S83–S95.
- 11 von Mises, R. (1957) *Probability, Statistics and Truth*, 2nd edn. New York: Dover Publications.
- 12 Pearson, E. S. (1955) Statistical concepts in their relation to reality. *J. R. Statist. Soc. B*, **17**, 204–207.
- 13 Pearson, K. (1911) *The Grammar of Science*. New York: Cosimo.
- 14 Pollster.com (2004) Should pollsters weight by party identification? Pollster. (Available from [http://www.pollster.com/faq/should\\_pollsters\\_weight\\_by\\_party\\_identification.php](http://www.pollster.com/faq/should_pollsters_weight_by_party_identification.php).)
- 15 Porter, T. M. (1996) *Trust in Numbers: the Pursuit of Objectivity in Science and Public Life*. Princeton: Princeton  
16 University Press.
- 17 Rubin, D. B. (1984) Bayesianly justifiable and relevant frequency calculations for the applied statistician. *Ann.*  
18 *Statist.*, **12**, 1151–1172.
- 19 Sheiner, L. B. (1984) The population approach to pharmacokinetic data analysis: rationale and standard data  
20 analysis methods. *Drug Metablsm Rev.*, **15**, 153–171.
- 21 Silberzahn, R., et al. (2015) Crowdsourcing data analysis: do soccer referees give more red cards to dark skin  
22 toned players? Center for Open Science. (Available from <https://osf.io/j5v8f/>.)
- 23 Simmons, J., Nelson, L. and Simonsohn, U. (2011) False-positive psychology: undisclosed flexibility in data  
24 collection and analysis allow presenting anything as significant. *Psychol. Sci.*, **22**, 1359–1366.
- 25 Steegen, S., Tuerlinckx, F., Gelman, A. and Vanpaemel, W. (2016) Increasing transparency through a multiverse  
26 analysis. *Perspectives on Psychological Science*.
- 27 Tibshirani, R. J. (2014) In praise of sparsity and convexity. In *Past, Present, and Future of Statistical Science* (eds  
28 X. Lin, C. Genest, D. L. Banks, G. Molenberghs, D. W. Scott and J. L. Wang), pp. 505–513. London: Chapman  
29 and Hall.
- 30 Tukey, J. W. (1962) The future of data analysis. *Ann. Math. Statist.*, **33**, 1–67.
- 31 Vermunt, J. K. and Magidson, J. (2016) *Technical Guide for Latent GOLD 5.1: Basic, Advanced, and Syntax*.  
32 Belmont: Statistical Innovations.
- 33 Wang, W., Rothschild, D., Goel, S. and Gelman, A. (2015) Forecasting elections with non-representative polls.  
34 *Int. J. Forecast.*, **31**, 980–991.
- 35 Wasserman, L. (2006) Frequentist Bayes is objective (comment on articles by Berger and by Goldstein). *Baysn*  
36 *Anal.*, **1**, 451–456.
- 37 Weinberger, D. (2009) Transparency is the new objectivity. *Everything is miscellaneous blog*, July 19th. (Available  
38 from <http://www.everythingismiscellaneous.com/2009/07/19/transparency-is-the-new-objectivity/>.)
- 39 Yong, E. (2012) Nobel laureate challenges psychologists to clean up their act. *Nat. News*, Oct. 3rd. (Available  
40 from <http://www.nature.com/news/nobel-laureate-challenges-psychologists-to-clean-up-their-act-1.11535>.)
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- 48