

# Robust modeling in cognitive science

Michael D. Lee<sup>1</sup>, Amy Criss<sup>2</sup>, Berna Devezer<sup>3</sup>, Christopher Donkin<sup>4</sup>, Alexander Etz<sup>1</sup>, Fábio P. Leite<sup>5</sup>, Dora Matzke<sup>6</sup>, Jeffrey N. Rouder<sup>1</sup>, Jennifer S. Trueblood<sup>7</sup>, Corey N. White<sup>8</sup>, and Joachim Vandekerckhove<sup>1,†</sup>

This manuscript was compiled on March 19, 2019

In an attempt to increase the reliability of empirical findings, psychological scientists have recently proposed a number of changes in the practice of experimental psychology. Most current reform efforts have focused on the analysis of data and the reporting of findings for empirical studies. However, a large contingent of psychologists build models that explain psychological processes and test psychological theories using formal psychological models. Some, but not all, recommendations borne out of the broader reform movement bear upon the practice of behavioral or cognitive modeling. In this article, we consider which aspects of the current reform movement are relevant to psychological modelers, and we propose a number of techniques and practices aimed at making psychological modeling more transparent, trusted, and robust.

Cognitive modeling | Reproducibility | Open science | Robustness | Model comparison

You never want a serious crisis to go to waste . . . This crisis provides the opportunity for us to do things that you could not before.

Rahm Emmanuel, 1998

The field of psychology has recently questioned whether its findings are as reliable as they need to be to build a useful and cumulative body of knowledge. The growing lack of trust is sometimes called a “crisis of confidence” (Pashler & Wagenmakers, 2012, p. 528). A retrospective by Spellman (2015) identified a set of five causes for this crisis. A first rare but worrying culprit has been the manipulation and fabrication of empirical data (Simonsohn, 2013; Wagenmakers, 2012). A second more common problem has been the failure of established empirical findings to replicate in careful and systematic attempts (Alogna et al., 2014; Klein et al., 2014; Shanks et al., 2013; Open Science Collaboration, 2012). A third problem involves increasing recognition of the inherent but undisclosed flexibility in data collection and analysis, sometimes called “researcher degrees of freedom” (Simmons, Nelson, & Simonsohn, 2011). A closely-related fourth problem is the possibility of selective reporting and hypothesizing after empirical results are known, sometimes called “HARKing” (see Bones, 2012; Kerr, 1998; and Figure 1). Finally, Spellman (2015) noted the difficulties of obtaining other researchers’ data for re-analysis, verification, and conducting meta-analyses (Vanpaemel, Vermorgen, Deriemaeker, & Storms, 2015; Wicherts, Borsboom, Kats, & Molenaar, 2006).

In reaction to the crisis in confidence, there has been an effort

to identify and enforce good practices for analysis and reporting of experimental data. The practice of pre-specifying data collection and analysis plans, long required in clinical trials, has been proposed in psychology to limit both HARKing and undisclosed flexibility. This practice has become collectively known in psychology as *preregistration* (e.g., Matzke et al., 2015; Munafò et al., 2017; Nosek, Ebersole, DeHaven, & Mellor, 2018; Wagenmakers, Wetzels, Borsboom, van der Maas, & Kievit, 2012). A new publication format known as *registered reports* has been adopted by more than 100 psychology journals as a way to incorporate these ideas directly into the research and publication pipeline<sup>1</sup> (Chambers, 2013; Chambers, Dienes, McIntosh, Rotshtein, & Willmes, 2015; Hardwicke & Ioannidis, 2018). Psychologists have also recognized the importance of replication as a tool for verifying scientific claims (Open Science Collaboration, 2012, 2015), and vigorously debated what role replication plays in a healthy science (see Zwaan, Etz, Lucas, & Donnellan, 2018, and its associated commentaries). In addition, psychologists have pushed for open data, open code, and open materials to allow for better verification and reanalysis of study results. For example, the Transparency and Openness Promotion (TOP) guidelines (Miguel et al., 2014; Nosek et al., 2015) is a collection of eight key open science practices structured into three levels of increasing stringency. The TOP guidelines have been implemented by more than 5,000 scientific organizations and more than 1,000 journals spanning many scientific disciplines.

**The crisis of confidence reaches beyond experimental psychology.** The focus of the crisis of confidence in psychology has been in experimental psychology, involving the analysis of empirical data using standard statistical methods. Often, however, psychology seeks to understand its data and theories using models (Farrell & Lewandowsky, 2018; Sun, 2008). Modeling and model-based inference is closely related to the standard statistical data analysis routinely used in experimental psychology. While familiar data analysis methods like regression and analysis of variance are often thought of as procedures, they can also be thought of as statistical models used to perform inference.

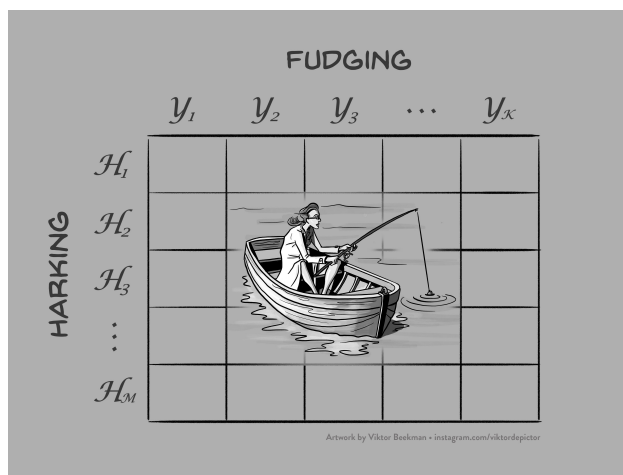
From this perspective, the only difference between statistical

<sup>1</sup>Department of Cognitive Sciences, University of California, Irvine; <sup>2</sup>Department of Psychology, Syracuse University; <sup>3</sup>Department of Business, University of Idaho; <sup>4</sup>School of Psychology, University of New South Wales; <sup>5</sup>Department of Psychology, The Ohio State University at Lima; <sup>6</sup>Department of Psychology, University of Amsterdam; <sup>7</sup>Department of Psychology, Vanderbilt University; <sup>8</sup>Department of Psychology, Missouri Western State University

<sup>†</sup>To whom correspondence should be addressed. E-mail: joachim@uci.edu.

This article is the product of the *Workshop on Robust Social Science* held in St. Petersburg, FL, in June 2018. The workshop was made possible by generous funding from the National Science Foundation (grant #BCS-1754205) to Joachim Vandekerckhove and Michael Lee of the University of California, Irvine. Alexander Etz was supported by NSF GRFP #DGE-1321846. Berna Devezer was supported by NIGMS of the NIH under award #P20GM104420. Dora Matzke was supported by a Veni grant (#451-15-010) from the Netherlands Organization of Scientific Research (NWO). Jennifer Trueblood was supported by NSF #SES-1556325. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health or the National Science Foundation.

<sup>1</sup>Note that registered reports involve more than preregistration: They also involve a journal’s guarantee that a paper will be published regardless of how the data turn out.



**Fig. 1.** A caricature of questionable research practices, conceived as fishing for research results by the manipulation of data (“fudging”) and selective consideration of hypotheses (“HARKing”). Artwork by Viktor Beekman based on a concept by Eric-Jan Wagenmakers. Reproduced under a CC-BY license. Source: <https://www.bayesianspectacles.org/library/>.

analysis and psychological modeling lies in the emphasis that psychological models place on substantive interpretation. The data-generating mechanisms in a psychological model can usually be interpreted in terms of psychological processes, such as storing an item in memory, attending to a stimulus feature, or making a decision. The parameters in a psychological model can usually be interpreted as unobservable psychological constructs governing the processes, such as the capacity of working memory, the level of selective attention, or the bias affecting a decision.

Such “psychology-descriptive” models can provide substantive insights, unlike some machine learning and statistical models that focus exclusively on prediction. They can broaden the types of experimental data that can be analyzed, including small data sets, and non-standard experimental designs. Finally, they can provide stronger tests of competing psychological theories, because the models correspond more closely to the theories and formalize more of the assumptions made by the theories (Vanpaemel, 2010).

The close relationship between standard data analysis and model-based analysis suggests that the critical re-examination of data analysis methods in psychology also has ramifications for psychological modeling. Accordingly, the goal of this article is to consider how the lessons learned from the crisis in experimental psychology could improve modeling practices in psychology and cognitive science. In our discussion, we divide good modeling practices into three general parts: those that apply before data have been collected, those that apply after data have been collected, and a set of general good practices throughout model-based research. We consider each of these parts in turn, illustrating the general issues with specific examples drawn from various sub-fields of psychological modeling. We conclude with a brief discussion of techniques to make psychological modeling more transparent, trusted, and robust.

### Good practices before data are collected

**Preregistering models, the players in the game.** Preregistering models and their predictions can be a useful scientific practice. One way to think of the practical benefits of preregistrations is that it can help a researcher in much the same way that preregister-

ing a dissertation research plan can help a graduate student. It provides an explicit and detailed plan of action at the beginning of the enterprise. Preregistration is not intended as a constraint on what *can* happen, and will generally not anticipate everything that *could* happen. The preregistration does, however, provide a clear statement of the motivating goals for the research, and the intended ways in which those goals will be met.

Preregistration is especially important in a confirmatory research setting, in which data are used to evaluate the adequacy of a model or to compare multiple models. As part of such a preregistration, it is important to be clear about what are core versus ancillary modeling assumptions, and how these relate to the research questions. Core assumptions are those that motivated the empirical test and will usually correspond to the major theoretical questions being addressed by the research. Ancillary assumptions involve various possible choices to non-core parts of the model. An example of this distinction is provided by the “Expected Utility Theory” case-study box.

Ideally, a preregistered model could take the form of the precise predictions that are made by the model. Bayesian methods, by requiring both a likelihood and a prior, automatically make comprehensive predictions about data, but it is usually possible to preregister some predictions using non-Bayesian methods as well. In addition, in most research situations involving model comparison, there are many possible models that could be included. As for the methods of analysis, if model comparison is to be used, it is important to preregister the models that will be compared to one another. This prevents “changing the players in the game” once data have been seen. That is, it prevents one from introducing a model that performs poorly on the data, to give the impression that the originally-proposed set of models fare relatively well, or from introducing a model that was directly inspired by the data, which will perform well on the data sample but may generalize poorly due to over-fitting.

It is important to emphasize that preregistration is not necessarily needed at all stages of the modeling process. A large part of developing psychological models is exploratory in nature. It may be more useful, in many cases, to engage in a practice we call *postregistration*. Postregistration involves keeping a comprehensive log of the modeling process that acts as an “activity log” documenting the process of model building and checking. This concept is discussed in the section “Exploratory analyses and postregistration” below.

### Preregistering evaluation criteria, the rules of the game.

There are usually many ways in which we can evaluate a model against data. All of the following metrics are used regularly to evaluate models:  $p$ -values, correlation coefficients, variance explained measures, sum of squared error, mean absolute deviation, proportion of agreement, maximum likelihood, normalized maximum likelihood, information criteria (AIC, BIC, DIC, WAIC, etc.), and Bayes factors (Shiffrin, Lee, Kim, & Wagenmakers, 2008). Many of these measures are similar to one another, or are exactly equivalent in special cases. In other cases, these metrics may give different and opposite results for a key research question for exactly the same models and data.

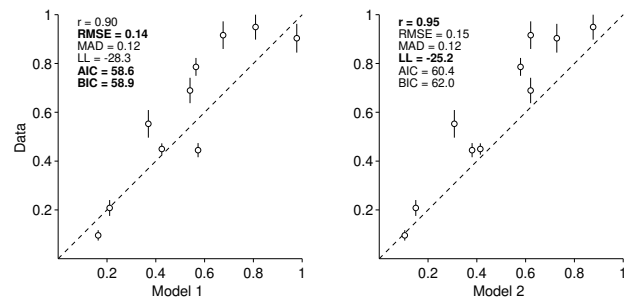
An illustrative example is given in Figure 2. Both panels in that figure display an artificial data set with 10 conditions. Each condition is characterized by its sample mean and sample standard deviation on the vertical axis. The horizontal axis depicts the predictions made by two fictional models, so that any deviation from the diagonal indicates the degree of model misfit. Various

### Example 1: Expected Utility Theory

As a concrete example of the difference between core and non-core assumptions, consider an example from psychological models of choice based on expected utility theory. The famous Allais (1953, 1979) paradox presents two problems involving choices between a “safe” gamble and a “risky” gamble. The problems are designed such that, according to expected utility theory, a decision maker should either choose the safe option for both problems, or the risky option for both problems. A core assumption of expected utility theory is that decision makers have a stable preference state for the safe or risky option that applies to both gambles. This leads to the prediction of the theory that a mixed response—choosing the safe option in one problem and the risky option in the other problem—cannot happen at the individual level. Observed behavioral data typically has at least some violations of this strict prediction, which is attributed to some form of error in the individuals’ decision-making processes (Birnbau & Quispe-Torreblanca, 2018). Assumptions about errors are good examples of ancillary assumptions. They are also a good example of the level of modeling detail typically needed to make complete predictions. Specifying how likely it is that errors will be made, and how frequent those errors could be, transforms the modeling predictions from a qualitative one of “the theory predicts this will not happen” to specific predictions about how many people will produce each of the possible types of behavioral patterns in an experiment. In this context, preregistration would be appropriate to test whether or not individuals show safe and risky options. A subsequent exploratory modeling exercise could be to develop a mechanism that describes the errors and when they occur.

summaries of the goodness-of-fit are also listed for each model. In terms of the product-moment correlation ( $r$ ) Model 2, in the right panel, is preferred. In terms of the root mean square error measure (RMSE) Model 1, in the left panel, is preferred. The models are very close in mean absolute deviation (MAD) measures. In terms of the log-likelihood Model 2 is preferred. Other metrics require an account of the complexity of the models. For the purposes of our example, assume that Model 1 has one parameter, and Model 2 has five parameters. Using this information, the Akaike Information Criterion (AIC) and Bayesian Information Criterion (BIC) both prefer Model 1. The difference based on the BIC, however, is larger, which may lead to stronger claims (Wagenmakers & Farrell, 2004). Which metric to choose in practical applications is a challenging statistical and methodological question that remains an active area of debate and research throughout the empirical sciences and statistics (Myung, Forster, & Browne, 2000; Navarro, in press; Shiffrin et al., 2008).

The purpose of the example in Figure 2 is to illustrate how even in very common situations different reasonable and widely-used metrics can suggest conflicting conclusions. This leads us to our recommendation that researchers preregister the methods of evaluation that will be used. Such preregistration prevents “changing the rules of the game”—whether intentionally or not—once data have been seen. A good preregistration should also provide an argument for the suitability of the chosen metric in terms of the relevant statistical and methodological considerations. Preregistration notwithstanding, once the data are collected it remains important to evaluate whether any assumptions of the analyses or model



**Fig. 2.** Example of different results provided by different metrics for assessing the adequacy of models. The two panels show the goodness-of-fit of two models to the same data. Each panel lists the adequacy of the model as measured by the product-moment correlation ( $r$ ), the root mean square error (RMSE), the mean absolute deviation (MAD), the log likelihood (LL), the Akaike Information Criterion (AIC), and the Bayesian Information Criterion (BIC).

comparison tools are violated. Preregistration is no substitute for good judgment and care must be taken not to fixate on the results of one model comparison metric for no other reason than that it was preregistered. Changing the model comparison metric after the data are seen might be advisable if it turns out that the data unexpectedly violate some assumption of the proposed analysis and selection criterion. For example, if a researcher who preregistered an analysis using a linear model and AIC discovers that their data exhibit large interindividual differences, it would be defensible for them to switch to an hierarchical model and Bayes factor.

**Registered modeling reports.** The advantage, especially in confirmatory research settings, of preregistering both models and their method of analysis, suggests the desirability of making these declarations in a systematic and routine way. Accordingly, we propose a new article format for modelers called *Registered Modeling Reports*, analogous to registered reports for experimental studies (Chambers et al., 2015; Hardwicke & Ioannidis, 2018). In a Registered Modeling Report, researchers pre-specify models, data collection mechanisms (whether experimental or observational), and analyses prior to data collection; then they write up the Introduction and Methods sections—and, if possible and relevant, the computer code to be used for model and data analysis—of an article for a “Stage 1” submission. The role of peer review then is to assess whether these specifications are principled and sufficient, and whether the study design and planned analyses are of the desired quality. A report that meets these criteria can be provisionally accepted for publication, contingent on the researchers following through with the registered model and methodology. After data collection, the authors complete the manuscript with a “Preregistered Results” section, an optional “Exploratory Analysis” section, and a “Discussion” section.

We propose the Registered Modeling Report format in order to insert the ideas of preregistration and transparency of modeling practices directly into the publication process. This potentially has a number of advantages. Most importantly, Registered Modeling Reports may help improve the research itself, by allowing reviewers to provide their expertise on experimental design and analysis before resources are invested in the collection of data. Registered Modeling Reports may also help streamline the review process by preventing reviewers arguing for additional models or analyses once the data have been collected.

While we think preregistration and Registered Modeling Reports are important new ideas in model-based inference, the limits of their



role should be understood. First, it is clear that not every modeling study is suited to preregistration or Registered Modeling Reports. Successful cognitive modeling almost always requires fine-grained iterative model development, and this process is not well-matched to a single preregistration or Registered Modeling Report. The iterative process of model development and evaluation may well be better documented during the research progress, and disseminated as part of postregistration. Secondly, the preregistration of a model and the way it will be used may not survive contact with data. Violations of protocol should be documented, and deviating from one's preregistration plans should be permitted. Preregistration should not prevent learning from the data at hand, and it does not prevent carrying out valuable exploratory analysis.

**What does not carry over.** On the other hand, some of the recommendations that have followed the crisis of confidence do *not* carry over to the standard practices of cognitive modelers. For example, unless a modeling project depends critically on a null hypothesis test and there is only one opportunity for data to be collected, a priori power analysis does not serve a necessary role. Similarly, since model construction does not necessarily rely on the availability of large data sets, we do not believe rules of thumb for sample sizes are useful considerations, especially if applied post hoc.

### Good practices after data are collected

**Utilities in model evaluation.** As empirical fields collect more phenomena, they sometimes develop checklists or benchmarks of qualitative properties that a good model should have. For example, Oberauer et al. (2018) present a set of phenomenological properties relevant to working memory research, and Epper and Fehr-Duda (2018) present a set of seven regularities involving risk taking and time discounting behavior. The “Context Effects in Decision Making” case-study box provides another concrete example of a checklist. These sorts of checklists—that characterize models in terms of discrete phenomena that are either present or absent, with no strong statements about their magnitude—encourage a falsificationist (Popper, 1959) perspective on model evaluation. While this may be appropriate, it is easy for checklists to miss or mischaracterize important aspects of empirical phenomena, and so provide incomplete or inappropriate benchmarks. For example, a checklist may neglect the role of individual differences, or ignore the joint prediction of other relevant behavioral data. Hence, for such multidimensional data it is appropriate to confirm the joint occurrence of the phenomena, and to consider their sensitivity to individual differences in the theory-building stage.

A checklist of observed phenomena may also set up inappropriate expectations for future model building. If the evidence for some phenomena on a list is in fact weak, but models are constructed with extra complexity to account for those phenomena all the same, the model is essentially overfit, meaning that it is overly attuned to the specific features of the previous data sets and will hence suffer in generalization tests. A conventional method of safeguarding against overfitting would be to use a model comparison metric that balances the quantitative fit of a model against the complexity of the model. However, it is less clear how to compute quantitative fit for a checklist of phenomena (but see Pitt, Kim, Navarro, & Myung, 2006).

Another standard modeling practice is to summarize the ability of a model to capture patterns in the data through an omnibus measure of goodness-of-fit, such as those considered in Figure 2.

### Example 2: Context Effects in Decision Making

Decades of research have been devoted to understanding the cognitive processes that give rise to context-sensitive behavior when people make decisions between multiple alternatives described by multiple attributes. The focus of this work has been on understanding three classic context effects: the *attraction* (Huber, Payne, & Puto, 1982), *compromise* (Simonson, 1989), and *similarity* effects (Tversky, 1972). These effects describe how preferences between two alternatives can change with the introduction of a new third alternative, and have been shown to occur in adults, children, monkeys, honeybees, hummingbirds, and even slime molds (Latty & Beekman, 2010). The effects are theoretically important because they challenge classical utility models of decision making (Luce, 1959) by showing that the relative preference for two options often depends on the utility of another “decoy” option. There are at least a dozen different computational models of these effects, with model evaluation focused on exploring the range of parameter values that can qualitatively produce the three effects. Typically, the modeling goal is to find a single set of parameters that can account for all three effects (e.g., Roe, Busemeyer, & Townsend, 2001; Usher & McClelland, 2004). However, the three effects are too fragile and subtle to serve as simple mandated checklists in this way. Very few participants produce all three effects within a single experiment, even though most people show the effects in isolation (Trueblood, Brown, & Heathcote, 2015). There are also large individual differences in the strength and co-occurrence of the effects (Liew, Howe, & Little, 2016). Thus, simply relying on a checklist of effects misses important aspects of the psychological phenomena being explored, and mischaracterizes the behavior for which an explanation is sought. A better approach is to evaluate detailed cognitive models of the decision processes involved, testing the accuracy of their predictions about the individual decisions that individual people make on these tasks (Evans, Holmes, & Trueblood, 2018; Turner, Schley, Muller, & Tsetsos, 2017).

A limitation of such an approach is that a single quantity may fail to capture the full richness of information that the data provide for evaluating a model, or the omnibus measure may be led astray by small blips in the data and uninteresting violations of ancillary assumptions.

One potential way to overcome the limitations of qualitative checklists and overly sensitive omnibus fit measures is to consider them as two end-points on a continuum of utilities (i.e., cost functions) for “scoring” a model against data. Checklists operate at a coarse resolution, measuring utility in terms of a few features, while fit gives consideration to every data point. Between these extremes lie utility functions that emphasize key qualitative points of comparison between models and data, but continue to consider all of the data in a fine-grained way. The “Absolute Identification” case-study box gives a concrete example from psychophysics of balancing qualitative and quantitative agreement between models and data.

The more widespread use of utility functions has the potential to strike an appealing balance between giving weight to qualitatively important data patterns, while still measuring overall quantitative agreement. Utilities to be used in confirmatory model evaluation become part of the “rules of the game” and as such should be determined—and preregistered—before the analysis is performed.

### Example 3: Absolute Identification

Despite its apparent simplicity, the task of attaching labels to stimuli varying on a single dimension (e.g., tones varying in pitch) reveals a number of fundamental limits in human information processing. For example, stimuli closer to the extremes of the range of stimuli presented are more accurately labelled than those in the middle. This well-replicated result, now called the *bow* effect, has become a benchmark finding that most researchers would agree any model of absolute identification must capture (Murdock, 1960). Similarly, models are expected to capture the relatively complex, but reliable, observed pattern of sequential effects typically observed in absolute identification studies. People show an assimilation effect which attracts their response towards the stimulus presented on the previous trial (Garner, 1953), but a weaker contrast effect that repels them from previously presented stimuli (Holland & Lockhead, 1968). Such benchmark findings have been used to evaluate a number of psychological models (e.g., S. D. Brown, Marley, Donkin, & Heathcote, 2008; Stewart, Brown, & Chater, 2005). Stewart et al. (2005) take stock of the benchmark effects observed within the absolute identification literature, asking whether different models can produce the same qualitative patterns as observed in data. Such an approach can be very useful, but they (and most others) also evaluate the quantitative agreement between the observed data and the predictions of models. Once the required qualitative properties have been established, quantitative fit remains important because of the additional challenge it poses to models, especially when capturing data at the level of individual participants.

**Bookend models.** Model comparison is inherently relative, and there is no way to measure absolute adequacy. For an example of the dilemmas this can cause, consider the case in which two models are proposed for a complex data set. Suppose that both models clearly fail in important ways, but it is still the case that one model decisively outperforms the other. One reasonable perspective in this situation is that “the second model is better than the first, but both are terrible, so what do we learn from the comparison?” An alternative reasonable perspective is that “the fact that one model is better means some additional insight might have been gained, and that can aid future development.” At a minimum, the worse performing model can now be rejected even though a good alternative remains undiscovered.

One practical way to address this dilemma is to include, when feasible, additional base-rate and catch-all models as “bookends.” This involves augmenting the set of models under consideration to include models that are much more parsimonious, and some that are much more complicated, than the substantive models of interest. If a model of interest outperforms the bookend models, this suggests that its success in accounting for the data does not come entirely from its parsimony or goodness-of-fit alone, but from striking a suitable balance. In this sense, comparison to bookend models provides a practical proxy for the assessment of absolute model adequacy. The “Memory Retention Functions” case-study box gives a concrete example of this use of bookend models.

**Prediction and generalization.** Beyond descriptive adequacy, prediction and generalization are important additional approaches for model evaluation. By prediction, we mean tests that measure the success of a model in accounting for unobserved data from

### Example 4: Memory Retention Functions

Models of memory retention characterize the change in probability of recall of an item or episode from memory as a function of time. Many functions, including various power functions and exponential decay functions, have been proposed to model this relationship (e.g., Rubin, Hinton, & Wenzel, 1999). The bookend approach would add to this set of serious theoretical competitors something like a null model that assumed memory for items was constant with respect to time, and a saturated model that allowed a free parameter for the probability of recall at every measured time point. The null model is presumably far too restrictive, and will under-fit the data. The saturated model is presumably far too complex, and will over-fit the data. Thus, for a theoretical model like a power or exponential function to be a serious contender, it should out-perform both of these bookends.

the same task. By generalization, we mean tests that measure the success of a model in accounting for unobserved data from a different but related task. This difference between prediction and generalization is emphasized by Busemeyer and Wang (2000), who argue for the merits of generalization tests. Successful prediction and generalization show a model to be robust, in the sense that the model is not over-emphasizing any idiosyncratic features of a particular set of data.

Some practices for prediction tests are well established, like cross-validation and accumulative prediction error (Shiffrin et al., 2008; Wagenmakers, Grünwald, & Steyvers, 2006). There are fewer examples of generalization tests in cognitive modeling (but see Criss, Malmberg, & Shiffrin, 2011, Guan, Lee, & Vandekerckhove, 2015, and Kılıç, Criss, Malmberg, & Shiffrin, 2017, for some recent examples). The “Serial Position Effects in Free Recall” case-study box provides one concrete example. Generalization tests should become more widespread as psychological modeling aims to demonstrate its robustness. The ability to make accurate predictions about what will happen in new and different psychological circumstances is a compelling way to demonstrate the explanatory power and range of applicability of a theory.

Other approaches related to prediction and generalization are emerging in contexts like machine-learning competitions. An early example was the Netflix competition (Bell, Koren, & Volinsky, 2010), which provided data on the ratings viewers gave to movies they watched, and tested the ability of algorithms to predict the ratings for withheld data. The final million-dollar prize was awarded for the first algorithm able to make a 10% improvement over Netflix’s own recommendation algorithm at the time. For the most part, the algorithms submitted to the competition were developed using statistical and machine-learning methods. We believe there should be a role for cognitive models in this competitive context. In the case of the Netflix competition, teams predicted human judgments of aesthetic stimuli – a quintessentially behavioral question common in the cognitive sciences. Contemporary data science competitions might benefit from, as components of a good entry, some measure of psychological modeling. For example, recent Kaggle competitions<sup>2</sup> include the “Dog Breed Identification Challenge” (determine the breed of a dog from an image), the “Toxic Comment Classification Challenge” (identify and classify inappropriate comments in an online setting), and the “Store Item Demand Forecasting Challenge” (predict three months of item sales at different stores). These Kag-

<sup>2</sup><https://www.kaggle.com/competitions>

### Example 5: Serial Position Effects in Free Recall

The serial position curve in free recall is one of the most robust findings in the study of memory. Items presented near the beginning and end of studied lists tend to be recalled better than items in the middle of the list, when there is no requirement to recall the items in order. The quantitative details of this qualitative regularity, however, depend on details of the task, including how many items are studied and the time interval between their presentation (Murdock, 1962). Shiffrin et al. (2008) present a case study of how a model of free recall can be evaluated in terms of its ability to generalize across these conditions. They focus on a hierarchical extension of the SIMPLE model (G. D. A. Brown, Neath, & Chater, 2007), using data from six experimental conditions to infer model parameters. The extension enabled predictions about the appropriate parameterization of the model in three extra conditions involving set sizes and presentation intervals that were different from those used to make the inferences. These predicted parameters, in turn, were used to generate the serial position curves that the model predicts. In effect, the extension served to broaden the model's account of free recall from the observed experimental tasks to new experimental tasks, so that evaluation against data from the new tasks would provide a strong test of the model.

gle competitions seem likely to benefit from psychological models of vision, language, and decision making, respectively. There are also contemporary competitions, such as the Choice Prediction Competition<sup>3</sup>, that are more explicitly focused on psychological theories and cognitive modeling challenges. A key element of all of these competitions is that the requirement of genuine prediction has to be carefully implemented in the assessment of competing models.

**Exploratory analyses and postregistration.** However a model is evaluated, the evaluation should ideally be augmented with exploratory analyses. Perhaps the most common exploratory analysis involves the discussion of model misspecification. All models are misspecified, and modelers often work through a sequence of models before arriving at the one ultimately presented – trying this functional form and that, allowing for individual differences or trial-to-trial effects, adding auxiliary assumptions for a new study design, and taking what was an auxiliary assumption and building it into a core assumption. Understanding the reasons for steps taken and the nature of the residual misspecification provides crucial information for guiding future model development that often goes unreported. In other words, knowing what did not work in model development and what still does not work in the final model, should be transparently reported to the field.

Model development is a creative activity that often proceeds in this incremental and exploratory fashion. A model is forged from data through a process of abductive reasoning, and it undergoes multiple cycles of empirical testing and adjustment over time.

In exploratory model development, we believe there is a useful expanded role for what we call “postregistration” documentation. Postregistration is part of an ongoing research effort and involves keeping an “activity log” documenting every model alteration tried during the study. This type of activity log is essentially a modeling lab notebook, not unlike a traditional lab notebook (Noble, 2009),

which is updated incrementally as the exploratory modeling proceeds. Modeling notebooks can be created using existing software tools such as Jupyter or Rmarkdown, and they can be made public at the time of publication or even as the research is being done. In a preregistered confirmatory setting, postregistration could focus on non-core modeling results, possibly in the published article, but possibly only in supplementary material. In either case, postregistration provides a mechanism for avoiding the modeling file-drawer effect, in which attempts at model development that fail are never made public (Rosenthal, 1979). The overarching aim of these additional reporting considerations is to inform the field, and to enhance the understanding of the model.

### Good practices throughout psychological modeling

Modelers should always endeavor to make their models available (Baumgaertner, Devezzer, Buzbas, & Nardin, 2018). The motivating goal of ensuring availability is to preserve the rights of others to reach independent conclusions about model-based inferences. A minimum standard, then, is to provide accessible modeling details that allow a competent person in the field to reproduce the results. This is likely to include mathematical and statistical description, an algorithm or pseudo-code, user documentation, and so on. Providing these details in a sufficiently precise form makes a model available, and means it is likely to be used and understood more broadly than by a specific researcher or a single lab.

**Making modeling robust.** Stephen Jay Gould (1996) pointed out that mistakes that favor a researcher's preferred conclusions tend to go uninvestigated, and so tend to remain (Rouder, Haaf, & Snyder, in press). A consequence of this “psychology of errors” is that mistakes in model implementation tend to be biased in favor of the model – that is, the results are not robust to who is performing the analysis. In computer science there are established “robust coding” techniques that can help researchers address their unconscious biases, including independent implementation and test-driven development (Beck, 2003). Nevertheless, making modeling robust to error is a challenging task. It is a special case of the more general challenge of establishing the robustness of model-based findings. We discuss one established modeling practice for increasing robustness in modeling, as well as two more recent ideas.

One established practice involves parameter recovery studies. These studies test the correctness of the computational implementation of a model through recovery studies that fit a model to data simulated from that model (Cook, Gelman, & Rubin, 2006; Heathcote, Brown, & Wagenmakers, 2015). The usual assumption is that it is desirable for a model to infer the parameter values that are known to have generated the simulated data. Model recovery studies can provide a way to understand the properties of a model. In particular, they can help diagnose issues like (weak) identifiability with respect to the type and amount of information likely to be available. These diagnoses in turn can help guide the decisions involved in experimental design. An extreme form of using models to guide experimental design involves the growing area of “optimal experimental design,” in which the predictions made by competing models are used to choose the conditions presented in an experiment, or even the stimuli presented on a trial-by-trial basis (Cavagnaro, Pitt, Gonzalez, & Myung, 2013; Zhang & Lee, 2010).

A more recent idea for combating errors involves the use of *blinding* in modeling (Dutilh et al., 2017; MacCoun & Perlmutter,

<sup>3</sup><https://cpc-18.com/>



### Example 6: The Worst Performance Rule

People with higher working memory capacity tend to respond relatively more quickly in elementary perceptual tasks, such as deciding whether a stimulus array contains more white or black dots (Jensen, 2006). According to the worst performance rule, the worst performance in these simple tasks is more predictive of high-order cognitive ability than best performance (Baumeister & Kellas, 1968). The evidence for this rule is that higher response time quantiles (i.e., slower responses) correlate more strongly with working memory capacity than lower response time quantiles (e.g., Unsworth, Redick, Lakey, & Young, 2010). Ratcliff, Schmiedek, and McKoon (2008) have argued that the worst performance rule can be explained by a drift diffusion model of the time course of making simple decisions. In particular, the diffusion model account of the worst performance rule posits that the same general processing speed—the “drift rate” parameter—facilitates performance in both simple perceptual tasks and complex cognitive tasks. Dutilh et al. (2017) tested these modeling claims using a confirmatory yet flexible two-stage modeling strategy. In the first stage, the modeler was provided with the choice response times and a randomly shuffled version of the working memory measurements. In this way, the critical correlation between person-specific working memory scores and person-specific choice response times was withheld from the modeler. After all of the modeling decisions were made, such as outlier exclusion and transformations, the analysis code was made public on the Open Science Framework. In the second stage, the working memory scores were unshuffled, and the analysis script was applied to the data to evaluate the preregistered hypotheses. The authors concluded that the results provided evidence against the worst performance rule. By blinding the modeler to the mapping between the key variables, in the form of the drift rate parameters and working memory scores, the two-stage analysis strategy allowed for flexibility in modeling while ensuring that the modeling decisions were not driven by expectations about the outcomes.

ter, 2017). The core requirement of blinded modeling is that the modeler is provided with most of the data, but that the data are scrambled or delabeled to make it impossible to determine if the outcome is desirable or undesirable with respect to a theory or model. Blinded modeling alleviates the psychology-of-errors bias, and thus provides a mechanism to increase confidence in the usefulness of model-based inference. In particular, it provides a strong test of selective influence (Voss, Rothermund, & Voss, 2004). “The Worst Performance Rule” case-study box provides an example of using blinding.

A model-based way of guarding against errors involves testing the robustness of results to small variations in the model definition. Most modeling applications in psychology involve using only one model to make inferences from data. It is the case, however, that the most important conclusions should be robust to the non-core details of the model. In the same way that we test the sensitivity of our conclusions to irrelevant variations in the priors, we can test their sensitivity to irrelevant variations in the likelihood (Farrell & Lewandowsky, 2018; Lee, 2018). This practice is sometimes called *likelihood profiling*. The “Predator Avoidance and Courtship in Butterflies” case-study box provides an example in which both priors and likelihoods are tested for robustness.

### Example 7: Predator Avoidance and Courtship in Butterflies

Finkbeiner, Briscoe, and Reed (2014) model approach-and-avoidance behavior (in predatory situations) and courtship behavior (in mating situations) in butterflies. They use Bayesian methods to implement models and evaluate them against the behavioral data. As part of testing the robustness of the modeling conclusions, they examined a set of variations on the original model, particularly with respect to the modeling assumption made about individual differences between butterflies. Different prior distributions on the level of variability are systematically tested, together with different assumptions about the shape of the distribution that characterizes individual differences. The observation that the important modeling results are robust to these changes suggests that they come from the data and core theoretical commitments of the model, rather than from the more arbitrary ancillary assumptions.

A more extreme version of this robust modeling approach uses multiple different models that formalize the same psychological theory (Dutilh et al., 2018). Analogously to the “many analysts” approach in data analysis, the goal of this approach is to test the variation in findings arising from different researchers tackling the same problem using different reasonable methods (Silberzahn et al., 2018). If different models converge on the same findings, it suggests the models capture the theory and the inferences are robust. If the results do not agree, rich diagnostic information is provided to investigate the models and psychological phenomena involved. We believe this crowd-sourced approach to evaluating the robustness of findings is an important emerging capability, facilitated by the increasing speed and ease of distributed scientific interaction.

**More complete modeling.** Exploratory and confirmatory methods can and should be used at the same time, for the same research question, and even for the same model and data. The exploratory part of a confirmatory study allows the data to inspire further model development. The exploratory evidence provided by current data can be measured for any model, including one provided by the data. Results of exploratory analysis then inspire future confirmatory tests using independent data. It is critical, however, that exploratory evidence should not be misinterpreted as confirmatory evidence if the model was not anticipated before the data were seen.

We believe that one useful way to think of the distinction is in Bayesian terms. The canonical Bayesian approach to model selection is based on the Bayes factor, which is a ratio expressing how much evidence the data provide for one model over another (Kass & Raftery, 1995; Lee & Wagenmakers, 2013, Chapter 7; Vandekerckhove, Matzke, & Wagenmakers, 2015). The Bayes factor can be thought of as the change from prior model odds to posterior model odds, as in the following equation:

$$\frac{\text{posterior odds}}{\frac{p(\mathcal{M}_m | \mathbf{y})}{p(\mathcal{M}_g | \mathbf{y})}} = \frac{\text{Bayes factor}}{\frac{p(\mathbf{y} | \mathcal{M}_m)}{p(\mathbf{y} | \mathcal{M}_g)}} \times \frac{\text{prior odds}}{\frac{p(\mathcal{M}_m)}{p(\mathcal{M}_g)}}. \quad [1]$$

In a confirmatory research setting, claims are sought about the relative probability of models, based on the data. These are claims about posterior odds, and thus require both prior odds and the Bayes factor measure of evidence. Thus, following the logic of

Equation 1, it is critical that the prior odds be declared (i.e., preregistered) before the critical data are seen. Alternatively, one might conduct a sensitivity analysis examining the range of prior odds for which the data lead to high posterior odds, thereby giving a lower bound on the amount of prior skepticism that would be required to negate the evidence in the data.

In an exploratory research setting, however, models are often inspired by the collected data. In this case, it is difficult to make claims about the prior probabilities of models. It does remain reasonable, however, to measure the *evidence* the current data provide for the newly-developed model, relative to other established models. This is what is measured by the Bayes factor, and it is validly calculated for any model with respect to any data.

From this perspective, what sets apart the exploratory settings is the need for extra care in expressing the knowledge claims. It is logical to say “these data are this many more times likely to arise under this model than that model.” It is not logical in the exploratory setting to say “this model is this many times more likely than that model, based on the data.” The former statement is appropriately cautious: it expresses only the strength of evidence and does not involve prior probabilities on the models. The latter statement is inappropriately bold: it expresses strength of belief and would have required prespecified priors on the models. After all, if a model was only inspired by examination of the data, it seems likely that its (implicit) prior probability was not high, and so the (implicit) posterior probability of the model also is not high. In the absence of prior probabilities, exploratory model development should restrict knowledge claims to those consistent with the Bayes factor interpretation of evidence.

**Solution-oriented modeling.** Ultimately, the test of the usefulness of a theory or model is whether it works in practical applications, and people have confidence in models that can be demonstrated to work. Applications of established models will often combine exploratory and confirmatory approaches. Verifying that the data in the domain are well captured by the model provides a test of model robustness and generalizability. Forcing a model to tackle real-world problems encourages solution-oriented science that may inspire future model development and evaluation (Watts, 2017).

In this vein, an important class of applied models come in the form of measurement models. The goal of these models is not necessarily to provide detailed accounts of cognitive phenomena, but to provide a useful “close enough” model that can infer context-relevant features of a person, stimulus, task, or some combination of all of these (Marsh, Morin, Parker, & Kaur, 2014). Measurement models have been historically important as the underpinning of the field of psychometrics and psychological assessment, and are of growing importance with real-world applications in the emerging field of cognitive data science. The “Cognitive Psychometric Models” case-study box provides an example of this sort of applied measurement modeling.

## Conclusion

Psychology’s crisis of confidence provides a challenge to the broader field, but it also provides an opportunity to improve psychological modeling in particular. In this article, we have attempted to identify a number of these opportunities, and highlighted emerging modeling practices and useful new ideas for psychological modeling. In particular, we have tried to highlight four key ideas that we offer as take-home recommendations.

## Example 8: Cognitive Psychometric Models

A recent advance in the development of measurement models is the practice of *cognitive psychometrics*, in which generic models of cognitive processing are applied in a measurement context (Batchelder, 2010). In one of the earliest such projects, Gerrein and Chechile (1977) used a model of a working-memory task to study the dynamics of alcohol-induced amnesia. The model was a multinomial processing tree, which is essentially a decision tree that the participant in a task is assumed to traverse in order to indicate a response. The model consisted only of a few chained probability statements and was not intended to test one or another hypothesis directly. Instead, it served mainly to restate the observed data in terms of process parameters rather than counts. The conclusions of interest—that alcohol intoxication impairs not only storage but also retrieval in working memory—were drawn based on patterns of change of these parameters across participant groups.

First, *preregistering models*, the predictions they make, and how they will be evaluated, is likely to improve the confidence the field has in results and conclusions of confirmatory model tests. Secondly, making models available and *post-registering exploratory model development* increases transparency and could speed model development. Thirdly, undertaking *detailed evaluation* of models improves the understanding of their strengths and weaknesses. Finally, we believe that *Registered Modeling Reports* could incentivize the field to test models that make risky predictions, providing strong tests of theory and potentially rapid progress (Platt, 1964).

## References

- Allais, M. (1953). Le comportement de l'homme rationnel devant le risque: Critique des postulats et axiomes de l'école Américaine. *Econometrica*, 21, 503–546.
- Allais, M. (1979). The foundations of a positive theory of choice involving risk and a criticism of the postulates and axioms of the American School. In M. Allais & O. Hagen (Eds.), *Expected utility hypothesis and the Allais paradox* (pp. 27–145). Dordrecht, The Netherlands: Riedel.
- Alogna, V., Attaya, M. K., Aucoin, P., Bahník, Š., Birch, S., Birt, A. R., ... others (2014). Registered replication report: Schooler and Engstler-Schooler (1990). *Perspectives on Psychological Science*, 9, 556–578.
- Batchelder, W. H. (2010). Cognitive psychometrics: Using multinomial processing tree models as measurement tools. In S. Embretson (Ed.), *Measuring psychological constructs: Advances in model-based approaches*. Washington, DC: American Psychological Association Books.
- Baumeister, A. A., & Kellas, G. (1968). Reaction time and mental retardation. *International Review of Research in Mental Retardation*, 3, 163–193.
- Baumgaertner, B., Devezzer, B., Buzbas, E. O., & Nardin, L. G. (2018, November). A Model-Centric Analysis of Openness, Replication, and Reproducibility. *ArXiv e-prints*, arXiv:1811.04525.
- Beck, K. (2003). *Test-driven development: by example*. Addison-Wesley Professional.
- Bell, R. M., Koren, Y., & Volinsky, C. (2010). All together now: A perspective on the Netflix prize. *Chance*, 23, 24–29.



- Birnbaum, M. H., & Quispe-Torrealblanca, E. G. (2018). TEMAP2.R: True and error model analysis program in R. *Judgment and Decision Making*, 13, 428–440.
- Bones, A. K. (2012). We knew the future all along: Scientific hypothesizing is much more accurate than other forms of precognition—a satire in one part. *Perspectives on Psychological Science*, 7, 307–309.
- Brown, G. D. A., Neath, I., & Chater, N. (2007). A temporal ratio model of memory. *Psychological Review*, 114(1), 539–576.
- Brown, S. D., Marley, A. A. J., Donkin, C., & Heathcote, A. (2008). An integrated model of choices and response times in absolute identification. *Psychological Review*, 115, 396–425.
- Bussemeyer, J. R., & Wang, Y.-M. (2000). Model comparisons and model selections based on generalization criterion methodology. *Journal of Mathematical Psychology*, 44, 171–189.
- Cavagnaro, D. R., Pitt, M. A., Gonzalez, R., & Myung, I. J. (2013). Discriminating among probability weighting functions with adaptive design optimization. *Journal of Risk and Uncertainty*, 47, 255–289.
- Chambers, C. D. (2013). Registered reports: A new publishing initiative at cortex. *Cortex*, 49, 609–610.
- Chambers, C. D., Dienes, Z., McIntosh, R. D., Rotshtein, P., & Willmes, K. (2015). Registered reports: Realigning incentives in scientific publishing. *Cortex*, 66, A1–A2.
- Cook, S. R., Gelman, A., & Rubin, D. B. (2006). Validation of software for bayesian models using posterior quantiles. *Journal of Computational and Graphical Statistics*, 15(3), 675–692.
- Criss, A. H., Malmberg, K. J., & Shiffrin, R. M. (2011). Output interference in recognition memory. *Journal of Memory and Language*, 64(4), 316–326.
- Dutilh, G., Annis, J., Brown, S. D., Cassey, P., Evans, N. J., Grasman, R. P. P. P., ... Donkin, C. (2018). The quality of response time data inference: A blinded, collaborative assessment of the validity of cognitive models. *Psychonomic Bulletin & Review*. doi:
- Dutilh, G., Vandekerckhove, J., Ly, A., Matzke, D., Pedroni, A., Frey, R., ... Wagenmakers, E.-J. (2017). A test of the diffusion model explanation for the worst performance rule using preregistration and blinding. *Attention, Perception, & Psychophysics*, 79, 713–725.
- Epper, T., & Fehr-Duda, H. (2018). *The missing link: Unifying risk taking and time discounting* (Tech. Rep. Nos. Department of Economics Discussion Paper 2018–12). St. Gallen, Switzerland: University of St. Gallen.
- Evans, N. J., Holmes, W. R., & Trueblood, J. S. (2018). *Response time data provides critical constraints on dynamic models of multi-alternative, multi-attribute choice*. Retrieved from <https://osf.io/h7e6v/> (Manuscript submitted for publication.)
- Farrell, S., & Lewandowsky, S. (2018). *Computational modeling of cognition and behavior*. Cambridge University Press.
- Finkbeiner, S. D., Briscoe, A. D., & Reed, R. D. (2014). Warning signals are seductive: Relative contributions of color and pattern to predator avoidance and mate attraction in heliconius butterflies. *Evolution*, 68, 3410–3420.
- Garner, W. R. (1953). An informational analysis of absolute judgments of loudness. *Journal of Experimental Psychology*, 46, 373–380.
- Gerrein, J. R., & Chechile, R. A. (1977). Storage and retrieval processes of alcohol-induced amnesia. *Journal of Abnormal Psychology*, 86(3), 285.
- Gould, S. J. (1996). *The mismeasure of man*. WW Norton & Company.
- Guan, M., Lee, M. D., & Vandekerckhove, J. (2015). A hierarchical cognitive threshold model of human decision making on different length optimal stopping problems. In R. Dale et al. (Eds.), *Proceedings of the 37<sup>th</sup> Annual Conference of the Cognitive Science Society*. Austin, TX: Cognitive Science Society.
- Hardwicke, T. E., & Ioannidis, J. (2018). *Mapping the universe of registered reports*. BITSS. Retrieved from [osf.io/preprints/bitss/fzpcy](https://osf.io/preprints/bitss/fzpcy) doi:
- Heathcote, A., Brown, S. D., & Wagenmakers, E.-J. (2015). An introduction to good practices in cognitive modeling. In B. U. Forstmann & E.-J. Wagenmakers (Eds.), *An introduction to model-based cognitive neuroscience* (pp. 25–48). Springer.
- Holland, M. K., & Lockhead, G. R. (1968). Sequential effects in absolute judgments of loudness. *Perception & Psychophysics*, 3, 409–414.
- Huber, J., Payne, J. W., & Puto, C. (1982). Adding asymmetrically dominated alternatives: Violations of regularity and the similarity hypothesis. *Journal of consumer research*, 9(1), 90–98.
- Jensen, A. R. (2006). *Clocking the mind: Mental chronometry and individual differences*. Elsevier.
- Kass, R. E., & Raftery, A. E. (1995). Bayes factors. *Journal of the American Statistical Association*, 90, 377–395.
- Kerr, N. L. (1998). Harking: Hypothesizing after the results are known. *Personality and Social Psychology Review*, 2, 196–217.
- Kılıç, A., Criss, A. H., Malmberg, K. J., & Shiffrin, R. M. (2017). Models that allow us to perceive the world more accurately also allow us to remember past events more accurately via differentiation. *Cognitive psychology*, 92, 65–86.
- Klein, R. A., Ratliff, K. A., Vianello, M., Adams, R. B. J., Bahník, S., Bernstein, M. J., ... others (2014). Investigating variation in replicability: A “many labs” replication project. *Social Psychology*, 45, 142–152.
- Latty, T., & Beekman, M. (2010). Irrational decision-making in an amoeboid organism: transitivity and context-dependent preferences. *Proceedings of the Royal Society B: Biological Sciences*, 278(1703), 307–312.
- Lee, M. D. (2018). Bayesian methods in cognitive modeling. In J. Wixted & E.-J. Wagenmakers (Eds.), *The Stevens' Handbook of Experimental Psychology and Cognitive Neuroscience. Volume 5: Methodology* (Fourth ed.). John Wiley & Sons.
- Lee, M. D., & Wagenmakers, E.-J. (2013). *Bayesian cognitive modeling: A practical course*. Cambridge University Press.
- Liew, S. X., Howe, P. D., & Little, D. R. (2016). The appropriacy of averaging in the study of context effects. *Psychonomic bulletin & review*, 23(5), 1639–1646.
- Luce, R. D. (1959). *Individual choice behavior: A theoretical analysis*. New York: Wiley.
- MacCoun, R. J., & Perlmutter, S. (2017). Blind analysis as a correction for confirmatory bias in physics and in psychology. In S. Lilienfeld & I. Waldman (Eds.), *Psychological science under scrutiny: Recent challenges and proposed solutions* (pp. 297–322). John Wiley & Sons.
- Marsh, H. W., Morin, A. J., Parker, P. D., & Kaur, G. (2014). Exploratory structural equation modeling: An integration

- of the best features of exploratory and confirmatory factor analysis. *Annual review of clinical psychology*, 10, 85–110.
- Matzke, D., Nieuwenhuis, S., van Rijn, H., Slagter, H. A., van der Molen, M. W., & Wagenmakers, E.-J. (2015). The effect of horizontal eye movements on free recall: A preregistered adversarial collaboration. *Journal of Experimental Psychology: General*, 144, e1–e15.
- Miguel, E., Camerer, C., Casey, K., Cohen, J., Esterling, K. M., Gerber, A., ... others (2014). Promoting transparency in social science research. *Science*, 343, 30–31.
- Munafò, M. R., Nosek, B. A., Bishop, D. V., Button, K. S., Chambers, C. D., du Sert, N. P., ... Ioannidis, J. P. (2017). A manifesto for reproducible science. *Nature Human Behaviour*, 1, 0021.
- Murdoch, B. B. (1960). The distinctiveness of stimuli. *Psychological Review*, 67, 16–31.
- Murdoch, B. B. (1962). The serial position effect in free recall. *Journal of Experimental Psychology*, 64, 482–488.
- Myung, I. J., Forster, M. R., & Browne, M. W. (2000). A special issue on model selection. *Journal of Mathematical Psychology*, 44, 1–2.
- Navarro, D. J. (in press). Between the devil and the deep blue sea: Tensions between scientific judgement and statistical model selection. *Computational Brain & Behavior*.
- Noble, W. S. (2009). A quick guide to organizing computational biology projects. *PLOS Computational Biology*, 5(7), 1–5. Retrieved from <https://doi.org/10.1371/journal.pcbi.1000424> doi:
- Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D., Breckler, S. J., ... others (2015). Promoting an open research culture. *Science*, 348, 1422–1425.
- Nosek, B. A., Ebersole, C. R., DeHaven, A. C., & Mellor, D. T. (2018). The preregistration revolution. *Proceedings of the National Academy of Sciences*, 201708274.
- Oberauer, K., Lewandowsky, S., Avh, E., Brown, G. D., Conway, A., Covan, N., ... others (2018). Benchmarks for models of short term and working memory. *Psychological Bulletin*.
- Open Science Collaboration. (2012). An open, large-scale, collaborative effort to estimate the reproducibility of psychological science. *Perspectives on Psychological Science*, 7, 657–660.
- Open Science Collaboration. (2015). Estimating the reproducibility of psychological science. *Science*, 349(6251), aac4716.
- Pashler, H., & Wagenmakers, E.-J. (2012). Editors' introduction to the special section on replicability in psychological science: A crisis of confidence? *Perspectives on Psychological Science*, 7, 528–530.
- Pitt, M. A., Kim, W., Navarro, D. J., & Myung, J. I. (2006). Global model analysis by parameter space partitioning. *Psychological Review*, 113(1), 57.
- Platt, J. R. (1964). Strong inference. *Science*, 146(3642), 347–353.
- Popper, K. R. (1959). *The logic of scientific discovery*. London: Routledge.
- Ratcliff, R., Schmiedek, F., & McKoon, G. (2008). A diffusion model explanation of the worst performance rule for reaction time and iq. *Intelligence*, 36, 10–17.
- Roe, R. M., Busemeyer, J. R., & Townsend, J. T. (2001). Multialternative decision field theory: A dynamic connectionist model of decision making. *Psychological review*, 108(2), 370.
- Rosenthal, R. (1979). The file drawer problem and tolerance for null results. *Psychological Bulletin*, 86, 638–641.
- Rouder, J., Haaf, J. M., & Snyder, H. K. (in press). Minimizing mistakes in psychological science. *Advances in Methods and Practices in Psychological Science*. Retrieved from <https://psyarxiv.com/gxyc5/>
- Rubin, D. C., Hinton, S., & Wenzel, A. (1999). The precise time course of retention. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 25, 1161–1176.
- Shanks, D. R., Newell, B. R., Lee, E. H., Balakrishnan, D., Ekelund, L., Cenac, Z., ... Moore, C. (2013). Priming intelligent behavior: An elusive phenomenon. *PLoS ONE*, 8, e56515.
- Shiffrin, R. M., Lee, M. D., Kim, W.-J., & Wagenmakers, E.-J. (2008). A survey of model evaluation approaches with a tutorial on hierarchical Bayesian methods. *Cognitive Science*, 32, 1248–1284.
- Silberzahn, R., Uhlmann, E. L., Martin, D. P., Anselmi, P., Aust, F., Awtry, E., ... others (2018). Many analysts, one data set: Making transparent how variations in analytic choices affect results. *Advances in Methods and Practices in Psychological Science*, 1, 337–356.
- Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology: Undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychological Science*, 22, 1359–1366.
- Simonsohn, U. (2013). Just post it: The lesson from two cases of fabricated data detected by statistics alone. *Psychological Science*, 24, 1875–1888.
- Simonson, I. (1989). Choice based on reasons: The case of attraction and compromise effects. *Journal of consumer research*, 16(2), 158–174.
- Spellman, B. A. (2015). A short (personal) future history of revolution 2.0. *Perspectives on Psychological Science*, 10, 886–899.
- Stewart, N., Brown, G. D. A., & Chater, N. (2005). Absolute identification by relative judgment. *Psychological Review*, 112, 881–911.
- Sun, R. (2008). *The Cambridge handbook of computational psychology*. New York, NY, USA: Cambridge University Press.
- Trueblood, J. S., Brown, S. D., & Heathcote, A. (2015). The fragile nature of contextual preference reversals: Reply to tsetsos, chater, and usher (2015). *Psychological Review*, 122(4), 848–853. doi:
- Turner, B. M., Schley, D. R., Muller, C., & Tsetsos, K. (2017). Competing theories of multialternative, multiattribute preferential choice. *Psychological Review*, 125(3), 329–362. doi:
- Tversky, A. (1972). Elimination by aspects: A theory of choice. *Psychological review*, 79(4), 281.
- Unsworth, N., Redick, T. S., Lakey, C. E., & Young, D. L. (2010). Lapses in sustained attention and their relation to executive control and fluid abilities: An individual differences investigation. *Intelligence*, 38, 111–122.
- Usher, M., & McClelland, J. L. (2004). Loss aversion and inhibition in dynamical models of multialternative choice. *Psychological review*, 111(3), 757.
- Vandekerckhove, J., Matzke, D., & Wagenmakers, E.-J. (2015). Model comparison and the principle of parsimony. In J. R. Busemeyer, Z. Wang, J. T. Townsend, & A. Eidels (Eds.), *Oxford handbook of computational and mathematical psychology* (pp. 300–317). New York, NY, US: Oxford

University Press.

- Vanpaemel, W. (2010). Prior sensitivity in theory testing: An apology for the Bayes factor. *Journal of Mathematical Psychology*, 54, 491–498.
- Vanpaemel, W., Vermorgen, M., Deriemaeker, L., & Storms, G. (2015). Are we wasting a good crisis? The availability of psychological research data after the storm. *Collabra*, 1:3. doi:
- Voss, A., Rothermund, K., & Voss, J. (2004). Interpreting the parameters of the diffusion model: An empirical validation. *Memory & Cognition*, 32, 1206–1220.
- Wagenmakers, E.-J. (2012). A year of horrors. *De Psychonom*, 27, 12-13.
- Wagenmakers, E.-J., & Farrell, S. (2004, Feb 01). Aic model selection using akaïke weights. *Psychonomic Bulletin & Review*, 11(1), 192–196. Retrieved from <https://doi.org/10.3758/BF03206482> doi:
- Wagenmakers, E.-J., Grünwald, P., & Steyvers, M. (2006). Ac-cumulative prediction error and the selection of time series models. *Journal of Mathematical Psychology*, 50, 149–166.
- Wagenmakers, E.-J., Wetzels, R., Borsboom, D., van der Maas, H. L., & Kievit, R. A. (2012). An agenda for purely confirmatory research. *Perspectives on Psychological Science*, 7, 632–638.
- Watts, D. J. (2017). Should social science be more solution-oriented? *Nature Human Behaviour*, 1, 1–5.
- Wicherts, J. M., Borsboom, D., Kats, J., & Molenaar, D. (2006). The poor availability of psychological research data for reanalysis. *American Psychologist*, 61(7), 726-728.
- Zhang, S., & Lee, M. D. (2010). Optimal experimental design for a class of bandit problems. *Journal of Mathematical Psychology*, 54, 499–508.
- Zwaan, R. A., Etz, A., Lucas, R. E., & Donnellan, M. B. (2018). Making replication mainstream. *Behavioral and Brain Sciences*, 41.