

1 Between the devil and the deep blue sea: Tensions between
2 scientific judgement and statistical model selection

3 Danielle J. Navarro¹

4 ¹ University of New South Wales

5 Abstract

6 Discussions of model selection in the psychological literature typically frame the issues as a question of statistical inference, with the goal being to determine which model makes the best predictions about data. Within this setting, advocates of leave-one-out cross-validation and Bayes factors disagree on precisely which prediction problem model selection questions should aim to answer. In this comment, I discuss some of these issues from a scientific perspective. What goal does model selection serve when all models are known to be systematically wrong? How might “toy problems” tell a misleading story? How does the scientific goal of explanation align with (or differ from) traditional statistical concerns? I do not offer answers to these questions, but hope to highlight the reasons why psychological researchers cannot avoid asking them.

Keywords: model selection, science, statistics

7 Model selection seems to be an evergreen topic in mathematical psychology. Given
8 two or more competing theories about the world, each instantiated as parameterised com-
9 putational models that provide different accounts of a data set, how should we decide which
10 model is better supported by the data? Typically we formulate this as a statistical inference
11 problem, with various authors arguing for Bayes factors (e.g., Wagenmakers 2007), mini-
12 mum description length (e.g., Grunwald 2007), cross-validation (e.g., Browne 2000) and a
13 variety of other possibilities besides. To highlight the behaviour of different model selec-
14 tion methods, we often consider “toy problems”, simplified versions of serious inferential
15 scenarios designed to elicit different intuitions about whether the model selection proce-
16 dure behaves sensibly. The large-sample results presented by Gronau and Wagenmakers

I am grateful to many people for helpful conversations and comments that shaped this paper, most notably Nancy Briggs, Berna Devezer, Chris Donkin, Olivia Guest, Daniel Simpson, Iris van Rooij and Fred Westbrook.

Correspondence concerning this article should be addressed to Danielle J. Navarro, School of Psychology, University of New South Wales. Kensington NSW 2052, Sydney, Australia. E-mail: d.navarro@unsw.edu.au

17 (2018) fall within this tradition, highlighted by the Dennis Lindley quote that motivates
18 the work. The results are perhaps unsurprising given the known inconsistency of orthodox
19 cross-validation estimators (Shao 1993), but there is value in highlighting the issue to a
20 broader audience and noting that a Bayesian formulation does not remove this limitation.
21 To the extent that some psychologists are unaware of the need for care when using cross-
22 validation methods – as indeed they may be unaware of a need for caution with respect to
23 Bayes factors or any other model selection procedure – the paper strikes me as helpful and
24 timely.

25 As much as I enjoyed the paper, I wonder whether the simplicity of exposition comes
26 at a cost. As Vehtari, Simpson, Yau and Gelman (2018) note in their commentary, Gronau
27 and Wagenmakers’ examples apply leave-one-out cross-validation in a fashion that is rather
28 at odds with how its advocates recommend that it be used. The original paper constitutes a
29 strong argument against naive or accidental misuse of some cross-validation procedures, but
30 the implications for best practice seem much less obvious. Noting that other commenters
31 have discussed technical issues in detail, my goal in this paper is to take a slightly broader
32 view on the tensions between scientific judgement and statistical model selection.

33 **Mistaking the map for the territory**

34 The quote by Lindley asks us to consider the question “if you can’t do simple problems,
35 how can you do complicated ones?” While I understand and sympathise with the sentiment,
36 for my own part I would be tempted to reverse the warning: if we *only* solve simple problems,
37 we may never learn how to think about the complex ones. As someone who has tried to use
38 many different model selection tools over the years, I am of the view that the behaviour of
39 a selection procedure applied to toy problems is a poor proxy for the inferential problems
40 facing scientists. As such, if we are to motivate our approach to model selection by quoting
41 famous statisticians, my preference would be to start with George Box’s (1976, p 792)
42 comment on the dangers of selective worrying:

43 Since all models are wrong the scientist must be alert to what is importantly
44 wrong. It is inappropriate to be concerned about mice when there are tigers
45 abroad.

46 Everyone who develops model selection tools is of course aware that all models are wrong.
47 Scientists do not fully understand the phenomena we are studying (else why study them?)
48 and every formal model-based description of the phenomenon is wrong in an unknown,
49 systematic fashion. One consequence of this is that while it is easy to construct artificial
50 scenarios in which any given procedure misbehaves, it becomes hard to know what impli-
51 cations these simple scenarios have for the real world scientific problems they approximate.

52 To illustrate how easy it is to tell a misleading story, consider the behaviour of the
53 Bayes factor – a procedure I presume Gronau and Wagenmakers would endorse as sensible
54 – when presented with a minor variation of their Example 1. In this scenario there are two
55 models, a “general law” \mathcal{M}_1 which asserts that a Bernoulli probability θ equals 1; and an

56 “unknown quantity” model \mathcal{M}_2 that expresses uncertainty by placing a uniform Beta(1,1)
 57 prior over θ . Given a sample n successes (i.e., all observations are 1) the Bayes factor will
 58 select \mathcal{M}_1 with certainty as $n \rightarrow \infty$, and the variant of leave-one-out cross-validation they
 59 discuss does not. The behaviour of the Bayes factor seems desirable insofar as \mathcal{M}_1 is the
 60 true model in this scenario. However, it is not difficult to reverse this intuition and construct
 61 an example where this same certainty seems *undesirable*.

62 Consider the “negligible error” scenario in which \mathcal{M}_1 is *almost* correct: the general law
 63 holds, apart from a single failure. The probability of success is 1, in the sense that one failure
 64 (or indeed any finite number of failures) in an infinite sequence of successes forms a set of
 65 measure zero. The true probability of success in a frequentist sense is $\lim_{n \rightarrow \infty} (n-1)/n =$
 66 1, and similarly, the posterior expected value of θ for the unknown quantity model \mathcal{M}_2
 67 converges on $\theta = 1$ in the large sample limit. In any sense that a pragmatic scientist would
 68 care about, the general law would count as the “correct” account for the phenomenon.¹
 69 Nevertheless the general law model \mathcal{M}_1 does not have support at the data \mathbf{x} . So while
 70 $P(\mathbf{x}|\mathcal{M}_1) = 0$ for all n after the single failure has occurred, \mathcal{M}_2 assigns positive prior
 71 probability to the data

$$P(\mathbf{x}|\mathcal{M}_2) = \int_0^1 \theta^{n-1}(1-\theta)d\theta = B(n, 2) = \frac{(n-1)!1!}{(n+1)!} = (n(n+1))^{-1}$$

72 The Bayes factor $P(\mathbf{x}|\mathcal{M}_1)/P(\mathbf{x}|\mathcal{M}_2)$ is therefore 0, and selects *against* the general law \mathcal{M}_1
 73 with certainty even though \mathcal{M}_1 makes an “almost exactly true” prior prediction, whereas
 74 \mathcal{M}_2 assigns the same degree of prior belief to the true rule $\theta = 1$ as it does to the exact
 75 opposite rule, $\theta = 0$.

76 To a statistician the reason for this misbehaviour is obvious, and rather boring: a
 77 general law formulated as a model that does not accommodate measurement error (and
 78 therefore lacks support across most of the sample space) will behave poorly in a world such
 79 as our own that actually does have such errors. The fact that the Bayes factor produces
 80 counterintuitive inferences when asked to choose between extremely bad models is not *prima*
 81 *facie* evidence that we should discard Bayes factors. Rather, it requires that we recognise
 82 that Bayes factors can produce strange answers when none of the models are “true”. In this
 83 instance the problem arises because the large sample behaviour of the Bayes factor is to
 84 select the model whose prior predictive distribution $P(\mathbf{x}|\mathcal{M})$ is closest in Kullback-Leibler
 85 divergence to the true data generating mechanism,² and this is often *not* the criterion that
 86 a scientist cares about. In real life none of us would choose \mathcal{M}_2 over \mathcal{M}_1 in this situation,
 87 because from our point of view the general law model is actually “closer” to the truth than
 88 the uninformed model. In general Kullback-Liebler divergence is not a good proxy for “what
 89 scientists care about”, and so in this instance the scientist would (quite correctly) disregard
 90 the Bayes factor and make the sensible choice. Importantly though, the fact that the Bayes

¹While there are many people who assert that “a single failure is enough to falsify a theory”, I confess I have not yet encountered anyone willing to truly follow this principle in real life.

²For instance, Gelman, Carlin, Stern & Rubin (2004, p586-587) present an analogous convergence result for the posterior distribution $P(\theta|x)$ within a single model \mathcal{M} . The result generalizes to the Bayes factor by noting that the Bayes factor identifies a model with the prior predictive distribution $P(x|\mathcal{M})$. Substituting $P(x|\mathcal{M})$ for the role of $P(x|\theta)$ in their derivation produces the necessary result.

91 factor does something unhelpful in a contrived example designed to make it misbehave tells
92 us very little – one way or the other – about whether it is useful in real life. The example
93 I chose is silly, and its evidentiary value is minimal.

94 Viewed more generally, I find it difficult to know how to apply simple examples to
95 real world problems. There are no shortage of illustrations that particular model selection
96 procedures misbehave when applied to problems they are not built to solve. For instance,
97 in one of my early papers (Navarro 2004) I documented an issue with (a specific version
98 of) the minimum description length criterion developed by Rissanen (1996) and introduced
99 to psychology by Pitt, Myung and Zhang (2002). The particular issue, in which it is
100 possible for a nested model to be judged *more complex* than the encompassing model, arose
101 when trying to solve an actual psychological model selection problem (see Navarro, Pitt &
102 Myung 2004) in which we compared an exponential forgetting function $y = a \exp(-bt)$ to
103 the strength-resistance model $y = a \exp(-bt^w)$ proposed by Wickelgren (1972) and several
104 other models besides. Given that the exponential function is a special case of the strength-
105 resistance model, it is logically impossible for it to be more complex, and the behavior of
106 the minimum description length criterion here is self-evidently absurd. Does that mean that
107 this criterion is “worse” than simpler criteria such as such as AIC (Akaike 1973) and BIC
108 (Schwarz 1978), in which model complexity is assessed simply by counting the number of
109 parameters? To me this seems the wrong lesson to draw, given that AIC and BIC both have
110 numerous flaws of their own. Fault can be found with *any* formal criterion for statistical
111 inference, as is nicely illustrated by the many documented concerns with p-values listed in
112 the psychological literature going back at least to Edwards, Lindman & Savage (1963). As
113 any survey of the statistical literature will reveal (e.g., Vehtari & Ojanen 2012), even the
114 basic desiderata for what model selection is supposed to accomplish are not agreed upon.
115 Viewed from this perspective, showing that a particular procedure behaves strangely in an
116 artificial scenario is not without value, but one should be wary of reading too much into
117 such demonstrations.

118 **Escaping mice to be beset by tigers**

119 To the extent that I am arguing that playing with toys leads us to encounter mice, I
120 suppose it is incumbent on me to say something about tigers. To my mind, there is at least
121 one tiger in plain view, namely the implied claim that *scientific* model selection questions
122 are addressable with statistical tools. If scientific reasoning necessarily takes place in a
123 world where all our models are systematically wrong in some sense (often referred to as the
124 \mathcal{M} -open case), what do we hope to achieve by “selecting” a model? To me, it seems that
125 much of this is tied to the question of what we consider the function of a model to be. In
126 considering this question Bernardo and Smith (2000, p238) write

127 Many authors . . . highlight a distinction between what one might call *scientific*
128 and *technological* approaches to models. The essence of the dichotomy is that
129 scientists are assumed to seek *explanatory* models, which aim at providing insight
130 into and understanding of the “true” mechanisms of the phenomenon under
131 study; whereas technologists are content with *empirical* models, which are not

132 concerned with the “truth”, but simply providing a reliably basis for practical
133 action in predicting and controlling phenomena of interest.

134 Under a “technological view”, the primary role of a model is *predictive*, though the pre-
135 diction problem differs depending on which methods one prefers. For example, under the
136 Bayes factor approach a model is identified with its prior predictive distribution $P(\mathbf{x}|\mathcal{M})$,
137 whereas under a cross-validation approach one is more likely to focus on the posterior pre-
138 dictive distribution $P(\mathbf{x}'|\mathbf{x}, \mathcal{M})$, where \mathbf{x}' represents future data drawn from the (unknown)
139 true distribution. Nevertheless, in both cases the primary role of a model is operationalised
140 in terms of predictions about data. In contrast to the predictive perspective, the “scientific
141 view” as described by Bernardo and Smith (2000) places more emphasis on the interpretabil-
142 ity and explanatory value of $P(\mathbf{x}|\theta, \mathcal{M})$. Ultimately Bernardo and Smith (2000) conclude
143 that the distinction is not especially important: if scientific models are evaluated on their
144 ability to make predictions, then the “scientific view” reduces to the “technological view”
145 for most intents and purposes.

146 My view is a little different. It strikes me as notable that statistics papers typically
147 define the term “generalisation” in a way that differs markedly from how psychologists define
148 the term when studying human inductive reasoning (e.g., Lake, Salakhutdinov & Tenen-
149 baum 2015). In the statistical context, predictive generalisation performance is typically
150 assessed with respect to test data sampled from the *same* process as the training data (e.g.,
151 Vehtari & Ojanen 2012). In the literature on human reasoning, however, generalisation is
152 typically assessed by examining how people think about test items that are *systematically*
153 *different* to the data upon which they were trained, and cannot be (easily) described as re-
154 alisations of the “same” data generating process from which the training data arose. In my
155 opinion at least, scientific model selections problem seem to have more in common with the
156 latter than with the former. To illustrate this, consider the question of why we consider the
157 Rescorla-Wagner model of Pavlovian conditioning (Rescorla & Wagner 1972) to be such an
158 important milestone in the development of theories of learning. While the model did indeed
159 provide a good account of a range of existing conditioning phenomena, such as blocking
160 (Kamin, 1969), overshadowing (Pavlov, 1927), conditioned inhibition (Rescorla, 1969), and
161 contingency effects (Rescorla, 1968) the truly impressive contribution was not the ability to
162 predict new data from replications of these experiments but rather to successfully anticipate
163 new phenomena, such as overexpectation (Lattal & Nakajima, 1998) and super conditioning
164 (Rescorla, 1971). That is, one of the most important functions of a scientific theory is not
165 simply to predict new data from old experiments, but to encourage directed exploration of
166 new territory, as illustrated by the important role the Rescorla-Wagner model has played in
167 assisting neuroscientists to investigate reward prediction error signals (e.g., Schultz, Dayan
168 & Montague, 1997). Curiously, it has sometimes been argued (Devezer, Nardin, Baum-
169 gartner and Buzbas, under review) that the apparent paradox of scientific progress in the
170 absence of replication (Shiffrin, Borner & Stigler 2018) may be tied to exactly this kind of
171 theory-guided scientific exploration.

172 It is not that statisticians are unaware of these issues, of course. For example, in
173 a thorough survey on the literature on Bayesian prediction methods, Vehtari and Ojanen
174 (2012, p174-177) characterise the issue very cleanly, by noting that if the training data

175 are all conditioned on specific values \mathbf{v} for auxiliary or explanatory variables but the test
176 data depend on new values \mathbf{v}' , then the prediction problem changes considerably. If the
177 values of \mathbf{v}' can differ *systematically* from the known values \mathbf{v} – as might happen if a
178 researcher with different theoretical views designs a different experiment to one’s own, or
179 the task used to isolate a psychological process changes – I am skeptical that any statistical
180 framing of the problem is any more than an “in principle” solution. None of us are in a
181 position to know what future experiments we or others may run, and estimating the future
182 performance of a model with regards to data collected via unknowable experiments is likely
183 impossible. To pretend otherwise strikes me as a form of what Box (1976, p797-798) referred
184 to as *mathematistry*: using formal tools to define a statistical problem that differs from the
185 scientific one, solving the redefined problem, and declaring the scientific concern addressed.

186 To illustrate how poorly even the best of statistical procedures can behave when
187 used to automatically quantify the strength of evidence for a model, I offer the following
188 example. As part of an exercise evaluating category learning models, Lee and Navarro (2002)
189 collected similarity ratings for nine items that varied on two ternary-valued features, shape
190 (circle, square or triangle) and colour (red, green or blue). The optimal multidimensional
191 scaling solution for representing these items was estimated by solving a model order selection
192 problem, using the most reasonable statistical criterion we could think of at the time (see Lee
193 2001a, 2001b). The estimated solution embeds these nine items within a four dimensional
194 space: two dimensions are used to represent the colours (i.e., red, green and blue form
195 the vertices of a triangle), and two more are used to represent shape. No more than that
196 is *required* to describe the similarity judgments that people made: as a consequence this
197 stimulus representation ends up being the simplest adequate account of the data and is
198 arguably the statistically “correct” representation to estimate from these data.

199 Nevertheless, when we used this stimulus representation as part of a categorisation
200 task that used those same stimuli – shifting the context from \mathbf{v} to \mathbf{v}' as it were – categori-
201 sation models that relied on this representation to define a measure of stimulus similarity
202 behaved very poorly. These failures did not occur due to a statistical failure in our multi-
203 dimensional scaling procedure, they arose because of a substantive scientific concern that
204 relates to the difference between the two tasks. The four dimensional embedding space
205 does not allow dimensional attention rules (e.g., Kruschke 1992) to be applied to *specific*
206 feature values, because the features themselves are not represented explicitly as *dimensions*.
207 That is, because “circle-versus-not-circle” is not represented as a primitive feature within
208 this four-dimensional multidimensional scaling solution, a categorisation model that relies
209 on this representation cannot use it as the basis for selective attention, even though human
210 participants do precisely this. To generalise sensibly from the similarity judgment task to
211 the categorisation task, the required representation involved placing the same items on a
212 six dimensional hypercube³ (i.e., employing six binary-valued features: circle vs not-circle,
213 square vs not-square, etc).

214 Critically, the reason this seems to happen is that there are factors \mathbf{v}' that influence
215 the notion of “stimulus similarity” (e.g., learned dimensional attention based on feedback,

³For the purposes of full disclosure, I should note that the precise situation from Lee and Navarro (2002) is a little more complex than this description implies and other details about how we had to adapt a model from one context to be applicable to the other have been omitted.

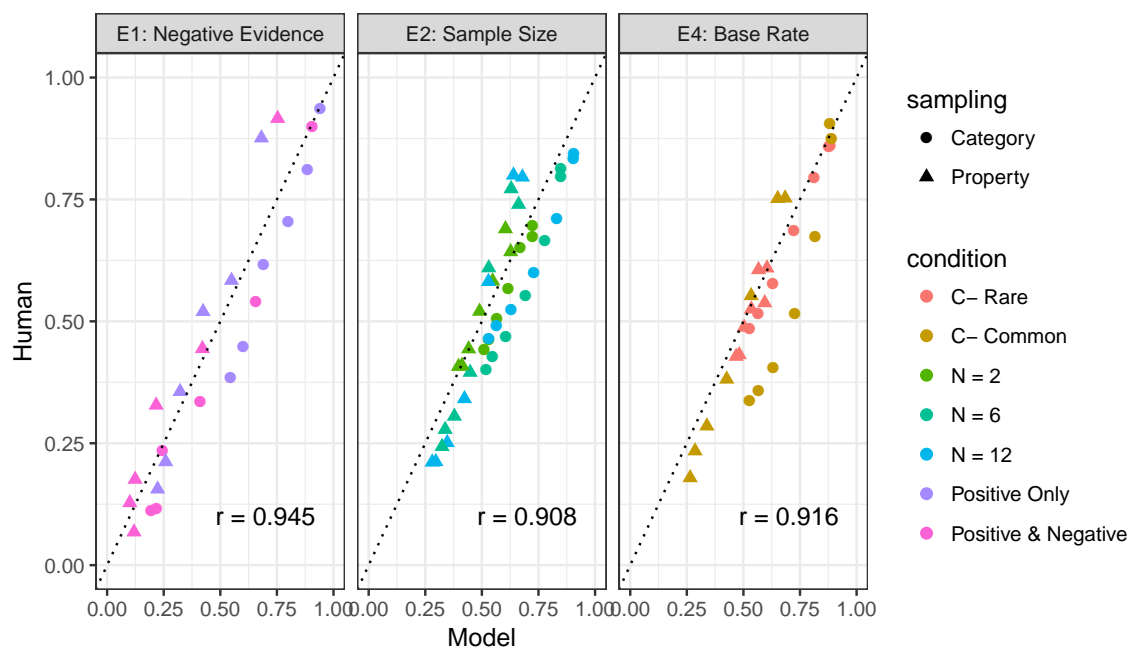


Figure 1. Model selection as viewed as a statistical problem typically emphasises quantitative measures of agreement between model predictions (or fitted values, x-axis) and human responses (y-axis). Even without any explanation given for the condition names or the experimental design, it is clear that the model in this figure provides a very good fit to the data. Nevertheless, knowing that the model fits depend on the values of parameters estimated from data, one might be tempted to ask if the researcher has encountered the Scylla of overfitting. Perhaps this apparent good performance is an illusion.

216 emphasis on differences between items) that applies in the categorisation task; and these
 217 are subtly different to the corresponding factors v (e.g., no feedback to direct attention,
 218 emphasis on commonalities among items) that apply to “stimulus similarity” in the direct
 219 elicitation task. In other words, because these auxiliary factors differ systematically be-
 220 tween the two tasks, even this “simple” generalisation turns out to be difficult and – while
 221 statistical measures of the adequacy of different similarity models were undoubtedly useful
 222 to us – it is unclear to me how we could have solved this model selection problem as a
 223 purely statistical exercise.

224 Between the devil and the deep blue sea

225 Gronau and Wagenmakers (2018) frame the question of model selection as a perilous
 226 dilemma in which one is caught between two beasts from classical mythology, the *Scylla* of
 227 overfitting and the *Charybdis* of underfitting. I find myself often on the horns of a quite
 228 different dilemma, namely the tension between the *devil* of statistical decision making and
 229 the *deep blue sea* of addressing scientific questions. If I have any strong opinion at all on
 230 this topic, it is that much of the model selection literature places too much emphasis on
 231 the statistical issues of model choice and too little on the scientific questions to which they

232 attach.

233 To again focus on my own papers rather than criticise others, consider the model fits
234 reported by Hayes, Banner, Forrester and Navarro (under review). In that paper we were
235 interested in how people’s inductive reasoning from data is shaped by what they know about
236 the process by which the data were selected, referred to as *sensitivity to sampling* in the
237 literature. This is a theme I have explored across multiple papers in the last several years.
238 To model sensitivity to sampling we relied on earlier work by Tenenbaum and Griffiths
239 (2001), as do most papers I have written on this topic (e.g., Navarro, Dry & Lee 2012,
240 Ransom, Perfors & Navarro 2016, Voorspoels, Navarro, Perfors, Ransom & Storms, 2015).
241 However, the task that we used in the Hayes et al. (under review) paper differs from
242 previous ones in many ancillary respects, and these ancillary details need to be formalised
243 in specific model choices. Some such choices (e.g., how smooth is an unknown generalisation
244 function?) can be instantiated as model parameters, but others (e.g., what class of functions
245 is admissible to describe human generalisation?) are not so simple. I think the choices I
246 made are sensible, but reasonable people might disagree.

247 How should I evaluate my modelling choices? A statistical perspective on this in-
248 ference problem might begin by estimating model parameters θ and producing a measure
249 of predictive performance. Setting aside the computational details of how one does this,
250 the result is likely to lead to a comparison between model predictions and human perfor-
251 mance similar to the one shown in Figure 1. Even without knowing the particular details
252 of the experiments, the scatterplot showing the fitted model values (x-axis) against the av-
253 erage reponse given by human participants (y-axis) across a large number of experimental
254 conditions strongly suggests that the model fits the empirical data well.

255 Perhaps it fits too well? When presented with such a figure, a reader familiar with
256 the model selection literature might be concerned that I have run afoul of the Scylla of
257 overfitting. This is not an unreasonable concern, but I find myself at a loss as to how cross-
258 validation, Bayes factors, or any other automated method can answer it. My scientific goal
259 when constructing this model was *not* to maximise the correlations as shown in Figure 1,
260 it was to *make sense* of the observed generalisation curves shown in Figure 2. The data in
261 Figure 2 are the same as those plotted in Figure 1, but drawn in a way that highlights the
262 empirical effects of theoretical interest. In each column there are multiple generalisation
263 curves shown, plotted separately for each experimental condition, with human data at
264 the top and model predictions at the bottom. It is clear from inspection that the data are
265 highly structured, and that there are systematic patterns to how people’s judgments change
266 across conditions. The scientific question of most interest to me is asking what theoretical
267 principles are required to produce these shifts. Providing a good fit to the data seems of
268 secondary importance. From visual inspection it is clear that the model captures *most*
269 patterns in the data, but not all. In particular, looking at the systematic model failure
270 in the second column from the right, the same reader might now be inclined to wonder
271 if I have fallen prey to the Charybdis of *underfitting*. So which of the mythical beasts,
272 Scylla or Charybdis, have I encountered? Would a cross-validation analysis or Bayes factor
273 calculation tell me? It seems unlikely.

274 To my mind, the bigger concern here is that to focus too heavily on the issue of

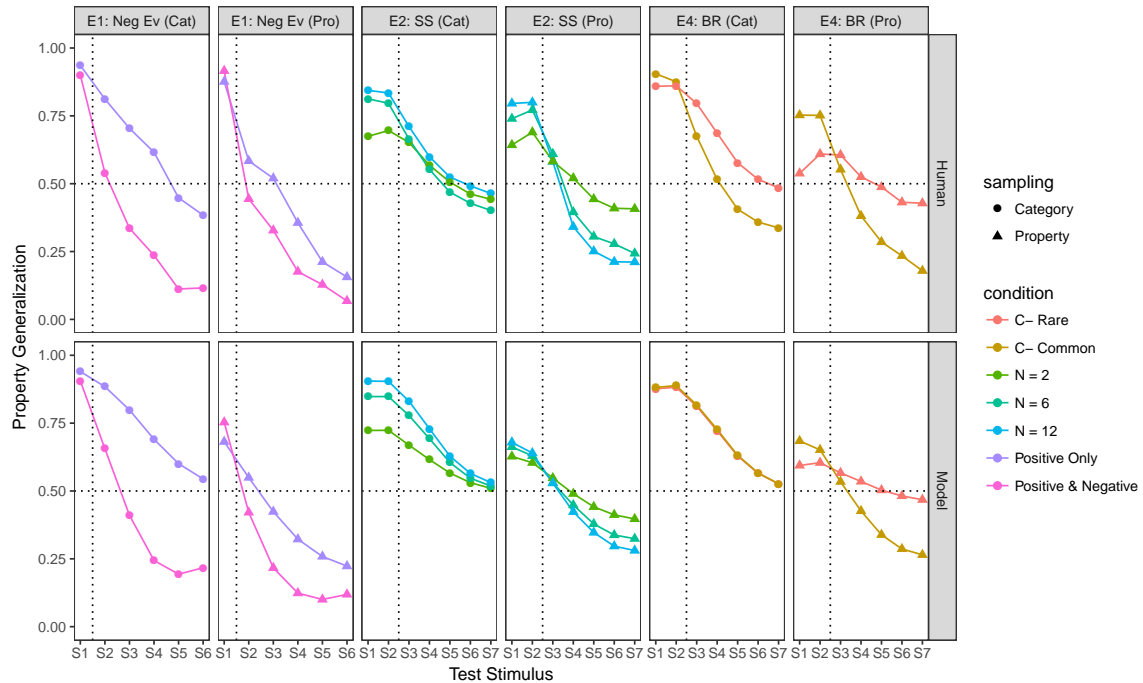


Figure 2. Scientific model selection is often more concerned with making sense of the systematic patterns observed in empirical data. This plots depict the extent to which people (top row) or a model (bottom row) will generalise (y-axis) from a small sample of training data to a novel item, shown as a function of the similarity of the novel item (x-axis) to the training data, with the most similar items shown on the left. Different panels (columns) and curves plotted separately as a function of three different experimental conditions reported by Hayes et al (under review). Even without a clear explanation of the different manipulations and their theoretical import, it is clear that the model provides a good account of the data in most conditions, but notably cannot reproduce the effect shown in the second panel from the right. One may be led to wonder if the researcher has encountered the Charybdis of underfitting. (Note: the data are the same as those plotted in Figure 1)

275 under/overfitting is to be seduced by the devil of statistical decision making. When we
 276 actually analysed the data, the allure of the deep blue sea of science led us to a different
 277 perspective. The approach we took was to ignore the quantitative fits almost entirely,
 278 and focus on the extent to which the key *qualitative* patterns in the data are an invariant
 279 prediction of the model across different choices of the parameter values θ . Loosely inspired
 280 by the “parameter space partitioning” idea introduced by Pitt, Kim, Navarro and Myung
 281 (2006), we defined a set of ordinal constraints in the data that any theoretical account would
 282 need to explain (e.g., increasing the number of observations caused a crossover effect under
 283 property sampling, column 4 from the left), and then showed that under most parameter
 284 values in the model, the predictions about these ordinal effects did not change. In other
 285 words – to recast this in the “scientific versus technological” language used by Bernardo and
 286 Smith (2000) – the scientifically important patterns are correctly predicted by $P(\mathbf{x}|\theta, \mathcal{M})$

287 regardless of the specific value of θ .

288 To my way of thinking, understanding how the qualitative patterns in the empirical
289 data emerge naturally from a computational model of a psychological process is much more
290 scientifically useful than presenting a quantified measure of its performance, but it is the
291 latter that we focus on in the “model selection” literature. Given how little psychologists
292 understand about the varied ways in which human cognition works, and given the artifi-
293 ciality of most experimental studies, I often wonder what purpose is served by quantifying
294 a model’s ability to make precise predictions about every detail in the data. Much as the
295 false confidence of the Bayes factor in the “negligible error” scenario I constructed at the
296 beginning is entirely an artifact of its sensitivity to a bad ancillary assumption made by one
297 of the models (that θ must be exactly 1 for a general law to hold), it seems to me that in
298 real life, many exercises in which model choice relies too heavily on quantitative measures
299 of performance are essentially selecting models based on their ancillary assumptions. It is
300 unclear to me if this solves a scientific problem of interest.

301 References

- 302 Akaike, H. (1973). Information theory and an extension of the maximum likelihood princi-
303 ple. In B. N. Petrov & F. Csaki (eds), *Second International Symposium on Infor-*
304 *mation Theory*, pp. 267-281. Budapest: Akademiai Kiado.
- 305 Bernardo, J. M. & Smith, A. F. M. (2000). *Bayesian Theory* (2nd edition). John Wiley &
306 Sons.
- 307 Box, G. E. P. (1976). Science and statistics. *Journal of the American Statistical Association*,
308 *71*, 791-799.
- 309 Browne, M. (2000). Cross-validation methods. *Journal of Mathematical Psychology*, *44*,
310 108-132.
- 311 Devezer, B., Nardin, L. G., Baumgaertner, B. & Buzbas, E. (under review). Discovery
312 of truth is not implied by reproducibility but facilitated by innovation and epis-
313 temic diversity in a model-centric framework. *Manuscript submitted for publication*.
314 arxiv.org/abs/1803.10118
- 315 Edwards, W., Lindman, H. & Savage, L. J. (1963). Bayesian statistical inference for psy-
316 chological research. *Psychological Review*, *70*, 193-242.
- 317 Gelman, A., Carlin, J. B., Stern, H. S., & Rubin, D. B. (2003). *Bayesian Data Analysis*
318 (2nd ed)
- 319 Grünwald, P. (2007). *The minimum description length principle*. Cambridge, MA: MIT
320 Press.
- 321 Gronau, Q. & Wagenmakers, E. J. (2018). Limitations of Bayesian leave-one-out cross-
322 validation for model selection.
- 323 Hayes, B.K., Banner, S., Forrester, S. & Navarro, D.J. (under review). Sampling frames and
324 inductive inference with censored evidence. *Manuscript submitted for publication*.
325 <https://doi.org/10.17605/OSF.IO/2M83V>

- 326 Kamin, L.J. (1969). Predictability, surprise, attention, and conditioning. In B. A. Campbell
327 and R. M. Church (Eds.) *Punishment and Aversive Behavior*. New York: Appleton-
328 Century-Crofts (pp 279-296).
- 329 Kruschke, J. K. (1992). ALCOVE: An exemplar-based connectionist model of category
330 learning. *Psychological Review*, *99*(1), 22-44.
- 331 Lake, B. M., Salakhutdinov, R. & Tenenbaum, J. B. (2015). Human-level concept learning
332 through probabilistic program induction. *Science*, *350*(6266), 1332-1338.
- 333 Lattal, K. M., & Nakajima, S. (1998). Overexpectation in appetitive Pavlovian and instru-
334 mental conditioning. *Animal Learning & Behavior*, *26*(3), 351-360.
- 335 Lee, M. D. (2001a). On the complexity of additive clustering models. *Journal of Mathe-*
336 *matical Psychology*, *45*, 131-148.
- 337 Lee, M. D. (2001b). Determining the dimensionality of multidimensional scaling models for
338 cognitive modeling. *Journal of Mathematical Psychology*, *45*, 149-166.
- 339 Lee, M. D. & Navarro, D. J. (2002). Extending the ALCOVE model of category learning
340 to featural stimulus domains *Psychonomic Bulletin & Review*, *9*, 43-58
- 341 Navarro, D. J. (2004). A note on the applied use of MDL approximations. *Neural Compu-*
342 *tation*, *16*, 1763-1768
- 343 Navarro, D. J., Dry, M. J. & Lee, M. D. (2012). Sampling assumptions in inductive gener-
344 alization. *Cognitive Science*, *36*, 187-223
- 345 Navarro, D. J., Pitt M. A. & Myung, I. J. (2004). Assessing the distinguishability of models
346 and the informativeness of data. *Cognitive Psychology*, *49*, 47-84
- 347 Pavlov, I. (1927). *Conditioned Reflexes*. London: Oxford University Press
- 348 Pitt, M. A., Myung, I. J. & Zhang, S. (2002). Toward a method of selecting among com-
349 putational models of cognition. *Psychological Review*, *109*, 472-491.
- 350 Pitt, M. A., Kim, W., Navarro, D. J. & Myung, J. I. (2006). Global model analysis by
351 parameter space partitioning. *Psychological Review*, *113*, 57-83.
- 352 Rescorla, R. A. (1968). Probability of shock in the presence and absence of CS in fear
353 conditioning. *Journal of Comparative and Physiological Psychology*, *66*, 1-5.
- 354 Rescorla, R.A. (1969) Conditioned inhibition of fear resulting from negative CS-US contin-
355 gencies. *Journal of Comparative and Physiological Psychology*, *67*, 504-509.
- 356 Rescorla, R. A. (1971) Variations in the effectiveness of reinforcement following prior in-
357 hibitory conditioning. *Learning and Motivation*, *2*, 113-123.
- 358 Rescorla, R. A. & Wagner, A. R. (1972) A theory of Pavlovian conditioning: Variations in
359 the effectiveness of reinforcement and nonreinforcement. In A. H. Black & W. F.
360 Prokasy (eds) *Classical conditioning II: Current Research and Theory* (pp 64-99).
361 New York: Appleton-Century-Crofts
- 362 Ransom, K., Perfors, A. & Navarro, D. J. (2016). Leaping to conclusions: Why premise
363 relevance affects argument strength. *Cognitive Science*, *40*, 1775-1796

- 364 Rissanen, J. (1996). Fisher information and stochastic complexity. *IEEE Transactions on*
365 *Information Theory* 42, 40-47.
- 366 Schultz, W., Dayan, P., & Montague, P. R. (1997). A neural substrate of prediction and
367 reward. *Science*, 275(5306), 1593-1599.
- 368 Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, 6,
369 461-464
- 370 Shao, J. (1993). Linear model selection by cross-validation. *Journal of the American Sta-*
371 *tistical Association* 88, 486-494.
- 372 Shiffrin, R. M., Borner, K. & Stigler, S. M. (2018). Scientific progress despite irreproducibil-
373 ity: A seeming paradox. *Proceedings of the National Academy of Sciences, USA*,
374 115, 2632-2639.
- 375 Tenenbaum, J. B. & Griffiths, T. L. (2001). Generalization, similarity, and Bayesian infer-
376 ence. *Behavioral and Brain Sciences*, 24, 629-640.
- 377 Vehtari, A., Simpson, D., Yao, Y. & Gelman, A. (2018). Limitations of “Limitations of
378 Bayesian leave-one-out cross-validation”
- 379 Vehtari, A. & Ojanen, J. (2012). A survey of Bayesian predictive methods for model
380 assessment, selection and comparison. *Statistics Surveys* 6, 142-228.
- 381 Voorspoels, W., Navarro, D. J., Perfors, A., Ransom, K. & Storms, G. (2015). How do
382 people learn from negative evidence? Non-monotonic generalizations and sampling
383 assumptions in inductive reasoning. *Cognitive Psychology*, 81, 1-25
- 384 Wickelgren, W. A. (1972). Trace resistance and decay of long-term memory. *Journal of*
385 *Mathematical Psychology*, 9, 418-455.