

# Scoring Model Predictions using Cross-Validation\*

Anna L. Smith      Tian Zheng      Andrew Gelman

Department of Statistics  
Columbia University

Draft of October 26, 2018

## Abstract

We formalize a framework for quantitatively assessing agreement between two datasets that are assumed to come from two distinct data generating mechanisms. We propose a methodology for prediction scoring which provides a measure of the distance between two unobserved data generating mechanisms (DGMs), along the dimension of a particular model. The cross-validated scores can be used to evaluate preregistered hypotheses and to perform model validation in the face of complex statistical models. Using human behavior data from the Next Generation Social Science (NGS2) program, we demonstrate that prediction scores can be used as model assessment tools and that they can reveal insights based on data collected from different populations and across different settings. Our proposed cross-validated prediction scores are capable of quantifying true differences between data generating mechanisms, allow for the validation and assessment of complex models, and serve as valuable tools for reproducible research.

Keywords: *complex models; cross-validation; model assessment; preregistration; reproducibility.*

## 1 Introduction

To begin, we provide a description of the motivation for our proposed methodology, stemming first from recent recommendations for more reproducible research and second from the question of how to appropriately validate complex models. We provide a working definition for data generating mechanisms, which we rely on throughout the paper, and

---

\*This material is based on research sponsored by the Defense Advanced Research Projects Agency (DARPA) agreement number D17AC00001. The content of the information does not necessarily reflect the position or the policy of the Government, and no official endorsement should be inferred.

28 differentiate between predictive and inferential DGMs. Finally, we provide a summary of  
29 existing methods for tackling these two issues.

## 30 1.1 Reproducible research

31 Open Science Collaboration (2015)’s attempt to replicate one hundred high-impact psy-  
32 chological studies marked the real beginning of the scientific community’s latest struggle,  
33 the “replication crisis”, although many had commented on similar phenomena previously.  
34 In essence, this article highlighted the unfortunate fact that many scientific studies fail to  
35 replicate in practice.

36 As a very simple example, suppose researchers at University A identify a positive effect  
37 for some new drug or treatment (e.g., they find evidence to reject a null hypothesis of no  
38 effect). Later, researchers at University B replicate University A’s study – closely following  
39 University A’s descriptions of subject recruitment (perhaps even increasing the sample  
40 size), experimental procedure, and analysis – but fail to identify a positive effect or even  
41 find the opposite, a negative effect. In cases like this one, we are faced with two separate  
42 studies that are concerned with the same scientific hypothesis, using the same experiments,  
43 but somehow reach different conclusions. This phenomenon has serious implications for the  
44 scientific community at large and for the general public; it impacts our ability to trust any  
45 single particular finding and results in the watering down of the credibility of science at  
46 large. Not surprisingly, this problem is not confined to one particular research domain  
47 (although much initial discussion focused on studies in psychology), but is a concern to  
48 scientific researchers in all fields.

49 Closely related to replicability is the issue of reproducibility. A study or experiment  
50 (including its specific analytic procedure) is said to be replicable if when the study is  
51 repeated, with fresh data or subjects, similar results are achieved. An analysis is said to  
52 be reproducible if when the same data is analyzed again, identical results are achieved.

53 Ideally, published scientific results should be the product of studies or experiments  
54 whose results can be independently replicated. This typically requires that the identified  
55 effect sizes are relatively large and can be accurately measured, samples sizes are relatively  
56 large, and the entire data collection and experimental procedures are well documented, as  
57 well as the results of the reproducible analyses (e.g., analytic procedures are well docu-  
58 mented and explained and all code and data are made publicly available, if possible). In  
59 practice, ensuring replicability and reproducibility is often not straightforward. In fact,  
60 since the identification of this replication crisis, it has become clear that no simple solution  
61 exists. Instead, advances in this area will necessitate concerted efforts for methodological  
62 improvement and more research on reproducibility across many fields. Recently, a variety  
63 of new initiatives in this vein have been proposed. For example, some qualitative recom-  
64 mendations can be found in Spies (2018) and Stodden et al. (2016) provide suggestions for  
65 computational methods.

66 One recent recommendation involves the registration of hypothesis tests and scientific  
67 analyses prior to data collection (Humphreys et al., 2013; Gelman, 2013), so that one can

68 avoid the “garden of forking paths,” a description offered by Gelman and Loken (2014)  
69 for the analytic pipeline in which decisions on data coding, exclusion, and analysis are  
70 made contingent on the data, thus inducing problems of multiple comparisons even if only  
71 one analysis is done on the particular data at hand. Ideally, prior to collecting any data  
72 whatsoever, researchers first prepare a preregistered plan of all data collection methods,  
73 modeling and analytic procedures, hypotheses, as well as plans for handling any unexpected  
74 deviations from this regime. This preregistration is made publicly available in some way,  
75 so that the researchers are held accountable to their preregistered plan. Just as simple  
76 random sampling is a prerequisite for classical interpretations of sampling probabilities,  
77 standard errors and point estimates, preregistration ensures that the classical interpretation  
78 of hypothesis tests and the resulting p-values is appropriate. It is worth pointing out that  
79 this sort of preregistration does not preclude further exploratory analyses; The point of  
80 preregistration is not to restrict analyses but rather to provide more structure to analyses  
81 that are already planned. For example, after data collection, a researcher may notice a  
82 pattern or posit a new explanation that motivates additional analyses. Such additional  
83 exploratory data analysis (beyond preregistered plans) are generally desirable as they can  
84 lead to new discoveries or hypotheses and even inspire additional confirmatory research. In  
85 some fields, researchers have the option of submitting such a preregistration to a scientific  
86 journal whom, if the submission is accepted, will agree to publish the research prior to  
87 any data collection or results. Such a manuscript is called a registered report and usually  
88 undergoes a round of peer review prior to the journal’s agreement to publish. Not only do  
89 registered reports encourage reproducible research, but they also help journals avoid the  
90 negative impacts of publication bias.

## 91 1.2 Assessing preregistered hypotheses

92 One requirement of these reports (both in the case of *registered reports* or *preregistrations*)  
93 is that the researchers make predictions about the scientific hypotheses to be assessed  
94 and models to be fit once data collection is complete. Once the analysis is complete, the  
95 researchers are faced with a natural question: *how well do the preregistered predictions align*  
96 *with the observed data?* This question requires a methodology to *score* the predictions, in  
97 the face of the materialized observations, usually through some form of “prediction scoring.”  
98 Further, prediction scoring represents one of the few *quantitative* recommendations for  
99 improving the reproducibility of scientific research.

100 The form of such predictions will largely impact the types of scientific insights one  
101 can gain, as well as impact the procedure for prediction scoring. We will discuss such  
102 impacts further in Section 3. For now, consider the case where the researcher is able to  
103 make predictions about the resulting data at the *observation-level* (i.e., rather than at  
104 some higher *summary* or *model estimate* level). In practice, such predictions could come in  
105 the form of data from a previously conducted closely related study, as pilot data, or from  
106 simulated data that represents *a priori* knowledge on the true data generating mechanism.  
107 We may then think of prediction scoring as a measure of the agreement or distance between  
108 the *assumed* data generating mechanism (DGM) behind the preregistered predictions and

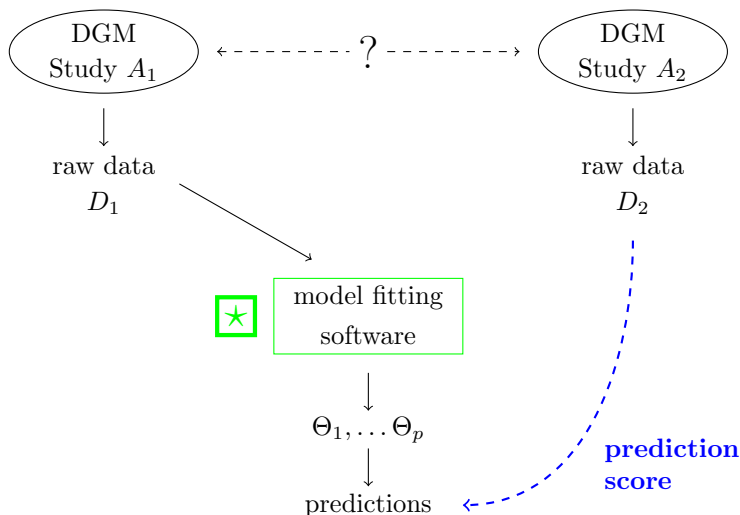


Figure 1: Prediction scores measure the agreement between predictions and realized data. In the case of preregistered predictions, Study  $A_1$  represents a set of pilot data, or data from an existing study, or simulated data.

109 the *true* data generating mechanism (DGM) behind the realized observations, as in Figure  
 110 1.

111 In this sense, the replication crisis and the related push for more reproducible research  
 112 motivate methodologies that can appropriately analyze and interpret scientific hypotheses  
 113 or models across different settings. When performing a preregistered study, how can we  
 114 quantitatively evaluate differences between the preregistered hypotheses or predictions and  
 115 the observed data? More generally, when we have access to data that comes from dif-  
 116 ferent settings of the same experimental framework (i.e., can be viewed as replications of  
 117 each other), can we quantify and evaluate differences across these settings? The first step  
 118 we propose is to view each of these (the set of predictions and the observed data) as re-  
 119 alizations from two distinct data generating mechanisms that describe the experimental  
 120 framework—the preregistered hypotheses or predictions follow from a DGM based on our  
 121 prior knowledge or pilot data and the observed experimental data follows from the true  
 122 observed DGM. In this sense, the evaluation of the preregistered hypotheses comes down  
 123 to identifying differences between the two DGMs.

### 124 1.3 Complex data generating mechanisms

125 In our approach, we will use the term “data generating mechanism” to refer to a particular  
 126 member of a family of probability distributions or equations that represent a set of (model)  
 127 assumptions. In this sense, we can use subsets of this family to represent beliefs about the  
 128 data generating mechanism that gave rise to the observed data, e.g., null and alternative  
 129 hypotheses. For example, suppose we believe two sets of experimental data are exponen-  
 130 tially distributed and we want to test hypotheses about the expected value of these two

131 distributions, represented by the parameter  $\lambda$ . The DGM family is a set of distributions for  
132  $\lambda$  and, for example, two members of this family that we might be interested in evaluating  
133 are  $p_1(\lambda \leq 5)$ , representing a null hypothesis, and  $p_2(\lambda > 5)$ , representing an alternative  
134 hypothesis. Alternatively, we can think of the DGM family as a unifying experimental  
135 framework, such that experiments within this framework can be viewed as replications of  
136 each other. Continuing our previous example, perhaps two groups of researchers each ran  
137 an experiment where they recorded the waiting time between participants' incoming calls  
138 on their personal phones and we want to compare the average waiting time across these  
139 experiments.

140 In practice, we can categorize data generating mechanisms as falling within one of  
141 two distinct forms:

- 142 1. *predictive*, where the data generating mechanism can be represented by a conditional  
143 probability distribution,  $p(y|x)$ , or
- 144 2. *inferential*, where the data generating mechanism can be represented by a distribution  
145 for a parameter or set of parameters,  $p(\theta)$ .

146 For example, in regression-style analyses, we are most interested in the conditional  
147 distribution of some response, given fixed covariate or predictor values. Both parametric  
148 and non-parametric regression-style models and other forms of predictive analysis which  
149 focus on sampling values (from a probability distribution) for a response or outcome vari-  
150 able,  $y$ , given fixed covariate or predictor values,  $x$ , fall in the first case. On the other hand,  
151 many analyses focus on a particular marginal distribution,  $p(\theta)$ . For example, we may use  
152 linear regression to model some phenomenon but are only interested in one particular slope  
153 parameter (i.e., the effect size of one particular predictor). In the example with phone call  
154 waiting times mentioned above, we have conceptualized the data generating mechanism as  
155 a set of distributions for a model parameter,  $\lambda$ . However, for many complex models or  
156 datasets with highly nuanced features, choosing a single parameter or summary statistic  
157 upon which to base evaluation of differences between DGMs is not straightforward. In  
158 these cases, relying on a single summary measure to capture all relevant (unknown) differ-  
159 ences between DGMs has the potential to oversimplify true differences or, in worst cases,  
160 fail to detect true differences altogether. For this reason, the prediction scoring methodol-  
161 ogy which we propose below is focused primarily on scoring differences between predictive  
162 DGMs.

## 163 1.4 Problem formulation

164 Under these settings, the natural *prediction scoring* question is translated to the following:  
165 are these experiments or realizations products of the same DGM or are they distinct in  
166 some way? While this is certainly a natural question, it is ill-posed for most experimental  
167 social science research settings. In almost all cases, the data generating mechanisms do  
168 in fact vary across experiments or settings, even if only slightly. Instead, we will focus  
169 on answering the following: How much do the data generating mechanisms differ across

170 settings (e.g., from preregistration to observed data) in a quantifiable way? Ultimately, we'd  
171 also like to be able to compare different ways that these differences across experimental  
172 frameworks (i.e., across experimental operationalizations of related scientific theories, ideas,  
173 or hypotheses) can be quantified.

174 In practice, just like any other statistical method based on sampled data, the observed  
175 differences between any two data generating mechanisms have two sources. The first is *by*  
176 *chance*, i.e., random variation in the data such as different samples of participants with  
177 different covariates and behaviors, or random variation in model estimates due to stochastic  
178 modeling algorithms. The second source is the *true differences* between the two DGMs.  
179 The latter is what we care to infer to derive scientific insight. Any methodology to compare  
180 these data generating mechanisms based on observed differences needs to be able to set  
181 apart differences due to these two sources. We will return to this point in Section 3.

182 Our proposed approach is to treat prediction scores as an *instrument* for quantify-  
183 ing the difference between our beliefs about the scientific process under study and reality,  
184 along a clearly specified “dimension”. With this language, we hope to evoke the type of  
185 instruments used by social scientists, where a variable is constructed to measure something  
186 abstract or unobserved in some sense and carries the same meaning under a general set-  
187 ting. For example, to measure a subject’s level of extraversion (which cannot be measured  
188 directly), a social scientist might design a survey that includes items that capture behavior  
189 indicative of extraversion. In a similar sense, we treat the true scientific process or data  
190 generating mechanism (DGM) as unobservable. This is a natural assumption, which sug-  
191 gests the construction of a numerical measure for the distance we are interested in (between  
192 our beliefs about that DGM and the true DGM) is not trivial and requires clear definitions  
193 and rigorous evaluation.

194 As shown in Figure 1, prediction scores shall directly measure the distance or dif-  
195 ference, via a predictive model, between preregistered predictions and observations from  
196 experimental data. The preregistered predictions are conditioned on our beliefs about the  
197 DGM through the identification of pilot or sample data, and, often less discussed, through  
198 the choice of a particular model. The model and the form of the prediction jointly define an  
199 aspect of the DGM for which the prior belief will be evaluated against the reality using our  
200 proposed prediction scores. In this sense, prediction scores can be thought of as capturing  
201 two levels of scientific insight: (1) how well the predictions match the materialized obser-  
202 vations and (2) how well our belief agrees with the reality in terms of the data generating  
203 mechanism.

## 204 **2 Background**

205 As best we can tell, current practice for prediction scoring in registered reports generally  
206 consists of making predictions in the form of directional hypotheses (in some cases, pre-  
207 dictions for the relative effect size are also included) for model parameters or summary  
208 statistics and assessing these predictions by performing a corresponding hypothesis test  
209 and checking for a significant effect (for some examples of published registered reports, see

210 the Zotero library maintained by the Center for Open Science: Mellor, 2018). Our proposed  
211 methodology will advance these methods in two main directions. First, it is general enough  
212 to accommodate parameter-level hypotheses (or other, higher-level summary statistics) as  
213 well as predictions at the individual-level or lowest level of analysis. Individual-level pre-  
214 dictions will allow for a more fine-grained assessment of the agreement between our prior  
215 beliefs and the true underlying DGM. We will also strongly encourage that these predic-  
216 tions incorporate appropriate measures of uncertainty, such as in the form of probabilistic  
217 forecasts. Second, our methodology will provide prediction scores on a continuous scale,  
218 which can be viewed as estimates of the distance between our prior beliefs and reality.  
219 Thus, we provide a quantitative measure of prediction performance rather than the simple  
220 binary detection of a significant effect.

221 In order to provide these advances, we pull ideas and insights from related research  
222 in the statistical literature; Below, we briefly describe statistical methodology for the eval-  
223 uation of probabilistic forecasts, Bayesian software-checking procedures, and approximate  
224 cross-validation (a more thorough discussion is available in the Supplementary Materials).

## 225 **2.1 Diagnostic plots for probabilistic forecasts and scoring rules**

226 In the statistical literature, perhaps the most applicable line of research to inform pre-  
227 diction scoring for preregistered hypotheses is the evaluation of probabilistic forecasts and  
228 the theoretical development of *scoring rules*. Gneiting and Katzfuss (2014) provide a nice  
229 summary of recent research in this area. First, let us point out that a scoring function  
230 measures the agreement between a point prediction and an observation while a scoring rule  
231 measures the agreement between a probabilistic forecast (a predictive probability distri-  
232 bution over future quantities or events of interest, such as a posterior predictive density  
233 from a Bayesian analysis) and an observation. Naturally, a probabilistic forecast contains  
234 much more information than a simple point prediction and, most importantly, provides  
235 a suitable measure of the uncertainty associated with the predictions. For this reason,  
236 we will focus on probabilistic forecasts (for a review of issues with point forecasts and  
237 scoring rules, see Gneiting, 2011). The importance of probabilistic forecasts as a tool for  
238 statistical inference is well-motivated by Dawid (1984)’s framework for prequential anal-  
239 ysis, which frames the creation of sequential probability forecasts (over time) as the true  
240 focus and underlying motivation for classical statistical concepts and theory. Of course,  
241 with the advent of rapidly increasing computational power, MCMC and other estimation  
242 techniques have greatly increased analysts’ ability to create probabilistic forecasts. In fact,  
243 in the past, much of the literature surrounding the evaluation of probabilistic forecasts  
244 came out of weather forecasting research. Currently, probabilistic forecasts have been used  
245 in applications ranging from climate models, flood risk, seismic hazards, renewable energy  
246 availability, economic and financial risk management, election outcomes, demographic and  
247 epidemiological projections, health care management, and preventive medicine.

248 Gneiting et al. (2004) and Gneiting et al. (2007) define two important characteristics  
249 of probabilistic forecasts: sharpness and calibration. In this context, sharpness is (solely) a  
250 property of the predictive distribution and refers to the concentration of the distribution.

251 For a real-valued variable, we could measure the sharpness of the probabilistic forecast by  
252 considering the average width of prediction intervals. On the other hand, calibration is a  
253 property of both the predictive distribution and the materialized observations or events.  
254 A probabilistic forecast is calibrated if the distributional forecast is statistically consistent  
255 with the observations; In other words, the observations should be indistinguishable from  
256 random draws from the predictive distribution. Gneiting et al. (2007) outline various lev-  
257 els of calibration—probabilistic, exceedance, and marginal (listed here from most to least  
258 strict)—as well as provide diagnostic tools for identifying these properties in practice. It  
259 should be noted that calibration is defined in terms of asymptotic consistency between  
260 random variable representations for the probabilistic forecast and the true underlying dis-  
261 tribution for the observations (i.e.,  $F$  is a CDF-valued random variable representing the  
262 probabilistic forecast and  $G$  is a CDF-valued random variable representing the true data  
263 generating mechanism). Thus, in practice, these random variables are themselves unobserv-  
264 able and diagnostic approaches using sample versions (using empirical CDFs) are necessary  
265 to assess the calibration of a particular forecast.

266 To check for probabilistic calibration, histograms (or empirical CDFs, if the sample  
267 size is small) of the PIT (probability integral transform) values can be verified for unifor-  
268 mity (this idea can be traced as far back as Rosenblatt, 1952; Pearson, 1933, and perhaps  
269 earlier). In meteorological research, Talagrand et al. (1997) proposed a verification rank  
270 histogram or Talagrand diagram (Anderson, 1996; Hamill and Colucci, 1997) to assess the  
271 calibration of ensemble forecasts and Shephard (1994, page 129) has used a similar dia-  
272 gram to assess samples from an MCMC algorithm. However, in the introduction of their  
273 paper, Gneiting et al. (2007) demonstrate that merely checking for the uniformity of PIT  
274 values is insufficient for distinguishing the ideal forecaster from three (poorer) competitor  
275 forecasts. Instead of relying solely on the PIT diagnostic, the authors highlight additional  
276 diagnostics (described below) and advocate maximizing the sharpness of the predictive dis-  
277 tribution, subject to calibration, as mentioned above. To check for marginal calibration,  
278 Gneiting et al. (2007) suggest plotting differences between the average predictive CDF and  
279 the empirical CDF for the observations versus  $x$ . If the probabilistic forecast is marginally  
280 calibrated, we would expect to see only minor fluctuations about zero. Exceedance cali-  
281 bration does not allow for an obvious sample analogue.

282 Additionally, scoring rules allow us to assess calibration and sharpness simultaneously.  
283 Taking a decision theoretic perspective, we can think of a scoring rule as a loss function. In  
284 this sense, we can interpret the scores as penalties that the forecaster wishes to minimize.  
285 In terms of the choice of a particular form for a scoring rule, one natural restriction is  
286 that the truth or true forecast should receive an optimal score. This is precisely what is  
287 meant by proper scoring rules (some examples include the logarithmic score, the quadratic  
288 score, the spherical score, the continuous ranked probability score, and the Brier score).  
289 In fact, Gneiting and Raftery (2007) point out that the log Bayes factor is equivalent  
290 to a logarithmic scoring rule in the no-parameter case (i.e. forecasts do not depend on  
291 parameters to be estimated from the data). This implies that the log Bayes factor can be  
292 used to compare competing forecasting rules, and not only to compare models. When the  
293 forecasting rules are specified only up to unknown parameters which will be estimated from



294 the data, the authors outline a variation of cross-validation that could be used to replace  
295 the logarithmic score with other proper scoring rules, to estimate a predictive Bayes factor  
296 of some kind. While there are some connections to Bayesian methods, the literature on  
297 scoring rules and the evaluation of probabilistic forecasts generally assumes a frequentist or  
298 classical perspective. While the discussion is typically focused on predictions for continuous  
299 variables, Czado et al. (2009) provide extensions of many of these ideas for count variables.

300 This literature provides a sound framework for comparing probabilistic forecasts or  
301 predictions (such as from preregistration materials) to observed data, where each compet-  
302 ing forecast could correspond to different modelling choices or assumptions. The diagnostic  
303 tools and recommendations for scoring rules outlined above allow this comparison to be  
304 nonparametric and thus, enable the comparison of non-nested, highly diverse models. How-  
305 ever, each of these diagnostic measures is necessarily model-based in that any diagnostic  
306 plot or set of scoring rules depends on the model assumptions used to create the probabilis-  
307 tic forecast. This complicates the interpretation of the scores or diagnostics themselves,  
308 as they measure not only differences between our prior beliefs and the realized data (i.e.  
309 between the preregistered predictions and observations) but also any differences between  
310 the modelling choices and the true underlying data generating mechanisms. We will pro-  
311 pose a prediction scoring framework that uses cross-validation to remove the dependence  
312 on model-based differences which enables us to quantitatively measure true differences  
313 between our prior beliefs and the realized data.

## 314 **2.2 Bayesian software-checking**

315 Although perhaps not obvious at first glance, recent proposals for algorithm-checking of  
316 Bayesian model fitting software (Cook et al., 2006; Talts et al., 2018) can also provide  
317 interesting insights in the prediction scoring setting. These proposals recommend simulat-  
318 ing fake data conditional on random draws from the prior distribution, running the model  
319 fitting software to obtain draws from the posterior distribution, and using a summary  
320 measure to diagnose the alignment between the draws from the posterior distribution and  
321 the random draws from the prior distribution. Based on the self-consistency property of  
322 the marginal posterior and the prior distribution, these draws should be indistinguishable  
323 from one another. To diagnose this alignment, Cook et al. (2006) suggest computing em-  
324 pirical quantiles, comparing the random draw from the prior distribution to the posterior  
325 distribution based on that particular draw. The authors suggest looking at histograms  
326 of these quantiles, demonstrating that if the software is working correctly, the quantiles  
327 should be approximately uniformly distributed. Talts et al. (2018) point out that the em-  
328 pirical quantiles are necessarily discrete and that artifacts of this discretization can lead to  
329 misleading diagnostic quantile histograms. Instead, the authors suggest computing rank  
330 statistics which will follow a discrete uniform distribution, if the software is correct. Ad-  
331 ditionally, Talts et al. (2018) provide a nice summary of the types of expected deviations  
332 from uniformity that one might observe in the diagnostic histograms with corresponding  
333 explanations of modelling choices or software errors that could lead to such deviations.

334 In terms of the prediction scoring setting, we can think of this software-checking

335 methodology as a special case where the chosen modelling strategy matches the underlying  
336 DGM exactly. We will borrow ideas from this methodology, such as the use of empirical  
337 quantiles and rank statistics and the self-consistency properties, to motivate our proposed  
338 prediction scoring framework.

### 339 **2.3 Bayesian model selection and approximate cross-validation**

340 As briefly mentioned previously, our proposed prediction scoring framework will utilize  
341 cross-validation to separate true DGM differences from purely model-based differences.  
342 Cross-validation, particularly for Bayesian analyses, has been a very active research area  
343 in recent years. First, we should point that many Bayesian model comparison summary  
344 statistics (such as AIC, DIC, WAIC) can be motivated by the estimation of out-of-sample  
345 predictive accuracy (see Vehtari et al., 2012, for a thorough review, from a formal deci-  
346 sion theoretic perspective), which of course is one of the goals of cross-validation as well.  
347 Gelman et al. (2014) provide a nice review of these model comparison summary mea-  
348 sures. As opposed to exact leave-one-out cross-validation (LOOCV), each of the Bayesian  
349 model summary statistics utilize the full predictive density and perform an adjustment  
350 (e.g., importance sampling, or division by an appropriate variance) to remove the effect  
351 of over-fitting, since no data was actually held out. The authors conclude the paper by  
352 citing cross-validation as their preferred method for model comparison, despite its high  
353 computational cost and requirement that data can be easily partitioned (i.e., partitioning  
354 is often not straight forward for dependent data). In this line of thought, Vehtari et al.  
355 (2017) develop an approximate version of leave-one-out cross-validation which implements  
356 Pareto-smoothing of the importance sampling weights to improve robustness to weak pri-  
357 ors or influential observations. Li et al. (2016) develop a version of cross-validation that  
358 can be applied to models with latent variables, which relies on an integrated predictive  
359 density. In application with competing probabilistic forecasts, Held et al. (2010) compare  
360 software fitting algorithms using approximate cross-validation and many of the diagnostic  
361 plots mentioned by Gneiting et al. (2007). Finally, Wang and Gelman (2014) and Millar  
362 (2018) address the problem of appropriate data partitioning and out-of-sample prediction  
363 error estimation for multilevel or hierarchical model selection using cross-validation and  
364 predictive accuracy. Wang and Gelman (2014) highlight the fact that model selection can  
365 be largely based on the size and structure of the hierarchical data.

366 This line of research, and its proposed improvements and extensions of cross-validation  
367 in various Bayesian settings, can certainly be incorporated in the prediction scoring method-  
368 ology that we propose. Our contribution will be to expand this literature, from the perspec-  
369 tive of the registered reports setting as well as from the unique perspective offered by the set  
370 of NGS2 experiments (described in greater below). We formalize the use of cross-validation  
371 to appropriately adjust agreement measures between preregistered predictions and realized  
372 observations. In other words, we will recommend a unique combination of cross-validation  
373 *and* external validation to provide meaningful prediction scores and to enable nonparamet-  
374 ric model assessment. Further, in the application to NGS2, we will demonstrate how these  
375 cross-validated prediction scores can be used to assess scientific hypotheses across distinct

376 experiments and data in a nonparametric way.

### 377 **3 Cross-validated prediction scoring**

378 In this section, we provide a general framework for our proposed prediction scoring method-  
379 ology. Our goal is to formalize the problem and provide concrete procedures that are general  
380 enough to be applicable to a variety of statistical models and analytic procedures.

#### 381 **3.1 General framework**

382 For any family of data generating mechanisms, we will be interested in estimating the  
383 distance between different members of the same family. The assumption of a meaningful  
384 distance between DGMs is an essential element of this methodology; in order to make  
385 quantitative comparisons between DGMs, or between experimental settings, or between  
386 preregistered and confirmatory hypotheses, we need to define a distance between DGMs.

387 **Definition 3.1.** For a particular family of data generating mechanisms, let the distance  
388 between any two members of the family be given by

$$\Delta_{DGM} = f(p_i, p_j)$$

389 where  $p_i$  and  $p_j$  are the  $i$ th and  $j$ th members of the particular DGM family and the choice  
390 of the function  $f$  is motivated by the form of the DGM family.

391 Specifying the form of this distance is not straightforward. For example, consider the  
392 case where we are interested in measuring the distance between two straight lines in a two-  
393 dimensional Euclidean space. Candidate measures might include calculating the difference  
394 in the slope or calculating the Euclidean distance within some window. Each of these  
395 measures is a sensible candidate but could result in wildly different conclusions. The issue  
396 of choosing an appropriate distance metric is not unique to the example of lines in Euclidean  
397 space; a variety of candidate measures exist for assessing the distance or disagreement  
398 between sets, or network objects, or points in space, or shapes, etc. Instead, we argue that  
399 the form of this distance in the prediction scoring framework should be motivated by the  
400 form of the data generating mechanism family. For example, for predictive data generating  
401 mechanisms, we might consider conditional KL-divergence (also called relative conditional  
402 entropy), whereas for the inferential case,  $L_p$  distance is a more natural metric. Recall  
403 from Section 1, we want to move away from the simple binary question of disagreement  
404 across DGMs (i.e., are the two DGMs different?) and instead promote the quantification  
405 of a distance between them (i.e., how far apart are the two DGMs?).

406 As mentioned previously, we will treat the prediction scores as an instrument for  
407 estimating this unobservable distance between data generating mechanisms. In essence, the  
408 prediction scores compare the difference between model-based predictions and real-world  
409 observations, and in many ways, can be viewed as a validation procedure. Traditionally,  
410 model validation is used to assess the predictive ability of the model. In this setting,

---

**Algorithm 1** Prediction Scoring for Predictive Inference

---

```
1: procedure CROSS-VALIDATION
2:   for  $k = 1, \dots, K$  do
3:      $x_{1,-k} \leftarrow$  dataset  $x_1$  with  $k$ th observation(s) removed
4:      $\hat{\theta}|_{x_{1,-k}, y_{1,-k}} \leftarrow$  estimate using model fitting software,  $p_*$ 
5:      $\hat{y}_{1k} \leftarrow$  prediction, given  $x_{1,-k}, \hat{\theta}|_{x_{1,-k}, y_{1,-k}}$ 
6:      $q_{1k} \leftarrow g(\hat{y}_{1k}, y_{1k})$ 
7: procedure VALIDATION
8:   for  $k = 1, \dots, K$  do
9:      $\hat{\theta}|_{x_1, y_1} \leftarrow$  estimate using model fitting software,  $p_*$ 
10:     $\hat{y}_{2k} \leftarrow$  prediction, given  $x_2, \hat{\theta}|_{x_1, y_1}$ 
11:     $q_{2k} \leftarrow g(\hat{y}_{2k}, y_{2k})$ 
12: procedure PREDICTION SCORING
13:    $\Delta_{pred} \leftarrow h(q_1, q_2)$ 
```

---

411 we are less interested in the fit of any particular model and more interested in learning  
412 about potential differences in the data generating mechanism(s) across experiments or  
413 settings. Most importantly, note that validation captures differences due to *both* random  
414 noise and true differences. Instead of relying solely on validation measures, we propose using  
415 cross-validation to properly calibrate the measurements from validation (see Algorithm  
416 1 and Figure 2 for a description of our proposed methodology). In this way, we can  
417 separate the differences due to random variation (as measured by cross-validation) from  
418 any true differences between the the data generating mechanisms. Further, note that any  
419 decisions or conclusions based on validation or cross-validation results alone include the  
420 assumption that the researcher’s chosen model is correct. In this sense, any observed  
421 (apparent) differences between the data generating mechanisms could be due solely to an  
422 inadequate model. Instead, comparing results across validation and cross-validation avoids  
423 this issue. Because both routines rely on the same model fitting software, comparisons  
424 across these routines should be less sensitive to poor modelling choices. In this sense, we  
425 are using cross-validation to calibrate the results of the validation procedure.

426 For DGMs belonging to the same family, let  $x_1, x_2$  represent datasets corresponding  
427 to DGMs one and two, respectively. Let  $z$  represent the quantity of inference and  $p_*$  be the  
428 model fitting software, described by model parameters,  $\theta$ . As described in Algorithm 1 and  
429 Figure 2, the prediction scores are calculated as a difference between the distribution of  
430 prediction (dis)agreement measures across cross-validation and validation. For both cross-  
431 validation and validation procedures, we can define prediction (dis)agreement statistics as  
432 follows:

$$q_{jk} = g(z_{jk}, \hat{z}_{jk})$$

433 where  $z_{jk}$  is the  $k$ th observation (or set of observations) for the  $j$ th dataset,  $\hat{z}_{jk}$  is a set  
434 of predictions for this observation(s), and  $g$  is the (dis)agreement measure. For cross-  
435 validation,  $\hat{z}_{jk}$  is estimated from a model that uses the  $j$ th dataset with  $k$ th observation  
436 (or set of observations removed). For validation,  $\hat{z}_{jk}$  is estimated from a model that uses

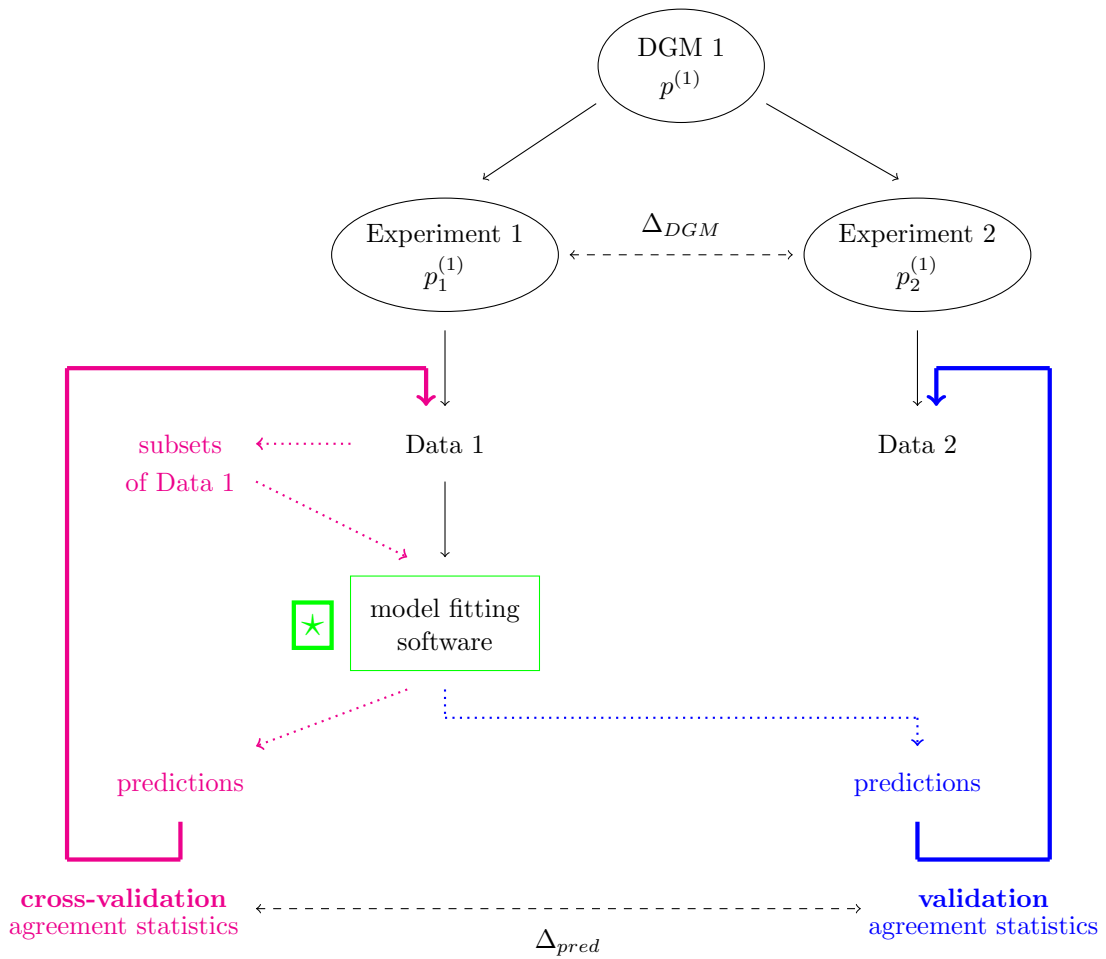


Figure 2: General outline of the proposed prediction scoring methodology for generic data generating mechanisms.

Model	$f$	$g$	$h$
linear regression <sup>1</sup>	conditional KL-divergence	empirical quantiles	KL-divergence
logistic regression <sup>2</sup>	$L^p$ distance		
logistic regression <sup>3</sup>	-	ROC curves	visual inspection
GP model <sup>3</sup>	-	MSE	difference

Table 1: Examples of choices of  $f, g$  and  $h$  used in this and related work. <sup>1</sup>Section 3.2; <sup>2</sup>Section 3.3; <sup>3</sup>Smith et al. (2018)

437 the  $(j - 1)$ th dataset, and plugs in any covariates or predictor variables observed in the  $j$ th  
438 dataset. The choice of  $g$  should be motivated by the model fitting software,  $p_{\star}^{(i)}$ , chosen by  
439 the researcher. For example, when using a linear regression model in focusing on inference  
440 for the conditional distribution  $p(y|x)$ , the predictions will be continuous and so quantiles  
441 are a natural choice. However, for logistic regression in the same predictive setting, the  
442 predictions will be probabilities (between 0 and 1) while the observations are binary. Some  
443 variant of the area under the curve (AUC) statistic would be a better choice for  $g$ .

444 Finally, with these sets of (dis)agreement measures, we can compute the prediction  
445 score:

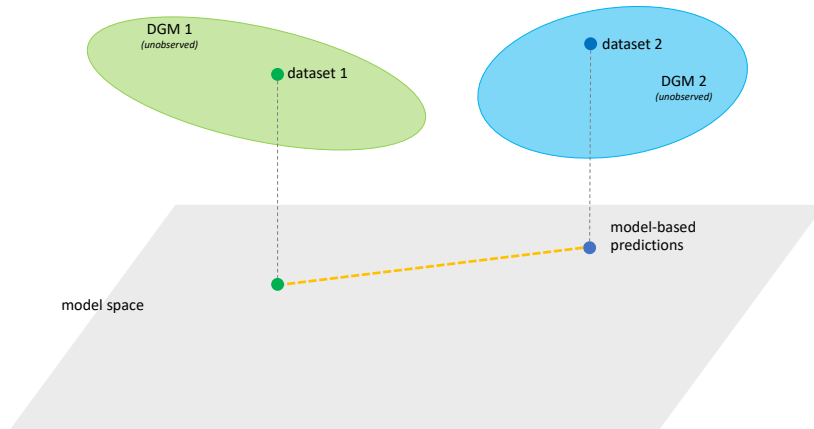
$$\Delta_{pred} = h(q_{j_1}, q_{j_2})$$

446 where  $q_{j_1}^{(i)}$  is the vector of cross-validation (dis)agreement statistics and  $q_{j_2}^{(i)}$  is the vector of  
447 validation (dis)agreement statistics for the  $i$ th experimental framework.

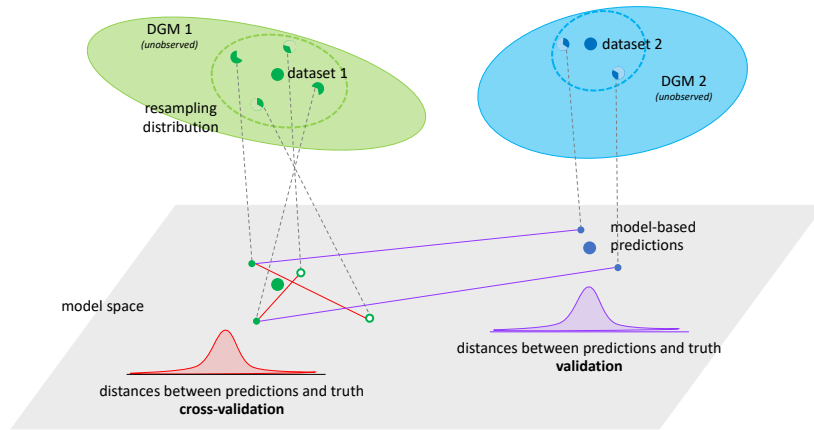
448 Note that for each particular application, appropriate choices for  $f$  (measure of the  
449 true difference between the data generating mechanisms),  $g$  ((dis)agreement statistic for the  
450 cross-validation and validation predictions), and  $h$  (measure of the difference between the  
451 distributions of (dis)agreement statistics) must be made. As we have suggested above, these  
452 choices should be well motivated by the particular application. More specifically,  $f$  should  
453 be motivated by the form of the family of data generating mechanisms being considered,  
454 and  $g$  should be motivated by the researcher’s model fitting software. Additionally, the  
455 choice of  $h$  should be motivated by both of these considerations and the subsequent choices  
456 for  $f$  and  $g$ . Although this methodology would be simpler if  $f, g$  and  $h$  were universally  
457 specified, it is important that they appropriately capture the important features of the  
458 data generating mechanisms and are suitable to whatever model fitting software is chosen  
459 by the researcher (see Table 1 for some specific examples). Further, note that this sort  
460 of conditional specification is not unlike the choice of an appropriate link function for  
461 generalized linear models. Appropriate forms of  $f, g$ , and  $h$  may be derived for more  
462 complex settings (e.g, dependent data, such as networks or time series) in the future.

## 463 3.2 Example: Linear regression

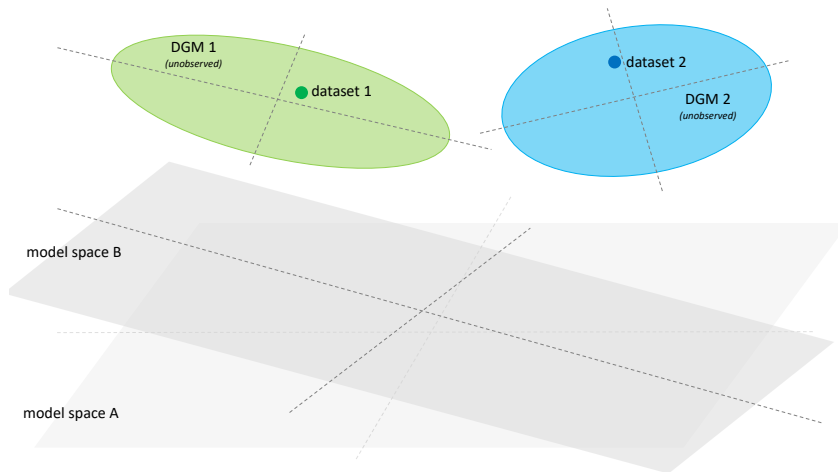
464 To better understand this methodology, we turn now to an example in the predictive case, a  
465 Bayesian linear regression model, documented in Figure 4. In this setting,  $p_j^{(i)} = p_j^{(i)}(y|x)$  is  
466 the conditional distribution of the outcome or response variable,  $y$ , given fixed values of the



(a) Differences between unobserved DGMs can be measured indirectly through model-based predictions.



(b) Cross-validation and validation account for sampling variability and separate model fit issues from true DGM differences.



(c) Prediction scores measure the distance between DGMs along the dimension of the model used to make predictions.

Figure 3: Geometric interpretation of the prediction scores.

467 predictors or covariates,  $x$ . The true difference between the data generating mechanisms  
468 is the conditional KL-divergence. In this Bayesian setting, predictions are draws from the  
469 posterior predictive distribution. For cross validation, this distribution,  $p_{*|-k}^{(1)}$  is conditioned  
470 on the set of data (covariates and responses) from experiment 1 with the  $k$ th subset re-  
471 moved and provides a prediction for the  $k$ th subset of responses, corresponding to the  $k$ th  
472 set of covariates in experiment 1. For validation, the posterior predictive distribution,  $p_{*|1}^{(1)}$   
473 is conditioned on the entire set of data (covariates and responses) from experiment 1 and  
474 provides a prediction for the  $k$ th response, corresponding to the  $k$ th covariate in experi-  
475 ment 2. Finally, we estimate the true difference between the data generating mechanisms  
476 by calculating the KL-divergence between the distributions of (dis)agreement statistics  
477 across cross-validation,  $\mathbf{q}_1^{(1)}$ , and validation,  $\mathbf{q}_2^{(1)}$ . Motivated by the Bayesian software-  
478 checking approaches of Cook et al. (2006) and Talts et al. (2018), a natural choice for the  
479 (dis)agreement statistics might be empirical quantiles or rank statistics.

### 480 3.3 Example: Logistic regression

481 To understand how this methodology can be used for inferential DGMs, consider the case  
482 where we assume the underlying process follows a simple logistic regression.

## 483 4 Probabilistic behavior of prediction scores

484 To understand how these prediction scores behave in practice and to get a sense of their  
485 asymptotic behavior, we have designed a simulation study that utilizes a simplified experi-  
486 mental design and models the outcome of interest with logistic regression. Many aspects of  
487 this simulation study were designed to complement related research that examines predic-  
488 tion scores for human behavior data in experimental social science research (Smith et al.,  
489 2018). We summarize the set up of this simulation study below and will detail how this  
490 study has been extended here to better examine the general probabilistic behavior of our  
491 prediction scoring methodology.

492 In Smith et al. (2018), we consider  $K = 5$  settings of a public goods game in which  
493 each participant has the opportunity to contribute (“cooperate”) or not (“defect”) to a set  
494 of pooled resources that will be multiplied and shared among all participants. Additionally,  
495 we imagine that some percentage of the total number of players,  $\pi = \{0, 0.25, 0.50, 0.75, 1\}$ ,  
496 are in fact bot participants whose behavior is strictly specified according to some set of  
497 algorithmic rules. The goal of these hypothetical experiments is to understand the ways in  
498 which participants’ decisions to cooperate are influenced by the presence of bots.

**True DGM.** Let  $y_{ijkt}$  be the decision to cooperate ( $y_{ijkt} = 1$ ) or defect ( $y_{ijkt} = 0$ ) for the  
 $i$ th individual in the  $j$ th cohort of the  $k$ th experimental setting for round  $t$ . Additionally,  
let  $z_{ijk}$  be an indicator of whether the  $i$ th participant in the  $j$ th cohort of the  $k$ th round  
is a human participant ( $z_{ijk} = 1$ ) or a bot ( $z_{ijk} = 0$ ). We will assume the true underlying data



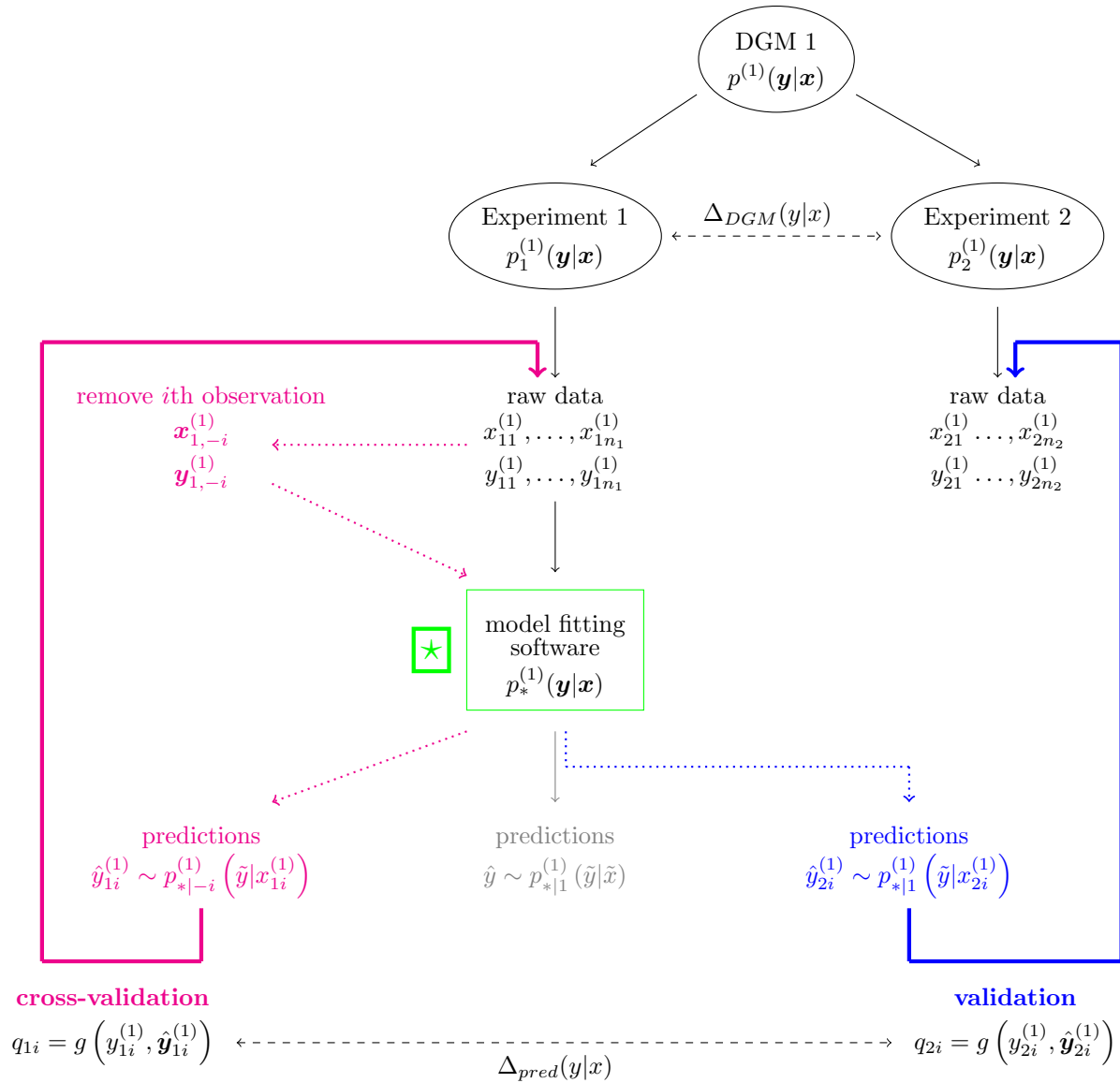


Figure 4: Outline of the procedure for a predictive data generating mechanism, such as linear regression.

generating mechanism is given by the following:

$$z_{ijk} \stackrel{iid}{\sim} \text{Bernoulli}(\pi_k)$$

Model 0:  $\text{logit}^{-1} [P(y_{ijkt} = 1 | z_{ijk} = 1)] = \beta_0 + \beta_1 t + \beta_2 y_{ijk,t-1} + \beta_3 \bar{y}_{.jk,t-1}$   
 $\text{logit}^{-1} [P(y_{ijkt} = 1 | z_{ijk} = 0)] = \beta'_0 + \beta'_2 y_{ijk,t-1}$

499 where  $\pi_k$  is the percentage of bots in the  $k$ th round,  $\beta_0$  and  $\beta'_0$  are baseline tendencies  
500 to cooperate,  $\beta_1$  captures any trend across the rounds,  $\beta_2$  and  $\beta'_2$  capture the tendency  
501 to switch between behaviors, and  $\beta_3$  represents the influence of team members' decisions.  
502 Values for these parameters for the simulated data are provided and motivated in Smith  
503 et al. (2018).

504 **Prediction scoring details.** In this setting, the DGMs being compared are predictive  
505 conditional distributions which we can refer to by  $p(\mathbf{y}_k | \mathbf{x}_k, \pi_k)$ . We perform this analysis in  
506 a Bayesian setting, so that predictions are draws from the posterior predictive distribution.

507 To compare these predictions to the set of true observations, we compute receiver  
508 operating characteristic (ROC) curves and the corresponding area under the curve (AUC)  
509 statistics (Davis and Goadrich, 2006). These measures are very popular model fit as-  
510 sessment tools for logistic regression. In order to compute these measures, we use  $L$ -fold  
511 cross-validation where  $L$  is chosen such that each partition contains roughly 500 observa-  
512 tions.

513 **Researcher models.** To uncover true differences across the experimental settings, we  
514 consider the following three researcher models:

$$\begin{aligned} \text{Model 1: } \text{logit}^{-1} [P(y_{ijkt} = 1)] &= \gamma_0 + \gamma_1 t, \\ \text{Model 2: } \text{logit}^{-1} [P(y_{ijkt} = 1)] &= \gamma'_0 + \gamma_2 y_{ijk,t-1}, \\ \text{Model 3: } \text{logit}^{-1} [P(y_{ijkt} = 1)] &= \gamma''_0 + \gamma_3 \bar{y}_{.jk,t-1}, \end{aligned}$$

515 where  $\gamma_0$  is a baseline tendency to cooperate,  $\gamma_1$  can capture some trends across the rounds,  
516  $\gamma_2$  represents the influence of of the most recent decision, and  $\gamma_3$  represents the influence  
517 of team members' decisions.

518 Smith et al. (2018) provide interpretations of visual differences in the ROC curves  
519 across the different models and experimental settings. To summarize these results, the  
520 prediction scores behave as expected; they appear similar when comparing data generated  
521 from the same DGM and appear more different as the distance between DGMs (here,  
522 measured simply in terms of  $|\pi_i - \pi_j|$ ) increases. This demonstrates that prediction scores  
523 can be used to uncover features of the DGM that vary across experimental settings, in a way  
524 that properly accounts for sampling variability. Additionally, the results of the simulation  
525 study indicate that Model 1 is the most sensitive to differences across the experimental  
526 settings. This is well-aligned with boxplots of the cooperation rate by round across each  
527 setting. In other words, when the model is aligned with true differences between the data

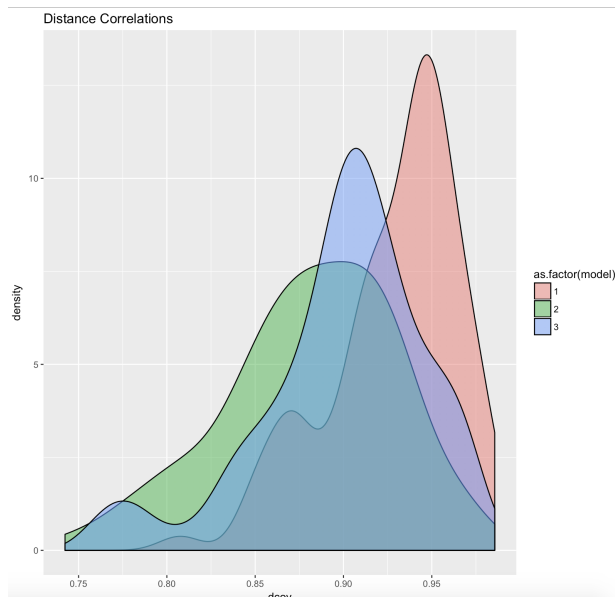


Figure 5: Distance correlations for prediction scores.

528 generating mechanisms, the distance between the cross-validation and validation statistics  
 529 reflects the true distance between the DGMs. In practice, relevant data patterns may  
 530 be much more nuanced (i.e., not obvious from simple summary plots) and the true data  
 531 generating mechanisms may be much more complex (i.e., it may be much more difficult to  
 532 specify a model that predicts well).

533 **Extension to study probabilistic behavior** . In order to get a sense of how these  
 534 prediction scores behave asymptotically, we repeat the above simulation study many times  
 535 and examine the relationship between the true distance between DGMs and our prediction  
 536 scoring estimates of that distance. This requires defining a true distance between the data  
 537 generating mechanisms. In this extended simulation, we consider two measures: (1) the  
 538 difference between the percentage of bots,  $|\pi_i - \pi_j|$ , and (2) the conditional KL-divergence,  
 539 calculated as follows:

$$KL(p_1, p_2) =$$

540 To evaluate whether or not the prediction scoring estimates are well-aligned with these  
 541 measures of the true underlying distance, we calculate distance covariances (Székely et al.,  
 542 2007). A distance covariance is a measure of dependence between two paired vectors that is  
 543 capable of detecting both linear and nonlinear associations. If the vectors are independent,  
 544 then the distance covariance is zero. We can treat each repetition of the above simulation  
 545 study (where we compute prediction scores across all possible pairs of  $\pi$ ) as a sample  
 546 which gives rise to a vector of prediction scoring distance estimates. Then we examine the  
 547 distribution of distance covariances, as a function of the (researcher) model used to make  
 548 predictions. After repeating this simulation 1000 times, we plot the distance covariances  
 549 in Figure 5.

## 5 Network experiments in cooperative games

The experiments proposed by the research teams in the (currently ongoing) Next Generation Social Science (NGS2) program present a great opportunity to evaluate the proposed prediction scoring framework. This program funds multiple research teams over two cycles of experiments and is designed as a methodologically-focused effort to develop a fundamental reimagining of the social science research cycle (Nosek et al., 2018). During each cycle, each research team will conduct distinct experimental social science studies regarding a shared research question. Prior to any data collection, each team will complete preregistration materials, which includes predictions for study outcomes. In the following, we briefly describe the Gallup teams' experiments for the first cycle of the program (for more detailed descriptions of each team's planned and completed research, see the preregistration materials which have been made publicly available on the Open Science Framework Nosek et al., 2018).

In the first cycle of the NGS2 program, the Gallup team provided an excellent application for our proposed prediction scoring methodology since their preregistered materials included pilot data from a previous study which informed their study hypotheses. This allows for an intuitive application of our proposed prediction scoring framework where we can compare the agreement between predictions for experimental data (based on the Gallup team's proposed modeling strategy and their identified pilot data) and the materialized observations from the experiment itself.

**Experimental setting** The first cycle of the NGS2 program focused on identifying pathways towards the formation of collective identity and cooperative decisions. To address this research question, the Gallup team considered the role of social networks in the development of large-scale cooperation among individuals in an economic game. They used a logistic regression model to examine individuals decisions (cooperation or defection) and showed that social networks which can be frequently updated by participants (rather than fixed throughout the course of the game or randomly updated) foster cooperative decisions in this setting. The Gallup teams experiments were designed to mimic the experiments performed by Rand et al. (2011), and whose data can serve as a set of preregistration data.

Experimenters randomly assigned participants to one of four conditions (see below) in a series of realizations of network experiments. In all conditions, subjects play a repeated cooperative dilemma (each game/session consists of multiple rounds) in an artificial social network created in the virtual laboratory. During each round of the game, each player can choose one of the following two actions: (1) cooperation: donate 50 units per neighbor, resulting in each neighbor actually gaining 100 units and (2) defection: donate nothing, resulting in neighbors getting nothing. After each round, players learn about the decisions of their neighbors and their own payoff. Additionally, the experimenters considered the following possible link-updating regimes for the social network in the game: (1) static or fixed links, (2) random link updating, where the entire network is regenerated at each round, (3) strategic link updating, where a randomly selected actor of a randomly selected pair may change the link status of that pair. The strategic link updating condition was

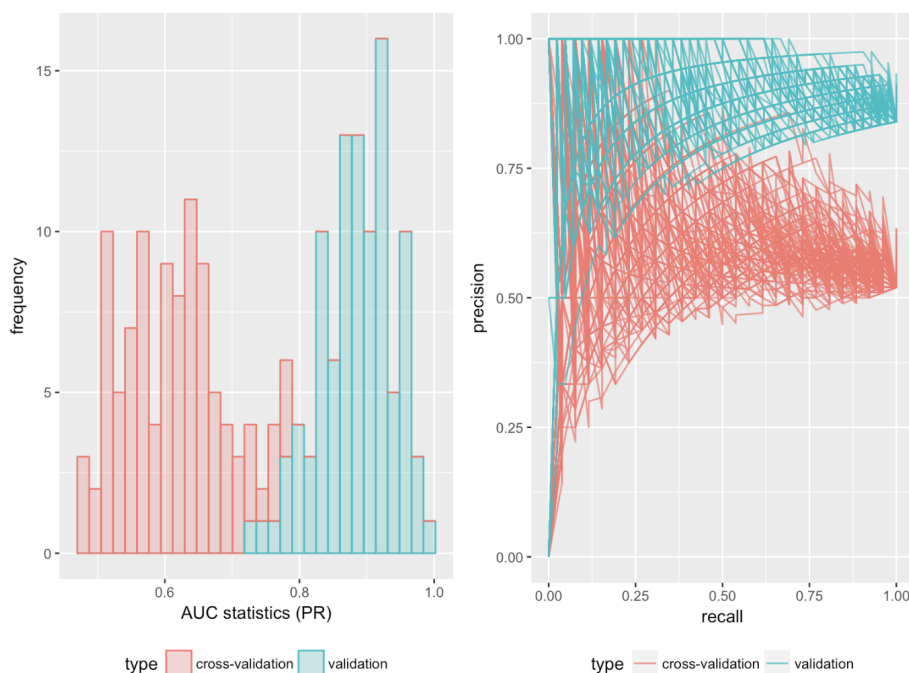


Figure 6: Prediction scores for Gallup’s Cycle 1 Hypothesis 1.4: rapidly updating networks support cooperation, relative to all other conditions.

591 further split into two categories: (a) viscous, where 10% of the subject pairs were selected  
 592 and (b) fluid, where 30% of the subject pairs are selected.

593 **Prediction scoring** Recall, that we have suggested using quantiles or rank statistics  
 594 as a disagreement statistic in our proposed prediction scoring methodology. However, for  
 595 logistic regression, observations and predictions will be collections of 0s and 1s. Thus,  
 596 using quantiles doesn’t make sense in this setting. Instead, we can use the ROC (receiver  
 597 operating characteristic) curve or precision-recall curve and the AUC (area under the curve)  
 598 statistic to measure the agreement between observations and predictions. These measures  
 599 are very popular model fit assessment tools for logistic regression. Thus, rather than  
 600 comparing quantile distributions across cross-validation and validation, we will compare  
 601 the distribution of AUC statistics across these settings. For this particular dataset, we will  
 602 calculate the AUC statistic for the precision-recall curve. Generally, the precision-recall  
 603 curve is preferred over the ROC curve whenever the data is imbalanced (see Davis and  
 604 Goadrich, 2006, for more discussion). Finally, we need to point out that the AUC statistic  
 605 is not defined for a single data point. Thus, we can not use leave-one-out cross-validation  
 606 in our prediction scoring routine. Instead we partition the dataset into  $k$  subsets and use  $k$ -  
 607 fold cross-validation, resulting in  $k$  AUC statistics. Similarly, when performing validation,  
 608 we must again partition the data into subsets.

609 As an example, consider Hypothesis 1.4 from the Gallup team’s preregistration mate-  
 610 rials. They hypothesized that rapidly updating networks would support cooperation more

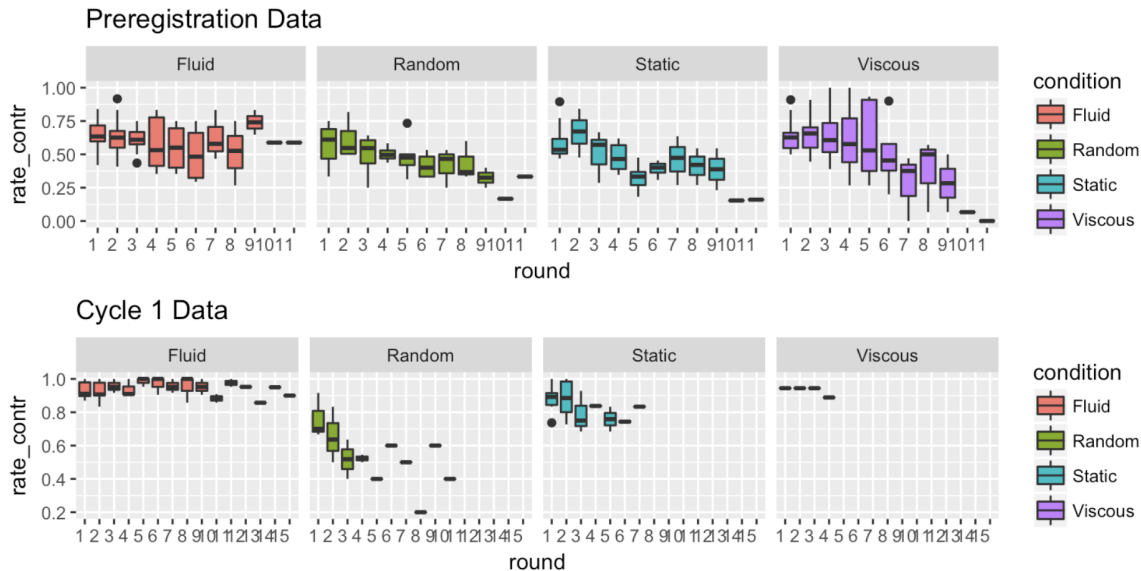


Figure 7: Boxplots of average cooperation levels across rounds of Gallup’s Cycle 1 games.

611 than any other condition. To evaluate the prediction scores, we have compared the distri-  
 612 bution of AUC statistics from cross-validation to those from validation as well as plotted  
 613 the corresponding precision-recall curves from each of the  $k$  subroutines of cross-validation  
 614 and validation (see Figure 6). The validation statistics measure differences due to both ran-  
 615 dom noise and true differences between the DGMs (i.e., between our prior beliefs about the  
 616 preregistered data and reality), while the cross-validation statistics only capture differences  
 617 due to random noise or disagreement between the underlying DGM and the chosen model.  
 618 In this case, we observe larger AUC statistics and better ROC curves in the validation  
 619 routine. This indicates that there is less variability in individuals’ behavior in the experi-  
 620 mental data, than in the preregistration data. In a sense, subjects in the Gallup experiment  
 621 are acting in more predictable ways than the subjects from the previous experiment. And  
 622 in fact, if we simply examine summary statistics of the in-game decisions themselves, we  
 623 can see the same type of pattern. In Figure 7, we provide boxplots of individuals’ average  
 624 cooperation levels across rounds of the game, where each color corresponds to a different  
 625 link-updating experimental condition. Comparing the preregistration data (top row) to  
 626 the experimental data (bottom row), we see that the boxplots are drastically narrower,  
 627 indicating that there is less variability in participant behavior.

628 This application serves as an illustration of how our prediction scoring can enable  
 629 interesting scientific insights. Further, we demonstrated how this methodology can be  
 630 adapted to appropriately address modelling choices (i.e., using distributions of AUC statis-  
 631 tics, rather than quantiles) and demonstrates the type of diagnostic plots that can be used  
 632 to interpret the resulting prediction scores.

## References

- 634 Anderson, J. L. “A method for producing and evaluating probabilistic forecasts from  
635 ensemble model integrations.” *Journal of Climate*, 9(7):1518–1530 (1996).
- 636 Cook, S. R., Gelman, A., and Rubin, D. B. “Validation of software for Bayesian models  
637 using posterior quantiles.” *Journal of Computational and Graphical Statistics*, 15(3):675–  
638 692 (2006).
- 639 Czado, C., Gneiting, T., and Held, L. “Predictive model assessment for count data.”  
640 *Biometrics*, 65(4):1254–1261 (2009).
- 641 Davis, J. and Goadrich, M. “The relationship between Precision-Recall and ROC curves.”  
642 In *Proceedings of the 23rd international conference on Machine learning*, 233–240. ACM  
643 (2006).
- 644 Dawid, A. P. “Present position and potential developments: Some personal views: Statis-  
645 tical theory: The prequential approach.” *Journal of the Royal Statistical Society. Series*  
646 *A (General)*, 278–292 (1984).
- 647 Gelman, A. “Preregistration of studies and mock reports.” *Political Analysis*, 21(1):40–41  
648 (2013).
- 649 Gelman, A., Hwang, J., and Vehtari, A. “Understanding predictive information criteria for  
650 Bayesian models.” *Statistics and computing*, 24(6):997–1016 (2014).
- 651 Gelman, A. and Loken, E. “The statistical crisis in science.” *American Scientist*, 102:460  
652 – 465 (2014).
- 653 Gneiting, T. “Making and evaluating point forecasts.” *Journal of the American Statistical*  
654 *Association*, 106(494):746–762 (2011).
- 655 Gneiting, T., Balabdaoui, F., and Raftery, A. E. “Probabilistic forecasts, calibration and  
656 sharpness.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*,  
657 69(2):243–268 (2007).
- 658 Gneiting, T. and Katzfuss, M. “Probabilistic forecasting.” *Annual Review of Statistics and*  
659 *Its Application*, 1:125–151 (2014).
- 660 Gneiting, T., Raftery, A., Balabdaoui, F., and Westveld, A. “Verifying probabilistic fore-  
661 casts: Calibration and sharpness.” In *Preprints, 17th Conf. on Probability and Statistics*  
662 *in the Atmospheric Sciences, Seattle, WA, Amer. Meteor. Soc*, volume 2 (2004).
- 663 Gneiting, T. and Raftery, A. E. “Strictly proper scoring rules, prediction, and estimation.”  
664 *Journal of the American Statistical Association*, 102(477):359–378 (2007).
- 665 Hamill, T. M. and Colucci, S. J. “Verification of Eta–RSM short-range ensemble forecasts.”  
666 *Monthly Weather Review*, 125(6):1312–1327 (1997).
- 667 Held, L., Schrödle, B., and Rue, H. “Posterior and cross-validators predictive checks: a  
668 comparison of MCMC and INLA.” In *Statistical modelling and regression structures*,  
669 91–110. Springer (2010).

- 670 Humphreys, M., Sanchez de la Sierra, R., and Van der Windt, P. “Fishing, commitment,  
671 and communication: A proposal for comprehensive nonbinding research registration.”  
672 *Political Analysis*, 21(1):1–20 (2013).
- 673 Li, L., Qiu, S., Zhang, B., and Feng, C. X. “Approximating cross-validators predictive  
674 evaluation in Bayesian latent variable models with integrated IS and WAIC.” *Statistics  
675 and Computing*, 26(4):881–897 (2016).
- 676 Mellor, D. “Registered reports library.” [www.zotero.org/groups/479248/osf/items/collectionKey/KEJP68G9](http://www.zotero.org/groups/479248/osf/items/collectionKey/KEJP68G9). Center for Open Science (2018).
- 678 Millar, R. B. “Conditional vs marginal estimation of the predictive loss of hierarchical  
679 models using WAIC and cross-validation.” *Statistics and Computing*, 28(2):375–385  
680 (2018).
- 681 Nosek, B. A., Spitzer, M., Russell, A., Tully, E., Rajtmajer, S., Ahn, S.-H., Zheng, T., Foy,  
682 D., Kluch, S. P., Stewart, C., and et al. “NGS2 DARPA Program.” (2018).
- 683 Open Science Collaboration. “Estimating the reproducibility of psychological science.”  
684 *Science*, 349(6251):aac4716 (2015).
- 685 Pearson, K. “On a method of determining whether a sample of size  $n$  supposed to have  
686 been drawn from a parent population having a known probability integral has probably  
687 been drawn at random.” *Biometrika*, 379–410 (1933).
- 688 Rand, D. G., Arbesman, S., and Christakis, N. A. “Dynamic social networks promote coop-  
689 eration in experiments with humans.” *Proceedings of the National Academy of Sciences*,  
690 108(48):19193–19198 (2011).
- 691 Rosenblatt, M. “Remarks on a multivariate transformation.” *The annals of mathematical  
692 statistics*, 23(3):470–472 (1952).
- 693 Shephard, N. “Partial non-Gaussian state space.” *Biometrika*, 81(1):115–131 (1994).
- 694 Smith, A., Zheng, T., and Gelman, A. “Evaluating drivers of human behaviors in experi-  
695 mental social science using prediction scoring.” (2018).
- 696 Spies, J. R. “Reproducibility Rubric.” (2018).
- 697 Stodden, V., McNutt, M., Bailey, D. H., Deelman, E., Gil, Y., Hanson, B., Heroux, M. A.,  
698 Ioannidis, J. P., and Taufer, M. “Enhancing reproducibility for computational methods.”  
699 *Science*, 354(6317):1240–1241 (2016).
- 700 Székely, G. J., Rizzo, M. L., Bakirov, N. K., et al. “Measuring and testing dependence by  
701 correlation of distances.” *The annals of statistics*, 35(6):2769–2794 (2007).
- 702 Talagrand, O., Vautard, R., and Strauss, B. “Evaluation of probabilistic prediction systems,  
703 paper presented at ECMWF Workshop on Predictability, Eur. Cent. for Med. Range  
704 Weather Forecasts.” *Reading, UK* (1997).



- 705 Talts, S., Betancourt, M., Simpson, D., Vehtari, A., and Gelman, A. “Validating  
706 Bayesian Inference Algorithms with Simulation-Based Calibration.” *arXiv preprint*  
707 *arXiv:1804.06788* (2018).
- 708 Vehtari, A., Gelman, A., and Gabry, J. “Practical Bayesian model evaluation using leave-  
709 one-out cross-validation and WAIC.” *Statistics and Computing*, 27(5):1413–1432 (2017).
- 710 Vehtari, A., Ojanen, J., et al. “A survey of Bayesian predictive methods for model assess-  
711 ment, selection and comparison.” *Statistics Surveys*, 6:142–228 (2012).
- 712 Wang, W. and Gelman, A. “Difficulty of selecting among multilevel models using predictive  
713 accuracy.” *Statistics at its Interface*, 7(1):1–88 (2014).