

The power of simulating data: a tool to design experiments, understand data limitations and improve scientific reasoning

Frederik Baumgarten^{1,2}, EM Wolkovich¹, invite Andrew Gelman

September 4, 2024

¹ Department of Forest and Conservation, Faculty of Forestry, University of British Columbia, 2424 Main Mall Vancouver, BC Canada V6T 1Z4.

² Swiss Federal Institute for Forest, Snow and Landscape Research WSL, Zürcherstr. 111, Birmensdorf 8903, Switzerland

Corresponding Author: Frederik Baumgarten; frederik.baumgarten@ubc.ca

Journal: statistical report in Ecology?

Abstract

Science continually seeks to explain patterns in nature that are often obscured by the noise of variation caused by a myriad of influencing factors. To disentangle pattern from noise and build upon the existing body of knowledge, researchers adhere to a scientific workflow. However, the replication crisis has highlighted serious problems about many aspects of the common scientific procedure: the formulation of meaningless (null) hypotheses, insufficient sample size, the problematic usage and interpretation of p-values that led to biased reportings and distortion of the literature.

Here, we propose the integration of data simulation into the scientific workflow to address these challenges by: 1) developing meaningful and testable hypotheses through the mathematical formulation of mechanistic models, 2) manipulating effect size, variance and sample size to generate synthetic/-fake datasets, 3) exploring the potential and limitations of both the underlying model and the data obtained, and 4) drawing conclusions that can reliably build upon the existing ‘body of knowledge’.

We expect that in a future dominated by powerful AI and machine learning algorithms applied to ever-larger datasets, meaningful, theoretically grounded hypotheses will be more important than ever for revealing underlying causalities. Without the understanding of causality, we may excel at making predictions within known limits, yet fail to understand the driving variables behind them.

Keywords: Bayesian statistics, statistical methods, quantitative ecology, hypothesis testing, scientific reasoning, experimental design, power analysis

1 Introduction

The constant search for patterns

A renowned researcher once explained to me the progression of a curve from a device attached to the bark of a tree next to us. As we watched the curve expand in real time on the screen, he gave a passionate and insightful explanation on the water relations of trees. A few moments later, however, the technician in charge entrusted me with the fact that the fluctuations on the graph emerged largely at random during the installation and calibration process and that the data displayed could not possibly represent anything biologically meaningful. This anecdote—which may sound familiar to many of us—illustrates 1) our tendency to interpret data in a way that supports our pre-existing beliefs, also known as confirmation bias (Nickerson, 1998), and 2) our incredible ability to recognise patterns even when there are none, or in other words, our poor intuition to judge what randomly generated data may look like. Human evolution has equipped us with a hypersensitive pattern detector and a set of cognitive biases (Tversky & Kahneman, 1974) that help us survive but get in the way of objective science. However, distinguishing between noise and pattern is perhaps the most important skill for any researcher.

How can we ensure scientific standards?

To overcome cognitive biases and objectify science, researchers usually follow a scientific workflow illustrated in Figure 1) that traditionally starts with identifying a research question and formulating a hypothesis and ends with a conclusion based on the statistical outcome given the data obtained and the hypothesis at hand (see for example Schwab *et al.* (2022)).

However, there is growing evidence that following this workflow is not enough. Serious concerns have been raised in the last decades that jeopardize/question scientific integrity and credibility and may slow down scientific progress despite of a record breaking and still increasing publishing rate (REF). Or as Ioannidis (2005) put it in his title: “Why most published research findings are false”.

Aim of this paper: An advocacy for data simulation

Here we propose to add a simulation step to the traditional scientific workflow that can mitigate current problems of replication and causal reasoning. We believe that translating descriptive (conceptual) hypotheses into mechanistic (mathematical) models and simulate data from them will i) greatly improve our mechanistic understanding of underlying processes, ii) lead to better experimental designs with adequate statistical power and iii) facilitate to move away from p-values and significance testing, towards evaluating model parameters such as effect sizes and error intervals.

Going through the traditional workflow we briefly outline some major concerns and pitfalls (see also Figure 1) and outline how data simulation through mechanistic models can help address and mitigate these issues. We provide concrete examples to showcase how this process may look like and include r-code for better illustration. We invite scientists across disciplines and career stages to consider and further explore this approach since we are convinced that simulation is a vital tool to do better science.

2 Problems with the current scientific workflow

Hypothesis and predictions

Beyond the purely exploratory nature of a study, a hypothesis well-founded in current theory is essential for causal inference and reasoning (Rajtmajer *et al.*, 2022). However, studies greatly differ in the usefulness of their stated hypotheses, which can range from non-existent to too vague, subjective,

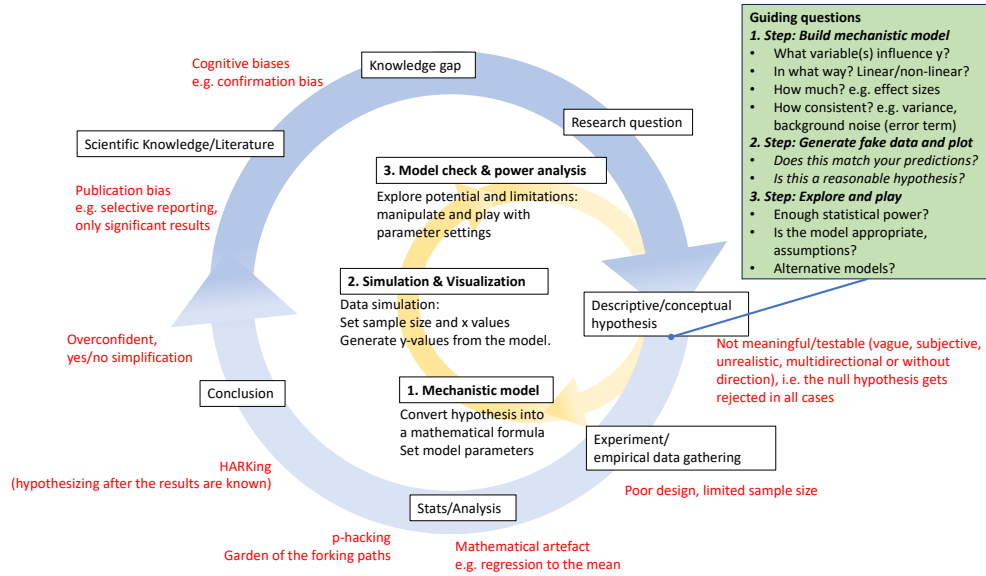


Figure 1: Schematic overview of the current scientific workflow (blue circle) with its related problems (red text) and how some of them can be anticipated with the help of data simulation (three-step-loop in yellow). The critical step is the translation of a verbal hypothesis into a mechanistic model expressed in mathematical terms. We believe that this 'simulation loop' fosters a better understanding of underlying processes founded in theory, clear predictions that can be evaluated and overall better science. The questions summarized in the green box should be addressed in the simulation loop before continuing on the blue circle, e.g. prior to the start of real data gathering/exploration.

untestable, multidirectional, or lacking any direction. Additionally, the common practice of testing to falsify or reject a null hypothesis has been shown to be of limited utility in most cases. This is because null hypotheses often oversimplify complex relationships into a binary outcome and are frequently rejected due to the myriad of influencing factors contributing to data variability (Rajtmajer *et al.*, 2022). Null hypotheses are particularly problematic when combined with the common practice of significance testing (NHST; null hypothesis significance testing), making them unsuitable as a cornerstone in scientific research. (Szucs & Ioannidis, 2017), see also the problems involved with p-values and 'significance' below.

Experimental design and sample size

Many studies are under-sampled. More precisely, they lack statistical power and are prone to find statistically significant results by chance when in fact there is none (Halsey *et al.* (2015) see below), which is a main cause of the replication crisis across many disciplines with vast number of studies not being able to reproduce previously published results (Ioannidis, 2005; Baker, 2016; Camerer *et al.*, 2018). Low sample size is of particular concerns when models are loaded with parameters including several interaction terms REF. On the other hand, the common emphasis on p-values gave rise to a multitude of studies that found highly significant findings with very high sample size, blending out the relevance of the effect size. If eating sausages significantly shortens your life span, the most important information is not the p-value but the effect size. We would like to know by how much. A two days shorter life expectancy is simply irrelevant but equivalent results are published frequently REF.

94 Problematic use of p-values

95 A widely recognized problem, despite its prime role in statistical inference arise from the (excessive) use-
96 age of p-values (Amrhein *et al.*, 2019b, 2017; Halsey *et al.*, 2015). Many studies misinterpret ‘statistical
97 significance’ as direct proof for both relevancy and certainty (REF). Even if p-values are interpreted
98 correctly the multitude of tests often performed in a single study can lead to a ‘fishing strategy’ also
99 known as p-hacking, that increases the chance of getting a statistically significant result (Stefan &
100 Schönbrodt, 2023). Finding a $p < 0.05$ during the process of exploring and analyzing data can occur
101 largely unintentional. Most analysis can easily find hundreds of possible ways and combinations to find
102 a potentially relevant pattern, e.g. by including/excluding variables or subgroups. Gelman & Loken
103 (2013) illustrated this problem with the metaphor of ‘the garden of the forking paths’ pointing at the
104 importance of study pre-registration and the awareness of the choices we make during the analysis
105 (Rubin, 2020; ?). Hypothesizing after the results are known (HARKing) is an additional pitfall that
106 sometimes dangerously attaches causal reasoning to significant results (?). Finally, these mentioned
107 problems feed into a reporting and publication bias (Yang *et al.*, 2023; Lin & Chu, 2018) that can
108 substantially distort the literature as was shown for meta-analyses (Van Zwet & Cator, 2021). Over-
109 all, there is a trend to abandon p-values and to retire ‘statistical significance’ as an outdated concept
110 (Amrhein *et al.*, 2019a; Berner & Amrhein, 2022; McShane *et al.*, 2019; Woolston, 2015; Wasserstein
111 *et al.*, 2019; Lee, 2016).

112 Mathematical artefacts

113 In some cases the way data are analyzed or displayed create patterns, that arise from mathematical
114 properties, sometimes hard to understand. A famous example is the regression to the mean that is
115 responsible for many reported effects that were later found to be overestimated or not present at all (e.g.
116 the famous and equally mystic effect of ‘unskilled and unaware of it’ by Kruger & Dunning (1999)).
117 A more recent example of mathematical artefact is the declining sensitivity of biological responses
118 to climate warming. Declining rates measured as day of advancement per degree of warming, might
119 largely reflect a mathematical artefact when variables are not properly transformed (Wolkovich *et al.*,
120 2021).

121 3 The three steps of data simulation

122 Here, we propose a small addition/adjustment to the traditional scientific workflow: The integration
123 of data simulation prior to data gathering or exploration (Figure 1). Concretely, this involves the
124 following three steps:

125 Step 1: From a conceptual hypothesis to a mechanistic model

126 Models are hypotheses. With each model we apply to our data, we examine/evaluate—often not with
127 full awareness—an additional hypothesis. To constrain and justify this process to a reasonable set of
128 hypotheses, we propose to build up mechanistic models by incorporating existing knowledge about the
129 influencing variables and their relationships in question. Which explanatory variables are influencing
130 y? Is it linear or non-linear? What do we know from the literature? This step nudges us to think
131 carefully about our hypothesis in terms of model parameters—a great way to translate subjective and
132 conceptual hypothesis to objective mathematical formulations. This step can be challenging sometimes
133 but is very vital: in our experience, a conceptual hypothesis often intuitively makes sense until we try
134 to translate it into a mechanistic model. Soon it will become clear, that the original hypothesis lacks
135 direction, concreteness, or foundation in basic knowledge. Beyond the classical descriptive hypothesis,
136 a mechanistic model forces us to think even further: What are the model parameters we would like
137 to estimate? And what values we expect for those parameters, e.g. for effect sizes and their variance.
138 This gives us the opportunity to nail a hypothesis with numbers and mathematical formulas.

Step 2: Simulate and visualize fake data

"The greatest value of a picture is when it forces us to notice what we never expected to see." This quote from Tukey (1977) perfectly captures the need for the second step. With your mechanistic model sorted out, you are ready to generate a synthetic/fake data set from it. Here is the important difference: you don't want to simulate your response variable directly from a suitable distribution. Instead, simulate your response variable using your set of exploratory variables (your model parameters), which you believe capture the ecological processes that ultimately shape the observed response. Choose a sample size, create data for your x-variables and add an error (the noise in your data, see more on that in the example below). By plotting your generated data, you will see precisely your hypothesis in a visualized form, nuanced not only by a clear direction of the effect, but also by its magnitude (e.g. effect size) and consistency (variance). We believe that many data papers would greatly benefit from visualizing their mechanistic model predictions to effectively and objectively validate and communicate their hypotheses.

Step 3: Power analysis and model exploration

We don't propose a common power analysis based on a significance threshold to calculate the probability of correctly rejecting a null hypothesis when it is false ($1 - \beta$). There is nothing wrong with that approach (besides the problems with significance testing as discussed above) and we even think this is a good start. However, here we propose to go further: Explore your model from your simulated data. Is it doing what you think it is? Does your predictions meet your expectations? Your data show no variation at all? That's likely unrealistic—so you might want to include an error term in your model. Re-generate your fake dataset by manipulating sample size, effect size(s), error-terms and relationships within the mechanistic model to explore how your model behaves. This is a playful, yet powerful way to design an experiment and understand how your model behaves. You will find out about the potentials and limits of both your model and your data set. Besides assessing the statistical power we suggest to focus on your parameters. You can even fit your model to your fake data to check if the parameter estimates turn out to be in the range you originally set them. That's a great way to check model assumptions before you fit empirical data (see Wolkovich *et al.* (2024) for more details). In the following we go through some concrete examples and show how this procedure can look like.

4 Practical example: Plant growth and nitrogen

Data simulation is best shown in practice. The example presented here along with more complex ones can also be performed through the R script available in the supplement.

Every study starts with a research question. So let's start with this one: What is the influence of nitrogen fertilization on plant growth? We might want to perform an experiment on sunflowers and fertilize some with nitrogen to assess the biomass after one season. So our dependent variable y is biomass and we might have a clunky hypothesis that plants grow more when fertilized with nitrogen.

What influences y in what way?

Here we can dive deep into the literature. What do we know that influences plant growth? Let's pick nitrogen concentration for simplicity to start with. Let's further assume there is a positive linear relationship of plant growth and nitrogen concentration. So we can express this relationship in a basic linear model with an intercept a (e.g. the biomass at zero nitrogen), a slope b (e.g. the increase of biomass per 1 unit increase in nitrogen) and an error term that captures the natural variability (otherwise your predictions will fall on a perfect line):

```
# Building up a mechanistic model
y = a+b*x + error #model formula with x being nitrogen concentration
error = rnorm(n,0,sigma) #error distribution
```

184 Note that we chose an error distribution that is drawn from a normal distribution with mean 0 and a
185 standard deviation of sigma.

186 Setting model parameters

187 This is the where things get interesting and concrete: What do you think your model parameter should
188 look like? In other words, what biomass do we expect under no nitrogen fertilization? What is the
189 increase in biomass per unit increase in nitrogen? And how variable is this relationship, e.g. what is
190 the background noise? Lets put some numbers here:

```
191     # specifying model parameters:
192     a <- 30 #intercept
193     b <- 7 #slope (effect size)
194     sigma <- 5 #standard deviation of the error distribution.
```

195 Simulate!

196 In order to simulate we only need to set sample size and create a meaningful range of x values, i.e.
197 nitrogen concentrations:

```
198     n <- 50 #sample size
199     x<-rnorm(n,10,4) #create data for nitrogen concentration
200     y<-a+b*x + rnorm(n,0,sigma) #generate y values
201     fake<-data.frame(x,y)
```

202 Visualize and assess your fake data

203 Now lets plot the generated data and see if that matches our expectations. This is basically our
204 hypothesis as concrete as it can be.

```
205     plot(fake$x, fake$y, main="Fake data")
```

206 We can even fit our model and check if the model parameters are in the range we expect them to be
207 and plot them.

```
208     fit_1<-stan_glm(y~x, data=fake) # fit the model
209     print(fit_1, digits=2)
210     plot(fake$x, fake$y, main="Data and fitted regression line")
211     a_hat<-coef(fit_1)[1]
212     b_hat<-coef(fit_1)[2]
213     abline(a_hat, b_hat)
```

214 How to use this - Play!

215 Going back and forth we can use this tool to better understand what your model is doing. Is our effect
216 size reasonable? Is the sample size high enough to detect the effect with reasonable certainty? What if
217 the background noise is much higher? You may realized quickly that this loop will give you confidence
218 in your study design and analysis beyond a simple power analysis.

219
220 play with replication while holding variance and effect size constant. p-value figure

221

222 5 Artificial intelligence (AI) and the importance of causal rea- 223 soning in the future

224 The steep rise in machine learning has propelled recent advancements of artificial intelligence (AI).
225 Current applications excel human-level intelligence by far in many aspects (REF) offering new oppor-
226 tunities for science REF. Many disciplines have already benefited from these very recent advancements,
227 eg. XXX. However, technologies to date are limited to pattern recognition based on correlations only
228 (Pearl, 2019). The critical ingredient lacking is the ability of causal reasoning - a feature so far limited
229 to humans (Pearl, 2019). While it might be possible that we see further improvements towards logic
230 and reasoning in AI systems, the separation of causality from correlation will likely continue to play a
231 pivotal role in science and relies on human logical thinking. With ever growing datasets we will likely
232 find strong correlations that allow for accurate predictions within known limits but fail to answer the
233 why-questions. We believe that with the rise of AI we need to learn to formulate better questions and
234 hypotheses to improve causal reasoning and to build ontop of the house of knowledge.

References

- Amrhein, V., Gelman, A., Greenland, S. & McShane, B.B. (2019a). Abandoning statistical significance is both sensible and practical. Tech. Rep. e27657v1, PeerJ Inc.
- Amrhein, V., Greenland, S. & McShane, B. (2019b). Scientists rise up against statistical significance. *Nature*, 567, 305–307.
- Amrhein, V., Korner-Nievergelt, F. & Roth, T. (2017). The earth is flat ($p > 0.05$): Significance thresholds and the crisis of unreplicable research. *PeerJ*, 5, e3544.
- Baker, M. (2016). 1,500 scientists lift the lid on reproducibility. *Nature*, 533, 452–454.
- Berner, D. & Amrhein, V. (2022). Why and how we should join the shift from significance testing to estimation. *Journal of Evolutionary Biology*, 35, 777–787.
- Camerer, C.F., Dreber, A., Holzmeister, F., Ho, T.H., Huber, J., Johannesson, M., Kirchler, M., Nave, G., Nosek, B.A., Pfeiffer, T., Altmeld, A., Buttrick, N., Chan, T., Chen, Y., Forsell, E., Gampa, A., Heikensten, E., Hummer, L., Imai, T., Isaksson, S., Manfredi, D., Rose, J., Wagenmakers, E.J. & Wu, H. (2018). Evaluating the replicability of social science experiments in Nature and Science between 2010 and 2015. *Nature Human Behaviour*, 2, 637–644.
- Gelman, A. & Loken, E. (2013). The garden of forking paths: Why multiple comparisons can be a problem, even when there is no “fishing expedition” or “p-hacking” and the research hypothesis was posited ahead of time. *Department of Statistics, Columbia University*, 348, 3.
- Halsey, L.G., Curran-Everett, D., Vowler, S.L. & Drummond, G.B. (2015). The fickle P value generates irreproducible results. *Nature Methods*, 12, 179–185.
- Ioannidis, J.P.A. (2005). Why Most Published Research Findings Are False. *PLOS Medicine*, 2, e124.
- Kruger, J. & Dunning, D. (1999). Unskilled and unaware of it: How difficulties in recognizing one’s own incompetence lead to inflated self-assessments. *Journal of personality and social psychology*, 77, 1121.
- Lee, D.K. (2016). Alternatives to P value: Confidence interval and effect size. *Korean Journal of Anesthesiology*, 69, 555–562.
- Lin, L. & Chu, H. (2018). Quantifying Publication Bias in Meta-Analysis. *Biometrics*, 74, 785–794.
- McShane, B.B., Gal, D., Gelman, A., Robert, C. & Tackett, J.L. (2019). Abandon Statistical Significance. *The American Statistician*, 73, 235–245.
- Nickerson, R.S. (1998). Confirmation Bias: A Ubiquitous Phenomenon in Many Guises. *Review of General Psychology*, 2, 175–220.
- Pearl, J. (2019). The seven tools of causal inference, with reflections on machine learning. *Commun. ACM*, 62, 54–60.
- Rajtmajer, S.M., Errington, T.M. & Hillary, F.G. (2022). How failure to falsify in high-volume science contributes to the replication crisis. *eLife*, 11, e78830.
- Rubin, M. (2020). Does preregistration improve the credibility of research findings? *The Quantitative Methods for Psychology*, 16, 376–390.
- Schwab, S., Janiaud, P., Dayan, M., Amrhein, V., Panczak, R., Palagi, P.M., Hemkens, L.G., Ramon, M., Rothen, N., Senn, S., Furrer, E. & Held, L. (2022). Ten simple rules for good research practice. *PLOS Computational Biology*, 18, e1010139.

- 275 Stefan, A.M. & Schönbrodt, F.D. (2023). Big little lies: A compendium and simulation of p-hacking
276 strategies. *Royal Society Open Science*, 10, 220346.
- 277 Szucs, D. & Ioannidis, J.P.A. (2017). When Null Hypothesis Significance Testing Is Unsuitable for
278 Research: A Reassessment. *Frontiers in Human Neuroscience*, 11.
- 279 Tukey, J.W. (1977). Exploratory data analysis. *Reading/Addison-Wesley*.
- 280 Tversky, A. & Kahneman, D. (1974). Judgment under Uncertainty: Heuristics and Biases. *Science*,
281 185, 1124–1131.
- 282 Van Zwet, E.W. & Cator, E.A. (2021). The significance filter, the winner’s curse and the need to
283 shrink. *Statistica Neerlandica*, 75, 437–452.
- 284 Wasserstein, R.L., Schirm, A.L. & Lazar, N.A. (2019). Moving to a World Beyond “ $p < 0.05$ ”. *The*
285 *American Statistician*, 73, 1–19.
- 286 Wolkovich, E.M., Auerbach, J., Chamberlain, C.J., Buonaiuto, D.M., Ettinger, A.K., Morales-Castilla,
287 I. & Gelman, A. (2021). A simple explanation for declining temperature sensitivity with warming.
288 *Global Change Biology*, 27, 4947–4949.
- 289 Wolkovich, E.M., Davies, T.J., Pearse, W.D. & Betancourt, M. (2024). A four-step Bayesian workflow
290 for improving ecological science.
- 291 Woolston, C. (2015). Psychology journal bans P values. *Nature*, 519, 9–9.
- 292 Yang, Y., Sánchez-Tójar, A., O’Dea, R.E., Noble, D.W.A., Koricheva, J., Jennions, M.D., Parker,
293 T.H., Lagisz, M. & Nakagawa, S. (2023). Publication bias impacts on effect size, statistical power,
294 and magnitude (Type M) and sign (Type S) errors in ecology and evolutionary biology. *BMC*
295 *Biology*, 21, 71.