

Research



Cite this article: Buzbas EO, Devezer B, Baumgaertner B. 2023 The logical structure of experiments lays the foundation for a theory of reproducibility. *R. Soc. Open Sci.* **10**: 221042. <https://doi.org/10.1098/rsos.221042>

Received: 11 August 2022

Accepted: 2 February 2023

Subject Category:

Mathematics

Subject Areas:

statistics

Keywords:

reproducibility, replication, open science, metascience, experiment, statistical theory

Author for correspondence:

Erkan O. Buzbas

e-mail: erkanb@uidaho.edu

The logical structure of experiments lays the foundation for a theory of reproducibility

Erkan O. Buzbas¹, Berna Devezer^{1,2} and Bert Baumgaertner³,

¹Department of Mathematics and Statistical Science, ²Department of Business, and

³Department of Politics and Philosophy, University of Idaho, Moscow, ID 83844, USA

EOB, 0000-0003-1446-3447; BD, 0000-0002-5979-2781

The scientific reform movement has proposed openness as a potential remedy to the putative reproducibility or replication crisis. However, the conceptual relationship among openness, replication experiments and results reproducibility has been obscure. We analyse the logical structure of experiments, define the mathematical notion of idealized experiment and use this notion to advance a theory of reproducibility. Idealized experiments clearly delineate the concepts of replication and results reproducibility, and capture key differences with precision, allowing us to study the relationship among them. We show how results reproducibility varies as a function of the elements of an idealized experiment, the true data-generating mechanism, and the closeness of the replication experiment to an original experiment. We clarify how openness of experiments is related to designing informative replication experiments and to obtaining reproducible results. With formal backing and evidence, we argue that the current ‘crisis’ reflects inadequate attention to a theoretical understanding of results reproducibility.

1. Introduction

In a number of scientific fields, replication and reproducibility *crisis* labels have been used to refer to instances where many results have failed to be corroborated by a sequence of scientific experiments. This state of affairs has led to a scientific reform movement. However, this labelling is ambiguous between a crisis of practice and a crisis of conceptual understanding. Insufficient attention has been given to the latter, which we believe is a detriment to moving forward to conduct science better. In this article, we make theoretical progress towards understanding replications and

reproducibility of results (henceforth, ‘results reproducibility’)¹ by a formal examination of the logical structure of experiments.²

We view replication and reproducibility as methodological subjects of metascience. As we have emphasized elsewhere [10], these methodological subjects need a formal approach to properly study them. Therefore, our work here is necessarily mathematical; however, we make our conclusions relatable to the broader scientific community by pursuing a narrative form in explaining our framework and results within the main text. Mathematical arguments are presented in the appendices. Our objective is to build a strong, internally consistent, verifiable theoretical foundation to understand and to develop a precise language to talk about replication experiments and results reproducibility as well as to use this framework to study how openness in and of experiments is related to either of them. We advance mathematical arguments from first principles and proofs, using probability theory, mathematical statistics, statistical thought experiments and computer simulations. We ask the reader to evaluate our work within its intended scope of providing theoretical precision and nuanced arguments.

The following backdrop to motivate our research matters: a common concern voiced in the scientific reform literature and recent scholarly discourse regards various forms of scientific malpractice as potential culprits of reproducibility failures, and openness is sometimes touted as a remedy to alleviate such malpractices [11–15]. Some malpractice is believed to take place at the level of the scientist. For example, hypothesizing after the results are known involves presenting a *post hoc* hypothesis as if it were an *a priori* hypothesis, conditional on observing the data [16,17]. Another example is *p*-hacking, a statistically invalid form of performing inference to find statistically significant results [17–19]. Some is believed to operate at the community or institution level. For example, publication bias involves omitting studies with statistically non-significant results from publications and is primarily attributed to flawed incentive structures in scientific publishing [1,17]. Transparency in scientific practice in general and tools to promote openness in experimental reporting (such as preregistration, registered reports, and open laboratory books) in particular are often highlighted as potential remedies to curb such malpractice. Before we suspect malpractice of either kind and set out to correct the scientific record or demand reparations, however, it behoves the scientific community to gain a complete understanding of the factors that may account for a given set of results in a sequence of replication experiments. This way we can hope to understand what aspects of experiments need to be openly communicated and to what end.

If a result of an experiment is not reproduced by a replication experiment, before we reject it as a false positive or suspect some form of malpractice, we need to assess and account for: (i) sampling error, (ii) theoretical constraints on the reproducibility rate of the result of interest, conditional on the elements of the original experiment, and (iii) assumptions from the original experiment that were not carried over to the replication experiment. First of these is a well-known and widely understood statistical fact that describes why methodologically we can at best guarantee reproducibility of a result on average (i.e. in expectation). The second point about the theoretical limits of the reproducibility rate is not well understood, and we hope to address this oversight in this article. The last one has been brought up in individual cases but typically in an ad hoc manner, and we aim to provide a systematic approach for comprehensive evaluations of replication experiments. Since metascientific heuristics may lead us astray in these assessments, we need a fine-grained conceptual understanding of how experiments operate and relate to each other, and what role openness plays in facilitating replications or promoting reproducible results. Indeed a replication crisis and a reproducibility crisis are different things and should be understood on their own. We distinguish between replication experiments and results reproducibility, discuss precursors of each, and assess how openness of experiments relates to each separately.

¹We focus on the end products of experiments and the results and not other components of experiments that bring those products about. This choice is fitting given the etymology of the term ‘reproduce’ in the sense of producing a given result, and by its formal association with statistical theory. In our research, we have aimed to use *results reproducibility* or *reproducibility of results* consistently. This usage is not idiosyncratic or esoteric. Reproducibility has been defined in a similar (if less technical) fashion by other scholars [1–4]. Unfortunately, there is variation in usage of these terms in the metascientific literature [5,6]. For example, *replicability* may refer to what we call results reproducibility (e.g. [7]). Further, *reproducibility* may convey *computational reproducibility* of the results given the data; i.e. obtaining the same output when a computer code is re-run with fixed input (e.g. [8]). Here, we do not refer to computational reproducibility. Our context is statistical: the reproducibility of experimental results in replication studies. We sidestep the potential confusion by laying out the definitions as we have and adhering to them for the remainder of this article.

²Some of the ideas developed in depth here appeared in a preliminary form in [9].

In this article, we argue that ‘failed’ replications do not necessarily signify failures of scientific practice.³ Rather, they are expected to occur at varying rates due to the features of and differences in the elements of the logical structure of experiments. By using a mathematical characterization of this structure, we provide precise definitions of and clear delineation between replication, reproducibility and openness. Then, by using toy examples, simulations and cases from the scientific literature, we illustrate how our characterization of experiments can help identify what makes for replication experiments that can, in theory, reproduce a given result and what determines the extent to which experimental results are reproducible. In the next section, we define main notions that we use to build a logical structure of experiments which help us derive our theoretical results.

2. The logical structure of experiments

2.1. Definitions

The *idealized experiment* is a probability experiment: a trial with uncertain outcome on a well-defined set. A scientific experiment where inference is desired under uncertainty can be represented as an idealized experiment. The results from an experiment can be defended as valid only if the assumptions of the probability experiment hold. One useful set-up for us is as follows: given some background knowledge K (see table 1 for reference to all notation and terms introduced in this section) on a natural phenomenon, a scientific theory makes a prediction, which is in principle testable using observables, the data D . A mechanism generating D is formulated under uncertainty and is represented as a probability model M_A under assumptions A . Given D , inference is desired on some unknown part of M_A . The extent to which parts of M_A that are relevant to the inference are confirmed by D is assessed by a fixed and known collection of methods S evaluated at D (similar descriptions for other purposes can be found in [10,21]).

Definition 2.1. The tuple $\xi := (K, M_A, S, D)$ is an *idealized experiment*.

Definition 2.1 of ξ captures some key distinct elements of experiments whose population characteristics can in principle be tested. These elements are not necessarily independent of each other. For example, K may inform and constrain the sets of plausible M_A and S . Or it may be necessary for M_A to constrain S .

M_A includes the sampling design when sampling a population conforming A , which we assume to be independent of sampling design. For example, A may be the description of an infinite population of interest, which may be sampled in a variety of ways to yield distinct probability models M_A for the data depending on the sampling scheme.

We distinguish two elements of S : S_{pre} and S_{post} . S_{pre} is the *scientific* methodological assumptions made before data collection and procedures implemented to obtain D . S_{pre} captures assumptions in designing and executing an experiment such as experimental paradigms, study procedures, instruments and manipulations. Conditional on K and M_A , S_{pre} is *reliable* if the random variability in D is due only to sampling variability modelled by M_A . S_{post} is the *statistical* methods applied on D . If inferential, S_{post} is *reliable* if it is statistically consistent. S is reliable if and only if S_{pre} and S_{post} are reliable.

We also distinguish two elements of D : D_s and D_v . D_s is the structural aspects of the data, such as the sample size, number of variables, units of measurement for each variable, and metadata. D_v is the observed values, that is, a realization conforming D_s . Some statistical approaches to assess risk and loss focus on the reproducibility conditional on D_v , whereas others focus on averages over independent realizations of D_v .

Definition 2.1 of ξ allows us to scaffold other definitions as follows. An exact replication experiment ξ' must generate D' independent of D conditional on M_A in the values but with the same structure D_s .

Definition 2.2. The tuple $\xi' := (K', M'_A, S', D')$ is an *exact replication experiment* of ξ if $K' \supset K$, $M'_A \equiv M_A$, $S' \equiv S$, $D'_s \equiv D_s$, and D'_v is a random sample independent of D_v . If at least one of (M'_A, S', D'_s) differs from (M_A, S, D_s) or $K' \not\supset K$, then ξ' can at most be a *non-exact replication experiment* of ξ .

Definition 2.2 mathematically isolates ξ and ξ' from R , the result of interest as formally defined in definition 2.3. That is, ξ' does not need to have a specific aim to be performed or worked with as a

³We are not the first to take issue with the ‘replication crisis’ framing. We invite the interested reader to visit Feest’s [20] provocative and incisive assessment of why replication is overrated.

Table 1. Quick reference guide to notation and key terms.

symbol	name	description	formal definition/result
ξ	idealized experiment	scientific experiment represented as a probability experiment	definition 2.1
ξ'	replication experiment	idealized experiment aiming to reproduce R from another experiment by generating D' independent from D	appendix B, definition 2.2, result 2.9
K	background knowledge	state of scientific knowledge on the phenomenon of interest used to conceptualize, design and perform the experiment	appendix C, remark 4.1
M_A	assumed model	assumed mechanism generating the data	appendix D, result 4.2
A	population assumptions	population characteristics independent from sampling design, such as finiteness and continuity	
M	model specification	model properties that depend on researcher assumptions, such as sampling scheme	
S	method	fixed and known set of methods for collecting and analysing data	
S_{pre}	pre-data methods	scientific assumptions made before collecting data and procedures implemented to obtain D	result 4.3
S_{post}	statistical methods	statistical procedures applied on D	appendix E, result 4.4
D	data	application of S_{pre} to sample the population assuming M_A	
D_s	data structure	structural aspects of the data such as sample size and number and type of variables	appendix F, result 4.5
D_v	data values	observed values that signify a fixed realization of the data	result 4.6
R	result	decision rule which maps the application of S_{post} to D onto the decision space such as choice of a model over others, a parameter estimate, or rejection of a null hypothesis	definition 2.3, 2.4
ϕ	true reproducibility rate	limiting frequency of reproduced results in a sequence of replication experiments	appendix A, appendix G, definition 2.5,

(Continued.)

Table 1. (Continued.)

symbol	name	description	formal definition/result
ϕ_N	estimated	sample frequency of	result 2.7,
	reproducibility	reproduced results in a sequence of N	result 5.1,
	rate	replication experiments	remark 2.8
π -Open	openness	which elements of ξ are available to ξ'	definition 2.6

mathematical object. The benefits of this isolation will become clear in §3, where an unconditional ξ and its non-exact ξ' pair may become a ξ and its *exact* ξ' pair, conditional on R .

Often, however, we would perform experiments with a specific aim and would like to see whether the result of ξ is reproduced in ξ' . Depending on the desired mode of statistical inference, example aims include hypothesis testing, point or interval estimation, model selection or prediction of an observable. Further, when augmented with an R , K' must differ from K in a specific way. Encompassing all these statistical modes of inference, we introduce the notion of a *result* R , as a decision rule. For convenience, we assume that R lives on a discrete space here.

Definition 2.3. Let \mathcal{X} be the sample space and $\mathcal{R} = \{r_1, r_2, \dots, r_q\}$, $q \in \mathbb{Z}^+$ be the decision space. For sample size $n \in \mathbb{Z}^+$, the function $R: \mathcal{X}^n \rightarrow \mathcal{R}$ is a *result*.

R is obtained by mapping the application of S_{post} on D on to the decision space. If ξ' is aimed at reproducing R of ξ , it is conditional on R and leads us to the following connection between an idealized experiment and a result.

Definition 2.4. Let R and R' be results from ξ and ξ' , respectively. $R = r_o$ is *reproduced* by $R' = r_d$ if $d = o$, else $R = r_o$ is *not reproduced*.

In definition 2.4, reproducibility of R depends on the available actions r_1, r_2, \dots, r_q . The size of q is case specific. Examples are as follows. In a null hypothesis significance test, $q = 2$: the null hypothesis and the alternative hypothesis. In a model selection problem, we entertain q models and choose one as the best model generating the data. In a parameter estimation problem for a continuous parameter, we build q arbitrary bins, and call a result reproduced if the estimate from ξ' falls in the same bin as the result from ξ . How the bins are constructed in a problem affects the actual reproducibility rate of a result. However, for our purposes in this article, theoretical results hold for all cases regardless of this tangential issue.

The class of problems of interest to us here involves cases where, in a *sequence* of exact replication experiments, if S is reliable, we should expect a regularity in the results. That is, probability theory tells us that if the elements of an idealized experiment are well defined, then we should expect the results from a sequence of replication experiments to stabilize at a certain proportion, given the characteristics of an idealized experiment and the true data-generating mechanism. This notion is formalized in definition 2.5.

Definition 2.5. Let $\xi^{(1)}, \xi^{(2)}, \dots, \xi^{(N)}$ be a sequence of idealized experiments. The *reproducibility rate*

$$\phi = \lim_{N \rightarrow \infty} N^{-1} \sum_{i=1}^N \mathbf{I}_{\{R^{(i)}=r_o\}},$$

of a result, $R = r_o$ is a parameter of the sequence ($\mathbf{I}_{\{C\}} = 1$ if C , and 0 otherwise).

An advantage of definition 2.5 is that conditional on $R = r_o$ in ξ and a sequence of replication experiments $\xi^{(1)}, \xi^{(2)}, \dots, \xi^{(N)}$, the *relative frequency* of reproduced results ϕ_N converges to $\phi \in [0, 1]$ as $N \rightarrow \infty$. So, we immediately have $\phi_N = N^{(-1)} \sum_{i=1}^N \mathbf{I}_{\{R^{(i)}=r_o\}}$ as a natural estimator of ϕ . Further, we are formally comforted to know that $\lim_{N \rightarrow \infty} \mathbb{P}(\phi_N = \phi) = 1$. That is, with high probability, the estimated reproducibility rate ϕ_N from a sequence of replication experiments will get closer to the true reproducibility rate of the original experiment ϕ .

Finally, we turn to the last of our key concepts: *openness*. Openness refers to the accessibility of all necessary information regarding the elements of ξ by another idealized experiment ξ^* . This accessibility may be used for a variety of purposes. For example, S_{post} can be re-applied to D to verify R independently of ξ . In this capacity, openness facilitates the auditing of experimental results by way

of screening off certain errors, including human and instrumental (e.g. data entry and programming errors), that may be introduced in the process of obtaining R initially. On the other hand, openness may be needed to perform an exact ξ' by way of duplicating S_{pre} to obtain D' and S_{post} to obtain R' . In this capacity, openness makes exact ξ' possible.

Openness is critically related to reproducibility since the degree to which information is transferred from ξ to ξ' impacts the ϕ of a given result. However, not all elements of ξ need to be open for all purposes. Therefore, a nuanced understanding of openness requires evaluating it at a fixed configuration of the elements of ξ conditional on a specific purpose, rather than as a categorical judgement at the level of the whole experiment, as open or not. This leads us to think of openness element-wise, as in definition 2.6.

Definition 2.6. Let Π be the power set of elements of ξ and $\pi \in \Pi$. ξ is π -Open for ξ^* if $\pi \subset K^*$, where ξ^* is an idealized experiment that imports information from ξ .

A specific example of π -Open of definition 2.6 would be $\pi \equiv (M_A, S_{\text{pre}})$, where ξ^* gets all the information about the assumptions, model and pre-data methods from ξ but no other information. Another example of π -Open is the special case where ξ has all its elements open, such that $\pi \equiv (K, M_A, S, D)$. In this case, for convenience, we say ξ is ξ -Open for ξ^* .

2.2. Fundamental results on replications and reproducibility rate from first principles

Here, we present two results about reproducibility and some remarks, based on definitions 2.1–2.6. A well-formed theory of reproducibility requires results of these types: fundamental, mathematical and invoking a functional framework to study replications and reproducibility. They serve as theoretical benchmarks to check other results against. Technically oriented readers may refer to appendices A and B for a more detailed discussion and results complementary to the main argument.

We begin by noting that, given definition 2.5 and the discussion following it, it is not straightforward to say exactly what we gain if we were to update the estimated reproducibility rate based on the results obtained from performing more replications. Indeed, to understand the value of replication experiments in assessing the reproducibility of a result, a strong mathematical statement is required, which is our result 2.7.

Result 2.7. Let $\xi^{(1)}, \xi^{(2)}, \dots, \xi^{(N)}$ be a sequence of replication experiments with reproducibility rate ϕ given by definition 2.5. Then,

$$\mathbb{P}\left(\lim_{N \rightarrow \infty} \phi_N = \phi\right) = 1, \quad (2.1)$$

where ϕ_N is the sample reproducibility rate of result $R = r_o$ obtained from the sequence (proof in appendix A).

Result 2.1 is fundamental to study replications and reproducibility for a number of reasons:

1. It provides a basis for building trust in the notion of reproducibility from replication experiments. Roughly, it says that if we perform replication experiments and estimate the reproducibility rate of r_o by ϕ_N from these experiments, then we are *guaranteed* that deviations of ϕ_N from ϕ are going to *get small* and *stay small*.
2. It is almost necessary to move forward theoretically. It immediately implies that if the assumptions of an original experiment are satisfied in its replication experiments, then we are *adopting a statistically defensible strategy* by continuing to perform replication experiments and updating ϕ_N as a proportion of successes to assess the reproducibility rate. Therefore, result 2.7 gives us a theoretical justification of *why we should care* about performing more replication experiments whose assumptions are satisfied and be interested in estimating the reproducibility rate based on those replication experiments alone. Further, violating the assumptions of ξ in replication experiments implies that ϕ_N converges to some ϕ defined by the flaws underlying a non-exact sequence of replications of ξ rather than the reproducibility rate of r_o of interest.
3. As we will detail in result 2.9, a theoretically fertile way to study replication experiments is by defining a sequence of experiments as a stochastic process. The results from such processes almost always require the solid foundation provided by result 2.7.

Remark 2.8. The reproducibility rate given in definition 2.5 has excellent properties as shown by result 2.7. However, we keep in mind that definition 2.5 is only one way to measure reproducibility. It

is a counting measure which counts the reproduced results. Instead, a continuous measure as a degree of confirmation of a result might seem more proper to measure reproducibility. One has to be aware that just defining a reproducibility measure does not imply that it has desirable mathematical properties. It is easy to define meaningful continuous measures of reproducibility which might have pathological properties (e.g. that do not satisfy result 2.7), and these should be avoided (see appendix A for details).

In practice, S_{post} are functions of sample moments, such as the sample mean. In these cases, sometimes the Lindeberg–Lévy central limit theorem (CLT) and its extensions provide useful results about the properties of $\xi^{(1)}, \xi^{(2)}, \dots$. However, restricting S_{post} this way constrains the mathematical setting to study the statistical properties of $\xi^{(1)}, \xi^{(2)}, \dots$ or results reproducibility. For example, working with the CLT is challenging when S_{post} cannot be formulated as a function of a fixed sample size or to discuss the properties of a sequence of replication experiments directly, without referring to S_{post} as a means to estimate a particular R .

We provide a broad setting without these limitations by assuming that K requires only minimal validity conditions on M_A and S . Specifically, we let M_A be any probability model, subject only to some mathematical regularity conditions such as continuity of distribution functions, the existence of the mean and the variance of the variable of interest. We also let S_{post} be the sample distribution function.⁴ With the generality provided by these assumptions, we obtain one of our main theoretical results.

Result 2.9. The sequence of idealized experiments $\xi^{(1)}, \xi^{(2)}, \dots$ given by definition 2.5 is a proper stochastic process, seen as a joint function of random sample D and of each value in the support of data-generating mechanism, $x \in \mathbb{R}$ (see constructive proof in appendix B).

Result 2.9 is of fundamental importance to study results reproducibility mathematically because it allows us to apply the well-developed theory of stochastic processes to build a theory of results reproducibility. Two aspects of result 2.9 are noteworthy:

1. When we obtain a random sample in ξ and perform inference using a fixed value of a statistic such as a threshold, the sequence $\xi^{(1)}, \xi^{(2)}, \dots$ constitutes random variables independent of each other conditional on the true model generating the data. Obtaining the distributions implied by ξ helps us understand the statistical nature of replication experiments.
2. ξ' generates new data D' , and R' is conditional on D' . That is, when inference is performed for a particular replication experiment, the data are fixed. Most generally, conditional on D' if the empirical distribution function is R' , then the replication experiment estimates the model generating the data. Therefore, a replication experiment determines a sample-based estimate of a statistical model.

In the next section, we introduce a toy example as a running case study to instantiate our theoretical results on replications, reproducibility and openness.

3. A toy example

Our toy example involves an inference problem regarding a population of ravens, K . An infinite population of ravens where each raven is either black or white constitutes the population assumptions, A . Each uniformly randomly sampled raven can be identified correctly as black or white, which defines the pre-data methods, S_{pre} . The result of interest, R , is to estimate the (unknown) population proportion of black ravens, p , or some function of it.

We consider six distinct sampling scenarios, which lead to six distinct M_A , and thus six distinct idealized experiments. To avoid overly complicated mathematical notation, we denote the models by $\xi_{\text{bin}}, \xi_{\text{negbin}}, \xi_{\text{hyper}}, \xi_{\text{poir}}, \xi_{\text{exp}}$ and ξ_{nor} . These models represent the binomial, negative binomial, hypergeometric, Poisson, exponential and normal probability distributions for the data-generating mechanism, respectively. In specific examples, we also vary S_{post} , the point estimator of the parameter of interest to take values as maximum likelihood estimate (MLE), method of moments estimate (MME) and posterior mode (i.e. Bayesian inference). We further vary D_s via the sample size (i.e. $n \in \{10, 30, 100, 200\}$). We use these idealized experiments to illustrate our results in the rest of the article.

⁴We assume that the order in which the data values appear has no bearing on the inferential goal. The cases in which the order contains information are important for a variety of subject matters, but it is well known that the statistical techniques that deal with them are too specialized to be treated in a general set-up. An example is autoregressive models.

experiments (exact)	result	model	population assumptions	model	result	experiments (approximate)
ξ_{bin} binomial experiment	estimate proportion of black ravens p	sample n ravens	infinite population of black and white ravens population proportion of black is p	sample n ravens	estimate mean number of black ravens approximated by np	ξ_{poi} Poisson experiment
ξ_{negbin} negative binomial experiment	estimate proportion of black ravens p	sample w white ravens		measure time between sampled n ravens	estimate mean time to black ravens approximated by np	ξ_{exp} exponential experiment
ξ_{hyper} hypergeometric experiment	estimate proportion of black ravens p	sample n ravens from a finite subset of the population		sample n ravens	estimate mean number of black ravens approximated by np	ξ_{nor} normal experiment

Figure 1. Six idealized experiments ξ_{bin} , ξ_{negbin} , ξ_{hyper} , ξ_{poi} , ξ_{exp} , ξ_{nor} : The binomial, negative binomial, hypergeometric, Poisson approximation to binomial, exponential waiting times between Poisson events and normal approximation to binomial, respectively. All but ξ_{hyper} assume infinite population (A) of black and white ravens, with sampling designs resulting in distinct probability models (M_A). ξ_{hyper} assumes sampling from a finite subset of the population. All experiments aim at performing inference on result (R), which reduces down to an estimate of either the population proportion of black ravens or the mean number of black ravens in the population.

These six idealized experiments make the following sampling assumptions. ξ_{bin} stops when n ravens are sampled. ξ_{negbin} stops when w white ravens are sampled. ξ_{hyper} is a special case where the sampling has access only to a finite subset of the infinite population delineated by A . ξ_{bin} , ξ_{negbin} and ξ_{hyper} are often called *exact* models, in the sense that their M_A does not involve any limiting or approximating assumptions. On the other hand, ξ_{poi} approximates ξ_{bin} , where a large sample of n ravens is sampled when the proportion of black ravens p is small. The larger the n and the smaller the p such that np remains constant, the better the approximation. ξ_{exp} has the same approximative characteristics and parameter as ξ_{poi} . However, notably, ξ_{exp} records the time between observations instead of counting the ravens, so its S_{pre} is different from all other experiments. Finally, ξ_{nor} approximates ξ_{bin} where a large sample of n ravens with intermediate proportion of black ravens, p , holds.

As the result of interest, R , these six idealized experiments aim to estimate either the proportion of black ravens, p , in the population or the rate of black ravens sampled, $np \rightarrow \lambda$, a function of p , in the approximative models. Figure 1 shows distinctive elements of these six idealized experiments.

In §4, we use these six idealized experiments to show that *openness* connects to reproducibility in a variety of ways and to *reproduce* a given result, and *replication experiments* do not need to be *exact*. We show that conditional on a given result from an original experiment, *non-exact* replication experiments can serve as valid *exact* replication experiments, if the inferential equivalence holds between the original and the replication. We further show that the true rate of reproducibility of a sequence of exact replication experiments and a sequence of non-exact replication experiments are distinct (except trivially) for a given result.

4. Element-wise openness and assessing the meaning of replications

Tools and procedures have been developed to help facilitate openness in science [11,14,17,22]. Guidelines may argue for making as much information available as possible about an experiment or leave it to intuition to guide which elements of an experiment are relevant and need to be shared for replication. We are interested in better understanding what does and does not need to be made available, in service of which objective, and under what conditions. We perceive two main issues: what openness means for performing meaningful replications and how it impacts results reproducibility. We first evaluate the former. Then we show that a uniform, wholesale framing of openness is not the remedy to the reproducibility crisis that some take it to be.

ξ has elements involving uncertainty, such as D_o taken as a random variable. Uncertainty modelled by probability is always conditional on the available background information [23], and thus,

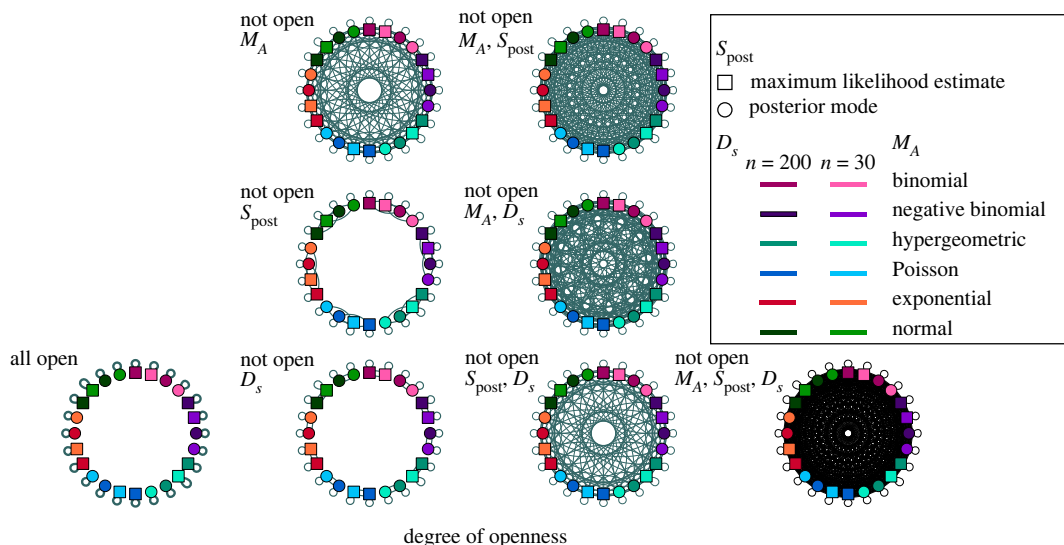


Figure 2. For the models in the toy example, degrees of openness (as given by definition 2.6) are depicted in eight networks, each consisting of the same 24 idealized experiments. Each idealized experiment is represented by a node in each network. These 24 experiments are obtained by a $6 \times 2 \times 2$ factorial design. The first factor, M_A , takes six values: binomial, negative binomial, hypergeometric, Poisson, exponential and normal. The second factor, S_{post} , takes two values: MLE and posterior mode. The third factor, D_s , takes two values: $n=30$ and $n=200$. Connections between nodes represent potential substitutions of non-open elements of idealized experiments. As more elements of an idealized experiment are non-open, the probability of choosing an exact replication decreases, as indicated by increased connectivity in the network.

reproducibility of R is always conditional on K . That is, ξ' must import sufficient information from ξ with respect to R of interest to assess whether R is reproduced in R' . A ξ' that aims to reproduce a given result from ξ may be performed in a variety of ways depending on which elements of ξ are open.

In the context of our toy example, figure 2 shows a network structure of some possible ξ as a function of which elements of ξ are open. Specifically, we consider variations of the six experiments introduced in §3 for two S_{post} (MLE and posterior mode) and two D_s ($n=30$ and $n=200$) yielding 24 distinct ξ , each denoted by a node in each network in figure 2. Given 1 of these 24 as ξ , all possible 24 experiments are either exact or non-exact ξ' . We use definition 2.6 and specify π to assess the degree of openness in these experiments. When ξ is ξ -Open, the probability of exact replication is 1, and every node of the network is only connected to itself. If ξ is π -Open, where π is a proper subset of ξ , then ξ' may be a non-exact replication of ξ in various ways because ξ' needs to substitute in a value for elements that are not in π . Therefore, the probability of ξ' being an exact replication of ξ is lower than when ξ is ξ -Open. In figure 2, we show the network structures that result from choosing non-open elements with equal probability among all substitutions considered for each element. The network complexity depends on the size of π . If it is large, the number of connections among the nodes in the network is small, and each connection is strong (e.g. strongest when all open). In contrast, if it is small, the number of connections among the nodes in the network is large because there are both multiple substitutions to be made and multiple possibilities for each, and each connection is weak (e.g. weakest when M_A , S_{post} , D_s not open in figure 2). Hence, as the size of π decreases, it becomes less probable to perform an exact replication of ξ . By looking at which elements of ξ are open to start with, we can assess how the sequence $\xi^{(1)}$, $\xi^{(2)}$, ... of replication experiments can be misinterpreted if the necessary elements were not open and/or got lost in translation. In the rest of this section, we organize our results by elements K , M_A , S , D .

4.1. Background knowledge, K

Providing an exact description of what goes into K is notoriously difficult. K , which is more of a philosophical element of ξ , typically carries over much more than what can be immediately gleaned over by a transparent and complete description of M_A , S and D . We understand K to contain theoretical assumptions, contextual knowledge, paradigmatic principles, a specific language and presuppositions inherent in a given field; in short, a lot of inherited cultural and historical meaning of

the kind Feyerabend refers to as *natural interpretations* in the *Against Method* ([24], p. 49). As Feyerabend explains, such natural interpretations are not easy to make explicit or even sometimes be aware of and thus, being open about them might not be a matter of choice. However, observations gain meaning only against this backdrop, and experiments can only be interpreted correctly by using the same language used to design them in the first place. Within ξ , this tends to happen implicitly, whereas when performing ξ' , there is no guarantee that all the relevant information in K will carry over to K' .

By using the binomial experiment in our toy example, we can illustrate why K is an integral part of ξ and what role it plays for ξ' . In ξ_{bin} our aim (R) is to estimate the proportion of black ravens (p) in an infinite population of ravens (A). M_A samples n ravens. In the case of our S_{pre} , we count black and white ravens by naked eye. In the case of our S_{post} , we use the maximum likelihood estimator of p . We set $n = 100$, which constitutes our D_s . This description of ξ_{bin} based on a specific configuration of M_A , S_{pre} , S_{post} , D_s could just as well be used to define an experiment in which scientists are interested in estimating the proportion of black *swans* in a population of black and white swans. While ξ_{bin} would still be mathematically well defined, its scientific content and context are not captured by any of these four elements. For that, we need K . Without K , we would have to consider an ξ'_{bin} about black swans as an acceptable replication of ξ_{bin} about black ravens, based on the mathematical structure alone. K , then, communicates scientific meaning across experiments.

As a more practical example of the import of K , we consider a recent ‘failed’ replication experiment. Murre [25] attempted to replicate a classical experiment by Godden and Baddeley [26] on context-dependent memory. Context-dependent memory refers to the hypothesis that the higher the match between the context in which a memory is being retrieved and the context in which the memory was originally encoded, the more successful the recall is expected to be. In the abstract, Murre [25] summarizes the results of the replication experiment as follows: ‘Contrary to the original experiment, we did not find that recall in the same context where the words had been learned was better than recall in the other context.’ Does this suggest that the results of the original experiment were a false positive—as replication failures are commonly interpreted? There are many reasons to not jump to that conclusion including sampling error and the fact that the context of the replication was different from that of the original [26] experiment. Specifically, unlike the original, the replication was being filmed as part of a TV programme. We will set these obvious concerns aside for a moment to focus on another. Ira Hyman explains the issue in a Twitter thread [27]. Hyman indicates that the phenomenon of context-dependent memory is conditional on the distinctiveness of the encoding context. That is, if distinct contexts are used over multiple trials, the chances that the context will be remembered with the encoded information increases. When the context is not distinctive enough or remains constant over trials, the effect disappears. Another known boundary condition for the phenomenon is the outcome variable: past research has shown that this works for retrieval tasks (e.g. free recall) and not recognition. The Murre [25] replication did not carry over these contextual details and changed the design in a way to not instigate context-dependent memory. As a result, the differences between R and R' become impossible to attribute to a single cause and fail to provide evidence that can confirm or refute the results of the original Godden and Baddeley [26] experiment. It is even questionable whether the Murre [25] experiment provided an appropriate test of the result of interest in the first place to be considered a meaningful or relevant replication.

This replication example on context-dependent memory appears to imply that a ξ' is meaningful or relevant with respect to a specific result R . By definition 2.2 and its interpretation, however, we know that mathematically, it is more convenient to separate the definition of ξ' from R . It follows that there are at least two aspects of assessing the meaning and relevance of a replication.

Firstly, while an operational definition of K is elusive, a useful way to think about K is ‘all the information in ξ that is not already in M_A , S and D' ’. At the minimum, for ξ' to be considered a *meaningful* replication of ξ , K' must import some information in K regarding the immediate scientific context of ξ . For this to hold, there is no need to invoke the notion of R .

Second, to assess the reproducibility of a given R , K' must import *relevant* information pertaining to R from ξ . That is, replication experiments unconditional and conditional on R are not the same objects. To emphasize the difference between them, we distinguish between *in-principle* and *epistemic* reproducibility of an R in remark 4.1 (for further details, see appendix C).

Remark 4.1. Let ξ be an idealized experiment and ξ' be its exact replication. Conditional on R from ξ , K' is necessarily distinct from K for epistemic reproducibility of R by R' , but not necessarily distinct for in-principle reproducibility of R by R' .

In practice, ξ' can never be an *exact* replication of ξ in an ontological sense. The ξ is a one-time event that has already happened under certain conditions and ξ' has to differ from ξ in some aspect. The best standard that ξ' can purport to achieve is to capture relevant elements of ξ in such a way that performing inference about R while adhering to A and sampling the same population is possible within an acceptable margin of error. However, every experiment is embedded in its immediate social, historical and scientific context, making it a non-trivial task for scientists to include all the relevant K when they report the experiment in an article and make explicit all the natural interpretations used to assign meaning to its results. As such, designing and conducting replication experiments cannot be reduced to a clerical implementation of reported experimental procedures. A comprehensive understanding of K is increasingly critical as ξ' diverges further away from ξ to be able to comprehend the nature and importance of the divergence for the interpretability of ξ' and for results reproducibility. For ξ' to serve their intended objective, information readily available from ξ needs to be supplemented by a careful historical and contextual examination of the relevant literature and the broader scientific background. Otherwise, ξ' may differ from ξ in non-trivial ways impacting the meaning of the evidence obtained and changing the estimated reproducibility rate.

4.2. Model, M_A

For ξ' to be able to reproduce *all possible* R of ξ , M_A must be specified up to the unknown quantities on which inference is desired. This specification must be transmitted to ξ' , such that M_A and M'_A are identical for inferential purposes mapping to R . If an aspect of M_A that has an inferential value mapping to R is not transmitted to ξ' and this inferential value is lost, then R cannot be meaningfully reproduced by R' . On the other hand, given an inferential objective mapping to a specific R , the aspects of M_A that are irrelevant to that inferential objective need not be transmitted to ξ' to meaningfully reproduce R by R' . Counterintuitively, to meaningfully reproduce R by R' , M_A and M'_A do not need to be identical, as given by result 4.2.

Result 4.2. M_A and M'_A do not have to be identical in order to reproduce a result R by R' . Under mild assumptions, the requirement for R to be reproducible by R' is that there exists a one-to-one transformation between M_A and M'_A for inferential purposes mapping to R (proof and details in appendix D).

As an example of result 4.2, consider ξ_{bin} and ξ_{negbin} in figure 1. Conditional on the objective of estimating p , the population proportion of black ravens, any of $(\xi_{\text{bin}}, \xi_{\text{bin}})$, $(\xi_{\text{bin}}, \xi_{\text{negbin}})$, $(\xi_{\text{negbin}}, \xi_{\text{bin}})$ and $(\xi_{\text{negbin}}, \xi_{\text{negbin}})$ can be effectively considered a pair (ξ, ξ') of an idealized experiment and its (*exact*) replication. The reason is that the quantity of interest p is an identifiable parameter in both experiments, although M_A and M'_A are not necessarily identical.⁵

In practice, when conducting a sequence of replication experiments, we would be interested in gauging the extent to which we can reproduce a specific result. Assuming that S are the same throughout all experiments, we expect the observed reproducibility rate of a sequence of experiments whose elements are chosen from $\xi_{\text{bin}}, \xi_{\text{negbin}}$ to converge on the same value, capturing the information on p , in the same way. However, result 4.2 does not imply that the (true) reproducibility rate of any two sequences of experiments involving any M_A and M'_A are equal to each other. In fact, the (true) reproducibility rates of two sequences are not equal, when non-exact replications are involved.

Openness of M_A to M'_A needs to be distinguished from the equivalence of M_A and M'_A . In ξ_{bin} and ξ_{negbin} , M'_A is not equivalent to M_A . However, the binomial and the negative binomial models become equivalent with respect to a certain inferential objective that allows for reproducing a specific R , which is estimating p . To establish this compatibility, M_A should be open to ξ' but does not need to be assumed in ξ' . Specifically, to set M'_A to be the negative binomial model in ξ' to reproduce the estimate of p in ξ , we need to know that ξ has used the binomial model. This ensures that ξ' can use a model that has the same parameter p with the exact same meaning as in ξ and same population assumptions A such that the inferential equivalence holds. A model that has different population assumptions A from ξ_{bin} and ξ_{negbin} is ξ_{hyper} . This difference matters for reproducing a specific R . ξ_{hyper} samples from an arbitrary finite subset of infinite population but still uses the same parameter p as ξ_{bin} and ξ_{negbin} . The estimate of p in ξ_{hyper} will be biased due to differences in A . Without access to full specification of M_A , this compatibility between M_A and M'_A or lack thereof cannot be established.

⁵Compare this statement to definition 2.2 of an exact and non-exact replication experiment unconditional on an inferential objective.

This point is illustrated in many-analyst studies [28,29] in which a fixed D is independently analysed by multiple research teams who are provided D and a research question that puts a restriction on which R would be relevant for the purposes of the project. The teams were not, however, provided a M_A , S_{post} or full specification of K . Teams used a variety of models differing in their assumptions about the error variance and the number of covariates (M_A) to analyse D . The results differed widely with regard to reported effect sizes and hypothesis tests. So even when D was open, the lack of specification with regard to M_A yielded largely inconsistent results. It is not because the same aspects of reality cannot be captured by different models but because researchers did not automatically agree on which aspects to capture in their models.

Taking stock, our ravens example is deliberately simple to help in our analysis. State-of-the-art models are often large objects. If M_A is large, it might not always be clear which class of models M'_A can be drawn from to be equivalent to M_A , and finding this class might be unfeasible. Then M_A needs to be both open to and photocopied by ξ to be able to reproduce the results of interest. This point is particularly important to communicate to scientists who primarily engage in routine null hypothesis significant testing procedures and may not be conventionally expected to transparently report their models.⁶

4.3. Method, S

4.3.1. Pre-data methods, S_{pre}

S_{pre} comprises a wide range of procedural components in ξ that feeds into collection of D_o . Examples of S_{pre} are determining types of observables, unobservables and constants; measurement and instrumentation choices; and sampling procedures such as random number generators used in computational methods.

Pertaining to mathematical features of the variables of interest, S_{pre} may capture their types or a particular scaling. For example, a variable can be assumed discrete, continuous, or both discrete and continuous for mathematical convenience. This choice determines whether we are bound by a counting measure or a Lebesgue measure. A variable can also be assumed categorical, ordinal, interval or ratio. Some variables or parameters are scaled to the interval $[0, 1]$ on the real line, to make their interpretation natural. All of these S_{pre} choices affect M_A and the consequent S_{post} .

Pertaining to operational features of the variables of interest, S_{pre} may capture the method of observation and measurement instruments. In our toy example, a raven can be observed for its colour by naked eye (S_{pre}), but another investigator may opt for a mechanical pigment test (S'_{pre}). What considerations should be given when making substitutions for S_{pre} ? One issue due to choices in operationalization is measurement error. Measurement error in observables, when not accounted for, might be a factor unduly exacerbating irreproducibility or inflating reproducibility [10,32,33]. Another issue arises due to arbitrary choice of experimental manipulations or conditions which might not be mathematically equivalent. For example, manipulations that are not tested for specificity may end up manipulating non-focal constructs or only weakly manipulate the focal construct (i.e. leading to small effect sizes) [34].

Even though knowing all these features is useful in understanding S_{pre} , there is a caveat. All aspects of S_{pre} must be fixed before realizing D_o , and it is challenging to assess *a priori* whether ξ and ξ' using different S_{pre} and S'_{pre} , respectively, can be equivalent to each other. Due to these complexities and ambiguities surrounding S_{pre} , openness of S_{pre} seems to be the easiest way to obtain an equivalent S'_{pre} in designing and performing ξ' . However, there are well-known examples to show that S_{pre} and S'_{pre} can be different and yet ξ and ξ' can be equivalent conditional on R , which leads us to result 4.3.

Result 4.3. S_{pre} and S'_{pre} do not have to be identical to reproduce a result R .

As an example of result 4.3, consider models ξ_{poi} and ξ_{exp} in figure 1. ξ_{poi} has a good approximative model to the model in ξ_{bin} if we think of sampling ravens continuously from a population where black ravens are rare. We assume $np \rightarrow \lambda$, where λ is the rate of sampling the black ravens (parameter of the Poisson model), and under this assumption, we focus on inference on λ . Now, as a thought experiment, let us assume that we do not have a device to count the number of black ravens past 1.

⁶Cooper and Guest [30] and Guest and Martin [31] make a similar point for computational reproducibility. They highlight the importance of making models available, and particularly clearly reporting model specifications and implementation assumptions so as to facilitate replication.

However, we have a chronometer. As a result of using the model in ξ_{poi} , we are, as a mathematical fact, also using the model ξ_{exp} , which measures the *time* between observing black ravens. Further, the two models have the same parameter, with the same interpretation. Therefore, if we were to measure the time between observing black ravens for a sample, then we can still perform inference on the rate of observing black ravens from the population. We note that ξ_{bin} , ξ_{negbin} , ξ_{hyper} , ξ_{poi} and ξ_{nor} operate under different assumptions, but are still *counting* ravens and interested in the number of black ravens. In contrast, ξ_{exp} is considerably different from these experiments. It is *not* counting ravens, but *measuring time*, which we would reasonably define as a continuous variable. While S_{pre} in ξ_{exp} differs considerably from all other experiments in our toy example, the exponential experiment would serve as a meaningful ξ' to reproduce R in any of them, at least approximately.

4.3.2. Statistical methods, S_{post}

Statistical methods, S_{post} , that are designed for a specific inferential goal, R , but do not return identical values when applied to a fixed D are common. Conversely, some statistical methods return identical values for a specific inferential goal, R , and they are mathematically equivalent conditional on D , even though they operate under distinct motivating principles. We have the following result.

Result 4.4. S_{post} and S'_{post} do not have to be identical to reproduce a result R by R' .

For the experiments ξ_{bin} and ξ_{negbin} in our toy example, the MLE and the MME of p are numerically equivalent (appendix E). This equivalence holds even when the interpretation of probability differs between methods. For example, MLE and the posterior mode in Bayesian inference under uniform prior distribution on parameters are equivalent regardless of all else.

At the minimum, for ξ' to be a meaningful replication of ξ conditional on R , the modes of inference should be equivalent. That is, the pair $(S_{\text{post}}, S'_{\text{post}})$ should belong to one of: point estimators, interval estimators, hypothesis tests, predictions, or model selection. Further, S_{post} should be open to ξ' , but it does not need to be duplicated to establish equivalence. For example, to use MME to estimate p in ξ' , we need to know that ξ has used MLE or MME. This way, we can ensure that ξ' will at least use a numerically equivalent estimator as the one used in ξ , even if not equivalent in principle. On the other hand, it is well known that a variety of S_{post} for the same mode of inference may yield different R . The many-analyst project by Silberzahn *et al.* [29] provides clear examples of this. Teams that were given a fixed D to analyse for a predetermined R (i.e. effect size as given by odds ratio) ended up implementing their choice of S_{post} . Even when their modelling assumptions matched, the results they reported varied. For instance, out of the teams that assumed a logistic regression model with two covariates, one pursued a generalized linear mixed-effects model with a logit link for S_{post} ([29], line 15 in table 3) and another pursued a Bayesian logistic regression ([29], line 16 in table 3). The confidence intervals around the effect size estimates reported by these two teams do not even overlap despite using a fixed D .

4.4. Data, D

4.4.1. Data structure, D_s

In statistics and philosophy of statistics, D' is often seen as the *new data* of the old kind in the sense that D_v and D'_v are independent of each other, but D_s and D'_s are identical. However, conditional on R , we have result 4.5.

Result 4.5. D_s and D'_s do not have to be identical in order to reproduce a result R by R' .

As an example of result 4.5, we consider the models in ξ_{poi} and ξ_{exp} in figure 1. Poisson model *counts* the black ravens as observable. It assumes that black ravens are observed with a constant rate. Exponential model measures the *time* between arrivals of black ravens. It also assumes that black ravens are observed with a constant rate. By referring to the unit of observations, we see that the data structures in ξ_{poi} and ξ_{exp} are distinct. And yet, the unknown parameter about which inference is desired is the same, λ —the rate of black ravens appearing in continuous sampling (appendix F).

As another example, note that the stopping rules of ξ_{bin} and ξ_{negbin} are different from each other. The stopping rule affects D_s because the maximum number of black ravens in ξ_{bin} is n , but in ξ_{negbin} it is the maximum number of black ravens in the population. And yet, the estimate of p is the same in both experiments.

Data sharing is sometimes viewed as a prerequisite for a reproducible science [8,13,35,36]. Our analysis suggests that this statement requires further qualification and calls for attention to D_s . Result 4.5 notwithstanding, changes in D_s are not trivial and they impact the true reproducibility rate. For example, ξ' might be designed to have a larger sample size than that of ξ . In this case, the variance of the sampling distribution of the sample mean decreases linearly with the sample size, and hence, it would be different for ξ and ξ' . Typically, larger sample sizes are pursued to increase the statistical power of a hypothesis test in ξ' . While such ξ' will indeed increase the power of a test, it also impacts the reproducibility rate. Counterintuitively, under some scenarios, this might play out as reproducing false results with increased frequency (see [10], for such counterintuitive results).

4.4.2. Data values, D_v

Having open access to D_v has no bearing on designing and performing a meaningful ξ' or on the reproducibility of R . Conditional on R , ξ' aims to reproduce R , not D_v . Therefore, reporting R from ξ is sufficient for ξ' to assess whether R is reproduced by R' . However, information from ξ can be reported in a variety of ways and does not necessarily contain R . We show this with an example. We consider a model selection problem with three models M_1, M_2, M_3 , where ξ and ξ' use some information criterion (IC) as S_{post} . Assume ξ reports selecting M_1 as R . This is all ξ' needs to import to know whether R is reproduced in R' . If R' reports M_1 as the selected model, then it is reproduced, else it is not. However, if which model is selected is not reported as R , ξ' needs values of IC from ξ for all M_1, M_2, M_3 , so that ξ' can redo the analysis of ξ to find out what R was. In the unlikely event that not even ICs are reported, ξ' would need D_v to re-perform the whole analysis of ξ by applying S_{post} to D to calculate ICs and then to obtain R .

Result 4.6. ξ does not have to be D_v -Open in order for ξ' to reproduce a result R .

That said, openness of D_v might facilitate auditing of R and vetting it for errors. There may be other benefits to open D_v such as enabling further research on D_v (e.g. meta-analyses). The distinction we draw matters particularly when there may be valid ethical concerns regarding data sharing [37]. Open D_v is best evaluated on its own merits as has been discussed extensively [38] but cannot be meaningfully appraised as a facilitator of replication experiments or precursor of results reproducibility. While some level of open scientific practices is necessary to obtain reproducible results, open data are not a prerequisite.

5. Exact versus non-exact replications: a simulation study on reproducibility rate

So far we have established that to reproduce R , all elements of ξ do not need to be open, and not all elements that are required to be open need to be duplicated for a meaningful ξ' . On the flip side, we also established that relatively simple openness considerations such as experimental procedures, hypotheses, analyses and data will not suffice to make ξ' meaningful. The challenge in making π -openness useful for replication experiments is to clearly identify and delineate the elements of the idealized experiment. For example, proper K is difficult to define and communicate with precision. Also, M_A is at times conflated with S_{post} and left opaque in reporting. As we discussed earlier, making K explicit and clearly specifying M_A up to its unknowns is critical when designing ξ' .

Hitherto, we focused on replication experiments and only alluded to results reproducibility when needed. In this tack, we have mathematically isolated ξ from R and made some statements about ξ unconditional, and then conditional on R to emphasize their difference. Now that we turn our attention to explicitly drawing the link from replications to reproducibility, we condition R on ξ .

Given a sequence of *exact* replication experiments $\xi^{(1)}, \xi^{(2)}, \dots$ and a result R from an original experiment ξ , do we expect to confirm R with high probability irrespective of the elements of ξ ? The answer is 'no' as shown elsewhere [10,21]. The true reproducibility rate of a result is a function of not only the true model generating the data but also the elements of the idealized experiment. ξ may be characterized by a misspecified M_A (e.g. omitted variables, incorrect formulation between variables and parameters), unreliable S_{pre} (e.g. measurement error, confounded designs, non-probability samples), unreliable S_{post} (e.g. inconsistent estimators, violated statistical assumptions), errors in D (e.g. recording errors), or large noise-to-signal ratio (e.g. large error variance and small expected value). All of these lead to the mathematical conclusion that the true reproducibility rate ϕ is specific

to each configuration of ξ and thus can take any value on $[0, 1]$. Therefore, ϕ tells us more about the experiment itself than some unobserved reality that is presumed to exist beyond it. Since we are now conditioning on ξ and questioning the reproducibility rate of R , the conclusion is that while a degree of openness may be able to address a ‘replication’ crisis by facilitating faithful replication experiments, it does not suffice to solve any alleged ‘reproducibility’ crisis.

Openness of elements of ξ facilitates ξ' , thereby allowing us to estimate ϕ of R by ϕ_N conditional on ξ . However, ϕ cannot be reasonably used as a target of scientific practices where each ξ is designed to maximize it. It does not make sense to think that a ξ that returns the highest reproducibility rate for a given R is scientifically most relevant or most rigorous experiment. For example, choosing an S_{post} that always returns the same fixed value regardless of D_o would yield $\phi = 1$. In fact, ϕ can be made independent of what it would be under sampling error.⁷

A reasonable expectation from ξ' is to deliver a scientifically relevant estimate of ϕ , given R . Openness plays an important role in this regard. In §4, we established that any non-open elements of ξ would need to be substituted for in ξ' , leading to a non-exact replication. The following result states how a sequence of non-exact replications alter the reproducibility rate.

Result 5.1. Assume a sequence $\xi^{(1)}, \xi^{(2)}, \dots, \xi^{(J)}$ of idealized experiments in which a result R is of interest. Then, the estimated reproducibility rate of R in this sequence converges to the mean reproducibility rate of R in J replication experiments. (See appendix G for proof.)

Result 5.1 states that the true reproducibility rate to which the estimated reproducibility rate of a sequence of non-exact replication experiments converges is the mean reproducibility rate of results from all experiments in the non-exact sequence and not the true reproducibility rate of a fixed original result. Hence, the reproducibility rate is a function of all elements of the idealized experiment, for both a fixed original experiment and all its replications. Each replication that is non-exact in a different way from others introduces variability, decreasing the precision of estimates given a fixed number of replications.

We illustrate the link between replication experiments and reproducibility rate with a simulation study. We consider a series of exact and non-exact replication experiments to analyse the variation in the reproducibility rate of a result as a function of the elements of ξ . We use sequences of two idealized experiments ξ_{poi} and ξ_{nor} which are approximate models to binomial from our toy example. For all conditions, we fix the true proportion of black ravens and the number of trials in the exact binomial model at 0.01 and 1000, respectively. These arbitrary choices make the true reproducibility rate distinct under ξ_{poi} and ξ_{nor} . As R , we choose a point estimate for the location parameter of the probability model. For convenience, we assume that the parameter estimates of the original experiments are equal to the true value. After each replication experiment, we determine whether this result is reproduced by R' based on whether it falls within some suitably scaled population standard deviation units of the true parameter value.

In exact replications, we vary M_A , S_{post} , D_s of the idealized experiment, each element taking two values. This results in a $2 (M_A) \times 2 (S_{\text{post}}) \times 2 (D_s)$ study design (eight conditions) for exact replications where (i) model assumed, $M_A \in \{\xi_{\text{poi}}, \xi_{\text{nor}}\}$, (ii) method as point estimate, $S_{\text{post}} \in \{\text{MLE}, \text{posterior model}\}$, and (iii) sample size, $D_s \in \{30, 200\}$. When S_{post} is the posterior model, we use conjugate priors: Gamma distribution with rate and shape parameters 5 (arbitrarily chosen) for ξ_{poi} , and normal distribution with prior mean 10 and prior precision 1 for ξ_{nor} . Figure 3a,b shows 100 independent runs of a sequence of 1000 exact replication experiments under these conditions, for ξ_{poi} and ξ_{nor} respectively.

In non-exact replications, we vary the set from which the replication experiment is uniformly randomly chosen from in each step. This results in additional three conditions: a set of all eight idealized experiments, a set of four idealized experiments with lowest reproducibility rates and a set of four idealized experiments with highest reproducibility rates. Figure 3c shows 100 independent runs of a sequence of 1000 non-exact replication experiments under these conditions.

We emphasize that all parameters of the simulation example in figure 3 are chosen so that the implications of differences between different models, methods and data structures make the link between replications and reproducibility explicit. It is certainly possible to choose these parameters to obtain any true reproducibility rate defined by a specific ξ since $\phi \in [0, 1]$.

Conditional on R , some conclusions from figure 3 are as follows.

⁷See [10] for examples of $\phi = 1$ under uncertainty.

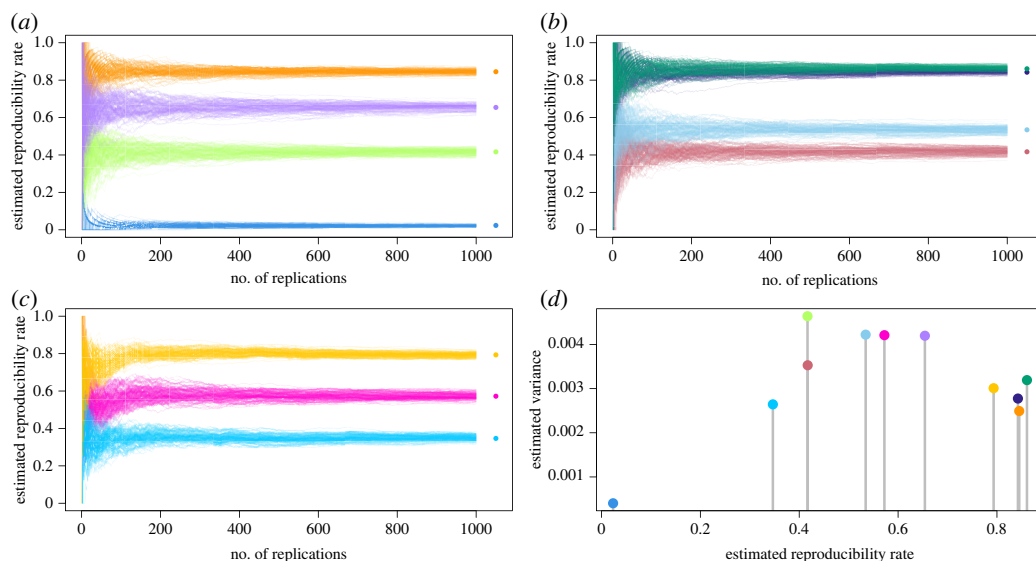


Figure 3. Reproducibility rates of a true result in sequences of 1000 exact (*a,b*) and non-exact (*c*) replication experiments. S_{post} is varied as MLE and posterior mode. D_s is varied as $n = 30$ and $n = 200$. Each condition is colour coded and consists of 100 independent runs. (*a*) M_A : Poisson. Orange; MLE, $n = 200$. Purple; posterior mode, $n = 200$. Light green; MLE, $n = 30$. Light blue; posterior mode, $n = 30$. (*b*) M_A : Normal. Dark green; posterior mode, $n = 200$. Dark blue; MLE, $n = 200$. Pale blue; posterior mode, $n = 30$. Rose; MLE, $n = 30$. (*c*) Three cases of 1000 non-exact replication experiments where they are chosen uniformly randomly from the set of all eight idealized experiments (magenta), four idealized experiments with lowest reproducibility rates (aqua blue), and four idealized experiments with highest reproducibility rates (yellow). (*a–c*) Asterisks denote the mean of the reproducibility rates of 100 runs at step 1000, an estimate of the true reproducibility rate for the sequence of idealized experiments. (*d*) Variances of all 11 exact and non-exact sequences at step 50 of the simulation with respect to the estimated reproducibility rate (see text for interpretation).

1. The true reproducibility rate depends on the true data-generating mechanism and the elements of the original experiment. Specifically, the true reproducibility rate in our simulation is a function of the true model generating the data, M_A , and also D_s such as the sample size, and S_{post} such as the method of point estimation. This can be seen from exact replication sequences of eight idealized experiments in figure 3*a,b*, with the true reproducibility rate for each experiment indicated by stars.
2. By weak law of large numbers, even if the true reproducibility rate is high (e.g. orange in figure 3*a* and green in figure 3*b*), the estimated reproducibility rate from a short sequence of exact replications has higher variance than the variance of the estimated reproducibility rate in a longer sequence. However, the estimated reproducibility rate from exact replications ultimately converges to the true reproducibility rate of an original result from a fixed ξ illustrating result 2.7.
3. Estimated rate of reproducibility from a sequence of non-exact replications may be drastically different from the true reproducibility rate of an original result. The sequence of idealized experiments shown in pink in figure 3*c* is a sequence of non-exact replications for any of the eight original idealized experiments in figure 3*a,b*. For example, assume the original experiment we aim to replicate is ξ_{poi} with S_{post} and D_s set to posterior mode and sample size $n = 30$, respectively. Blue sequences in figure 3*a* show that the true reproducibility rate of R (i.e. the estimate of location parameter) from these sequences of exact replication experiments is close to zero as shown by the convergence of 100 runs (i.e. blue star). If S_{post} and D_s were not open in this experiment, then we would have had to substitute for them and the pink sequences in figure 3*c* would serve as plausible replication experiments. In this case, we would estimate the reproducibility rate of R as approximately 60% (i.e. pink star).
4. In a sequence of replication experiments, the set we choose the experiments from matters for true reproducibility rate. An original idealized experiment and its non-exact replications belonging to a set of idealized experiments that have true reproducibility rates close to each other for a given R yield an estimated reproducibility rate that is closer to the true value of the original experiment. For example, the yellow and blue sequences in figure 3*c* come from a set of four idealized experiments with the lowest and highest reproducibility rates among all eight experiments, respectively. The set of experiments sampled in the blue sequence is compared with the set of

experiments sampled in the yellow sequence. The latter serves as a more relevant set of idealized experiments for replications of the orange and purple experiments in figure 3a, and the dark blue and green experiments in figure 3b, yielding a better reproducibility rate estimate for the original R . This pattern is an illustration of the broader theoretical result 5.1.

In practice, however, we do not have access to the true reproducibility rate of any idealized experiment to help determine our replication sets. We have to make our decision based on the elements of the idealized experiment instead, and that requires a thorough understanding of how each element of the idealized experiment impacts the reproducibility rate in a given situation.

5. The variance of the estimated reproducibility rate of results in a sequence of non-exact replications can be higher or lower than the variance of the estimated reproducibility rate in a sequence of exact replications of the original experiment. The pattern of variances we observe in figure 3d is a direct consequence of $n\phi$ following a binomial distribution and result 5.1. As a mathematical fact of the binomial distribution, its variance is maximum at $\phi=0.5$ and decreases as the probability of success, ϕ , gets closer to 0 or 1. Hence, we expect our estimates to vary greatly in a sequence of non-exact replication experiments with moderate true reproducibility rates. If a sequence of non-exact replications come from a homogeneous set of very high (or very low) true reproducibility rates, we expect our estimates to vary little. On the other hand, we expect highest variation in our estimates from exact replications if $\phi=0.5$ from the original experiment and from non-exact replications if they are highly heterogeneous in their true reproducibility rates.

In sum, the mere choice of the elements of ξ impacts both the level of the true reproducibility rate and the variance of the estimated reproducibility rate. Any divergence in ξ' may move the estimated reproducibility rate away from the true value for an original result and increase the variance of its estimates. In appendix H, we provide a broader example for result 5.1 in the context of linear regression models, under a model selection (rather than parameter estimation) scenario, where both true and false original results are considered. This simulation study demonstrates a similar pattern of results to those presented in figure 3. Combined, simulation results confirm that reproducibility rate can take any value on $[0, 1]$ depending on the elements of ξ even when the original experiment indeed captures a true result, there is no scientific malpractice, and meaningful replication experiments can be performed to reproduce R .

6. Discussion

In this article, we focused on scientific experiment as the critical unit of analysis, formalizing the logical structure of experiments towards building a theory of reproducibility. We clarified what makes for a *meaningful* replication experiment even when an exact replication experiment is not possible and established how openness of different elements of the idealized experiment contribute to it. We distinguished between the ability of a replication experiment to reproduce a result and the true reproducibility rate for that result. We showed that theoretically it is not possible to justify a *desired level* of reproducibility rate in a given line of research and to reach a high level of reproducibility rate via eliminating malpractice, requiring open procedures or data, or performing replication experiments.

We understand the potential lack of enthusiasm of the practitioner when they may find that the theory we develop does not have immediate application on their scientific practice. Our work is theoretical and is meant to present a framework to understand and study the objects and products of science. It is not meant to provide solutions to immediate problems scientists face. Practitioners often turn to theory for a clear answer to their difficulties in real-life studies. Our goal is not to provide these answers. We lay the groundwork that would potentially be needed to address such problems in the indefinite future, but the theoretical work is slow and incremental. In our simulations, we can create perfect conditions to illustrate our theoretical results because we set our own model parameters, and we know what our models are and what they mean, perfect transparency exists and there are no misunderstandings because every aspect of scientific objects is precisely known. All mathematical and statistical problems—excepting paper-and-pencil exact solutions—are primarily studied this way. Science as practised, on the other hand, is messy, ambiguous, loose, hard to define and communicate. There is no easy or direct translation of our work to myriad imperfections of the scientific practice. Our idealizations are removed from reality to make theoretical work possible in the first place. We are only laying the building blocks of such a theory to make practical implementations possible in the future. All this does not mean theory is currently of no practical relevance, however. In fact, we think

that without a thorough understanding of mathematical implications of reproducibility and replications, we cannot be ready to interpret the results and solve problems that arise in practice.

With this constraint in mind, we discuss key theoretical insights from our findings.

6.1. Reproducibility and the search for truth

A layperson understanding of reproducibility to the effect that ‘if we observe a natural phenomenon, we should be able to reproduce it and if we cannot reproduce it, our initial observation must have been a fluke’ is exceedingly misleading. A statistical fact is that reproducibility is not simply a function of ‘truth’. This was illustrated in [21] and proved in [10]: true results are not perfectly reproducible, and perfectly reproducible results are not always true (see appendix I for proof). *True reproducibility rate* of a result and the variability in its estimator are determined by many factors including but not limited to the true data-generating mechanism: The degree of rigour of the original experiment as assessed by the extent to which its elements are individually reliable and internally compatible with each other, the degree to which replication experiments are faithful to the original and how any discrepancies impact the results, the degree of rigour of the replication experiment wherever it diverges from the original and how we determine for a result to be reproduced. Factors such as effect size, sampling error, missing background knowledge and model misspecification [39,40] could render true results difficult to reproduce.

As a useful reminder, sampling error might be masked by the choice of method and other elements of the idealized experiment. A false result could be 100% reproducible due to the choice of estimation method. Therefore, judgements of reproducibility cannot exclusively be used to make valid inference on the truth value of a given result (see also [41], for a computational model with a similar conclusion).

Even if some form of a perfect experiment that captures ground truth and its exact replications exist, it might take many epistemic iterations of theoretical, methodological and empirical research to achieve them (see [42], p. 45, for a detailed discussion on epistemic iteration). We cannot expect to skip the arduous iterative process of doing science and hope to arrive at a non-trivially reproducible science with procedural interventions. In most fields and stages of science, focusing on maximizing reproducibility seems like a fool’s errand. For meaningful scientific progress, at the minimum, we should take care to properly analyse the elements of the original experiment to assess how they might impact the true reproducibility rate and analyse the discrepancies of replication experiment(s) from the original to gauge how our reproducibility estimates may vary from the true value of the original result’s reproducibility. In the course of ‘normal science’ (borrowing terminology from [43]), reproducibility of a result is more likely to tell us something about the experiments that generated the result and its reproducibility rate estimates than the lawlikeness of some underlying phenomenon.

6.2. Defining reproducibility

One aspect of reproducibility that often gets overlooked: how we define and quantify a result and its reproducibility also determines the true reproducibility rate (see also [44] for a discussion of different statistical methods to assess reproducibility and their limitations). For example, in a null hypothesis significance test, we might call a ‘reject’ decision in a replication experiment a successfully reproduced result if the original experiment rejected the hypothesis. On the other hand, we might instead look at whether effect size estimate of the replication experiment falls within some fixed error around the point estimate from the original experiment. Everything else being equal, the true reproducibility rates are expected to be different between these two cases using different reproducibility criteria.

Our findings hold under mathematical definitions of a result (definition 2.3) and of reproducibility rate (definition 2.5). In the absence of such theoretical precision, we often resort to heuristic, common sense interpretations of terms. In appendix A, we present a detailed argument on why and how theoretical precision matters and provide an example of a plausible measure of reproducibility without desirable statistical properties. Such lax standards in definitions invite unwanted or strategic abuse of ambiguities when interpreting replication results when we have a limited understanding of what we should expect to observe. Our surprise at ‘failed’ replication results or delight in ‘successful’ ones may not be warranted, and what we observe could simply be a theoretical limitation imposed by our definitions rather than a reflection of the true signal that presumably exists in nature. For an extreme example, consider the following: we might call a result as reproduced if the replication effect size estimate falls on the real line. That would trivially give us a 100% reproducibility rate.

Whenever we evaluate replications and estimate reproducibility, it is incumbent on us to understand how we define our results, how we determine reproducibility and how our measures should be expected to behave under specific conditions.

6.3. Reproducibility and openness

Open practices in science have been intuitively proposed as a key to solving the issues surrounding reproducibility of scientific results. However, a formal framework to validate this intuition has been missing and is needed for a clear discussion of reproducibility. The notion of idealized experiment serves as a theoretical foundation for this purpose. By using this foundation, we have distinguished the concepts of replication and reproducibility, showing how openness is related to meaningful replications. We have also distinguished between two types of reproducibility (appendix C). Whether elements from one experiment carry over to a replication experiment is only relevant to epistemic—as opposed to in-principle—reproducibility. In practice, however, resource constraints determine the availability and transferability of information between experiments. A realistic framework needs to provide a refined sense of which elements of an experiment need to be open to reproduce a given result, as opposed to simply saying ‘all of it’.

We have identified different levels and layers of openness and examined their implications. An experiment that is completely open in all elements does not necessarily lead to reproducible results and an experiment that does not open its data does not necessarily hinder replication experiments. Nevertheless, irreproducible results sometimes raise suspicion and discussions turn towards concerns regarding the transparency of research or validity of findings. These discussions are typically driven by heuristic thinking about replications. Such heuristics might not hold and can lead to erroneous inferences about research findings and researchers’ practices. To move the needle forward, we have provided a detailed evaluation of which elements of an experiment need to be made open relative to some objective, and which do not. For example, while necessary to audit the results of a given experiment, data sharing is not a prerequisite for performing replications or reproducing results (contrary to some suggestions, for example by [13]), but other elements of an experiment are. On the other hand, reporting model details, such as modelling assumptions, model structure and parameters, becomes critical for improving the accuracy of estimates of reproducibility. Notably, even in recent recommendations for improving transparency in reporting via practices such as preregistration, models are typically left out while transparency of hypotheses, and methods and study design are emphasized [45,46]. Also noteworthy is that some degrees of openness are difficult to attain, such as fully open background knowledge, often causing practical constraints to limit our choices for replication experiments.

When critical elements of an original experiment are not open, replication researchers would be forced to introduce substitutions in their experimental designs. Such substitutions, as we have illustrated, characterize non-exact replications and will probably alter reproducibility rates in different directions, contributing to the challenge of interpreting replication results. Strong theoretical foundations and well-defined shared empirical paradigms in a given area of research could help generate meaningful substitutions whose downstream consequences on inference are well understood.

6.4. Choosing non-exact replications

Assuming a sequence of perfectly repeatable experiments is a theoretical convenience—one that especially frequentist statistics enjoys greatly. In scientific practice, we lack the luxury provided by this assumption. Exact replications are practically impossible. Understanding the implications of result 5.1 is crucial in this respect. It states that any sequence of non-exact replications converges to a true reproducibility rate. This rate may or may not be scientifically meaningful for a specific purpose. Especially for a sequence of non-exact replications, it is hard to find a scientifically meaningful interpretation of what the reproducibility rate shows, even when it is high.

A proper understanding of the elements of the original experiment needs to precede any replication design. And wherever divergences from the original experiment are inevitable, we should strive to theoretically match new design elements to the original ones if our objective is to reproduce an original result. When that is not possible, simulations varying the degree and nature of these divergences would inform us on their impact on the reproducibility rate and can provide guidance in designing non-exact replication experiments. A lack of theoretical understanding in this regard poses significant constraints on the interpretability of replication results.

In cases where the original experiment suffers from design issues that make results predictably less reproducible, it is advisable to iteratively work toward improving the configuration of the idealized experiment first before attempting any non-exact replications [20]. If there is nothing there to revisit, we might be better off saving our scientific curiosity and resources for more fruitful avenues. In fact, there is room for major theoretical advancements on why and how to choose replications.

6.5. Reproducibility of a result versus accumulation of scientific evidence

We hope that advancing theoretical understanding of results reproducibility helps delineate how and why it is different from other quantities that aim to measure the accumulation of scientific evidence. The notion of reproducibility is unique in the sense that it is anchored on the results of an initial experiment. To the contrary, meta-analytic effect size estimates focus on an underlying true effect, after accounting for variation between studies being meta-analysed while robustness tests aim to assess to what extent estimated quantities of interest are sensitive to changes in model specifications. It is a widespread interpretation that reproducibility also speaks to the reliability or validity of an underlying true effect and can reasonably be used as a measure of evidence accumulation. It should be clear by now that this is a misconception. Truth certainly plays a role in reproducibility of a given result but not (always) too loudly, as reproducibility primarily captures patterns specific to the original experiment. A replication experiment in reference to an original result is a particular kind of an idealized experiment that has the capacity for achieving certain scientific objectives, such as confirming a theoretically precise prediction under well-specified conditions (i.e. attempting to account for sampling error as a last source of uncertainty after everything else has already been accounted for) or estimating the reproducibility rate of a particular result of a given experiment. For other scientific objectives, such as to make an initial scientific discovery, to pinpoint the conditions under which a precise and reliable signal can be captured, to aggregate evidence for a theorized phenomenon or to gauge the robustness or heterogeneity of an observed phenomenon across contexts, there are other idealized experiments better suited to the task than replications [20,41] such as systematic exploratory experimentation [47], metastudies [48], multiverse analyses [49], meta-analyses and continuously cumulating meta-analyses [50].⁸ The fact that scientists still care to meticulously design their experiments to be informative and meaningful has more to do with other scientific values and objectives than reproducibility.

In a sense, accumulation of scientific evidence in support of a finding requires epistemic iterations and confirmation by independent approaches and methods to achieve specific scientific objectives (e.g. discovering a new phenomenon, explaining a mechanism, predicting a future observation). This process leads to gradually eliminating uncertainty and enhancing our confidence in our theories and observations. On the other hand, attempts at reproducing a given result in replications prioritize understanding and fine-tuning the logical structure of experiments, which we see as human data-generating mechanisms. Proper appreciation of this aspect of reproducibility is capable of guiding us in the right direction in our struggle to design more rigorous and informative experiments under uncertainty.

6.6. Concluding remarks

The discourse on scientific reform and metascience has so far pursued a ‘crisis’ framing, focusing on behavioural, social, institutional and ethical failings of the scientific endeavour and calling for immediate institutional and collective action. Our analysis shows that neither elimination of scientific malpractice nor actively encouraging replication experiments would necessarily improve the reproducibility of results. Because irreproducibility, when formally defined, appears to be an inherent property of the scientific process rather than a meaningful scientific objective to pursue. While reproducibility rate is a parameter of the system and thereby a function of truth, that view of the concept misses the big picture—that reproducibility reflects the properties of experiments. We perceive two issues with advancing a replication/reproducibility crisis narrative:

1. Conflating replication and reproducibility creates an inaccurate impression that these two alleged issues of not being able to conduct informative replication experiments and not being able to reproduce results are indistinguishable issues that can be addressed via similar solutions.

⁸We have deliberately excluded multi-site replications from this list as there are reasons to suspect that, as they are practised, multi-site replications are not necessarily appropriate for the purposes of a robustness check for reasons detailed in [51]. This is largely on account of each replication being a non-exact replication in a unique and uncontrolled way.

2. Framing irreproducibility as a crisis implies that there is an ideal rate of reproducibility we should expect or strive to achieve in a given field at a given time, and we are falling short of this ideal standard.

Our mathematical results firmly argue against both of these misconceptions.

Shifting the discourse on scientific reform and metascience towards greater theoretical may help change the course of science. Instead of prioritizing crisis management measures, progress can be made by falling back on fundamental issues and working our way from the bottom up. That may require individual scientists to take a step back and reassess the way they have been practising science. Circling back to our original premise, we emphasize that the problem is conceptual: the logical structure of experiments is not well understood and how experiments relate to reality gets misconstrued. Experiments are data-generating machines, and each element outlined in this work determines what kind of data they will generate. Gaining clarity with regard to how experiments impact the observed reality and properly assessing the empirical value of a given experiment for a given objective should precede concerns regarding possible replications. Theory of reproducibility is a step in this direction.

Data accessibility. This article has no additional data.

Authors' contributions. E.O.B.: conceptualization, formal analysis, investigation, methodology, project administration, software, supervision, validation, visualization, writing—original draft and writing—review and editing; B.D.: conceptualization, investigation, methodology, supervision, validation, visualization, writing—original draft and writing—review and editing; B.B.: conceptualization, investigation, validation, writing—original draft and writing—review and editing.

All authors gave final approval for publication and agreed to be held accountable for the work performed therein.

Conflict of interest declaration. We declare we have no competing interests.

Funding. This study was supported by the National Institute of General Medical Sciences of the National Institutes of Health (Award no. P20GM104420).

Acknowledgements. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

Appendix A

Proof. Proof of result 2.7 and an example pertaining to remark 4.1 that meaningful continuous measure of reproducibility is nonetheless pathological. Result 2.7 is a consequence of Strong Law of Large Numbers. An easy proof relies on Kolmogorov's almost everywhere convergence which states that a sequence of independently and identically distributed random variables with finite mean converges almost surely to a constant if and only if that constant is the expected value of random variables. The sequence $\phi^{(1)}, \phi^{(2)}, \dots, \phi^{(N)}$ obtained from $\xi^{(1)}, \xi^{(2)}, \dots, \xi^{(N)}$ (respectively) satisfies Kolmogorov's. By definition 2.5, $\phi_N \in [0, 1]$ and ϕ_i are independent of each other and identically distributed and the expected value is $E(\phi_N) = \phi < \infty$, proving result 2.7.

Importantly, remark 4.1 cautions us that result 2.1 does not hold for all measures of reproducibility. A well-defined ξ and ϕ are prerequisites for result 2.7 to hold. We use a counterexample with a continuous measure of reproducibility to clarify this point. As opposed to a 0–1 measure such as ϕ_N , we consider a (maybe) more desirable measure of reproducibility rate, perhaps a degree of agreement between the results of ξ and ξ' to assess whether r_o from ξ is reproduced in ξ' . One way to represent this degree of agreement is to replace the indicator function in definition 2.5 with a function of a continuous random variable. For example, for a sequence of idealized experiments $\xi^{(1)}, \xi^{(2)}, \dots$ we might define $Y^{(i+1)}/Y^{(i)}$, where $Y^{(i)} \sim \text{Nor}(0, \sigma)$ is a centralized statistic from $\xi^{(i)}$, as score on how extreme is a specific result with respect to an original result $Y^{(o)}$. Here, $Y^{(i)}$ are independent and identically distributed random variables conditional on $\xi^{(i)}$. The set-up is such that if $Y^{(i+1)}/Y^{(i)} = 1$, then the results in $\xi^{(i+1)}$ and $\xi^{(i)}$ have exactly the same degree of agreement. Thus, one can define the reproducibility rate as follows:

$$\phi_N^* = N^{(-1)} \sum_{i=1}^N \left(\frac{Y^{(i+1)}}{Y^{(i)}} \right).$$

This measure of reproducibility rate might seem reasonable, but it is statistically unacceptable. To see this, we substitute ϕ_N with ϕ_N^* , and we see that equation (2.1) is not true and we do not have desirable statistical properties for our estimator of reproducibility [52, p. 12]. Consequently, the statistical justification for the

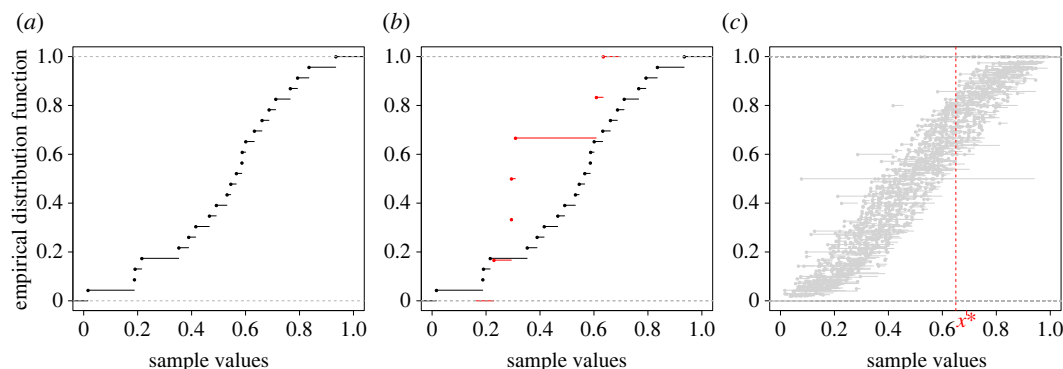


Figure 4. (a) Empirical cumulative distribution function (ECDF) of a sample of size 30, emphasizing that the ECDF is a right continuous function. (b) ECDF of the sample in (a) (black) and that of an independent sample of size 10 (red) emphasizing that the ECDF is a random variable whose probability distribution is determined by the sample values (and hence data-generating mechanism). (c) One hundred independent samples of varying sample size (grey) emphasizing that ECDF is a stochastic process. Red vertical line shows the distribution of ECDF conditional on value x^* .

concept of result reproducibility falls apart. This example shows that one has to define the parameter and its estimator of the reproducibility rate by obeying the constraints of statistically desired properties for reproducibility rate to be a useful concept. It is wise to check that a new concept defined in a developing field is statistically well behaved. Statistical nuances might get lost in applications with important consequences for results reproducibility.

Some additional statistical properties of ϕ_N given in definition 2.5 are as follows. The sampling distribution of ϕ_N is asymptotically normal with $E(\phi_N) = N\phi$ and $Var(\phi_N) = N\phi(1 - \phi)$ by the CLT. All else being equal, the results for which the true reproducibility rate is high or low have low variance for the estimator, and for the results for which the true reproducibility rate is around 0.5, the variance of the point estimator is large (largest when $p = 0.5$). Approximately 100% confidence intervals (and tests of approximately power 1) can arbitrarily be built, with the property that only finitely many of the confidence intervals do not contain the true reproducibility rate ϕ . This result, which fundamentally relies on the law of the iterated logarithm, constitutes a strong basis for statistical methods about ϕ . ■

Appendix B

Proof of result 2.9 (constructive): The sequence of idealized experiments $\xi^{(1)}, \xi^{(2)}, \dots$ given by definition 2.5 is a proper stochastic process, seen as a joint function of random sample D and of each value in the support of data-generating mechanism, $x \in \mathbb{R}$.

K , S and M_A are not stochastic, so we condition on them. $\xi^{(i)}$ draws a simple random sample $D^{(i)} = \mathbf{X}_n^{(i)}$ independent of all else. We note two facts for the proof:

- For fixed \mathbf{X}_n , the sample estimate of M_A is a well-defined probability model for all x . This set-up induces the set of proper probability distribution functions: right continuous cumulative distribution functions with left limits on $[0, 1]$. Three of these cumulative distribution functions are exemplified in figure 4a,b.
- For any fixed x in the support of the cumulative distribution function of M_A , the sample estimate of M_A as a function of the random data \mathbf{X}_n is a random variable, which makes ξ a random variable. This is exemplified in figure 4c, conditional on red line.

Together, (i) and (ii) imply that as a joint function of random D and x , ξ is a proper stochastic process ([52], chapters 1–3) on the space of right continuous functions with left limits on $[0, 1]$. Examples are all grey cumulative distribution functions (CDFs) depicted in figure 4(c).

Result 2.9 is a convenient way to study replications and reproducibility. It has a number of mathematical implications. Firstly, it established that ξ is a well-behaved stochastic process with a limiting distribution. It is of interest to know the limit of this process. It tells us to which point the sample reproducibility rate from replication experiments converge.

Technically, the sequence of probability measures defined for the stochastic process associated with $\xi^{(1)}, \xi^{(2)}, \dots$ on Borel sets with respect to the metric that we describe below has a limiting process that

convergences in distribution. Establishing this convergence helps us to understand the limiting behaviour of $\xi^{(1)}, \xi^{(2)}, \dots$, and characterizing this limiting behaviour. Donsker's Theorem characterizes the limiting process and states that ξ must converge to the Wiener measure. Thus, the probability distribution of the reproducibility rate converges to the normal distribution. Readers interested in the theory of convergence in stochastic processes may refer to [52], chapters 1–3 for details. We give a brief description of necessary background here. There are three essential elements to study the convergence of a proper stochastic process: (i) a proper field on which the process takes values (the class of sets of interest) and a metric associated with it to assess the convergence of the process, (ii) the probability measure that determines the behaviour of the process, and (iii) using (i) and (ii), a complete mathematical formulation of the stochastic process, which can be used to show convergence to some well-defined distribution.

We now consider a stochastic process as a function of $t \in [0, 1]$, a random point in the space of right continuous functions on $[0, 1]$ with left-hand limits. We let the supremum of the L1 norm between any two points in the space and the metric to assess the convergence to be the classical Kolmogorov–Smirnov distance. By $[nt]$, we denote the floor function, the integer part of nt . Given $\{X_n = (X_1, X_2, \dots, X_n); n \in \mathbb{Z}^+\}$, where X_i are independent of each other and identically distributed, we define the stochastic process defined on partial sums,

$$\frac{\sum_{i=1}^{[nt]} [X_i - E(X_i)] + [nt - [nt]] [X_{[nt]+1} - E(X_i)]}{\sqrt{n \text{Var}(X_i)}}.$$

For elements of this process, if we denote the probability distribution for a sample size n by P_n , then the limiting distribution is the well-known Wiener measure, \mathcal{W} . Some results follow from this.

ξ is most generic when M_A is *any* probability model. This induces S_{post} having the sampling distribution function of *any* statistic. In this most generic case, the distribution of the sample reproducibility rate ϕ_N for the sequence $\xi^{(1)}, \xi^{(2)}, \dots$ is asymptotically normal. To see this, we first let $X_n = M_A^{-1}(w)$, where $w \in [0, 1]$ so that we have the image of the statistical model and assume that ϕ_N evaluated at 0 and 1 is 0. The stochastic process

$$\sqrt{n} \{ \xi[M_A^{-1}(w)] - w \}$$

converges to a specific Wiener process, with bound endpoints, which is a Brownian bridge: the process is Gaussian with zero expectation and, for two points w_1 and w_2 , the covariance function $\text{Cov}(\mathcal{W}(w_1), \mathcal{W}(w_2)) = w_1(1 - w_2)$, with the ordering $w_1 \leq w_2$, and $w_i \in [0, 1]$.

By definition of this stochastic process and its convergence to a Brownian bridge, we see that for each fixed value of x , ξ is asymptotically normally distributed with mean M_A and variance $M_A(1 - M_A)/n$.

The result can also be studied fixing one dimension at a time, giving two corollaries. For random data X_n , the elements of the sequence of replication experiments $\xi^{(1)}, \xi^{(2)}, \dots$ are random variables and conditionally independent of each other. For fixed data, the elements of the sequence of replication experiments $\xi^{(1)}, \xi^{(2)}, \dots$ are probability models.

Appendix C

Details on remark 4.1: Let ξ be an idealized experiment and ξ' be its exact replication. Conditional on R from ξ , K is necessarily distinct from K for epistemic reproducibility of R by R' , but not necessarily distinct for in-principle reproducibility of R by R' .

We define and distinguish *in-principle* reproducibility and *epistemic* reproducibility conditional on a result R . It is clear that π -openness, where π is a non-empty set and is necessary to make the elements of ξ available for replication ξ' . Further, R also needs to be open for ξ' to be able to determine whether R' has epistemically reproduced R . So, information on R across the sequence of replication experiments is a logical necessity for *epistemic* reproducibility. As an example, consider two scenarios 1 and 2. In each scenario, there are two experiments, the originals (ξ_1 and ξ_2 , respectively) and their replications (ξ'_1 and ξ'_2 , respectively). Each experiment assumes an infinite population of black and white ravens (A). ξ_1 and ξ_2 have identical M_A , S and D_s . R is the estimate \hat{p} of the population proportion of black ravens p , obtained using an independent D_v . We assume that the number of black ravens b observed in ξ_1 and ξ'_1 , and ξ_2 and ξ'_2 are the same.

Closed scenario: The experiments are isolated from each other, and there is no information flow from ξ_1 to ξ'_1 . Thus, ξ'_1 can only match all the elements of ξ_1 that are relevant to \hat{p} either by *chance* or by an extreme precision of prior theoretical formulation. By our example, ξ_1 and ξ'_1 have identical M_A , D_s and S and

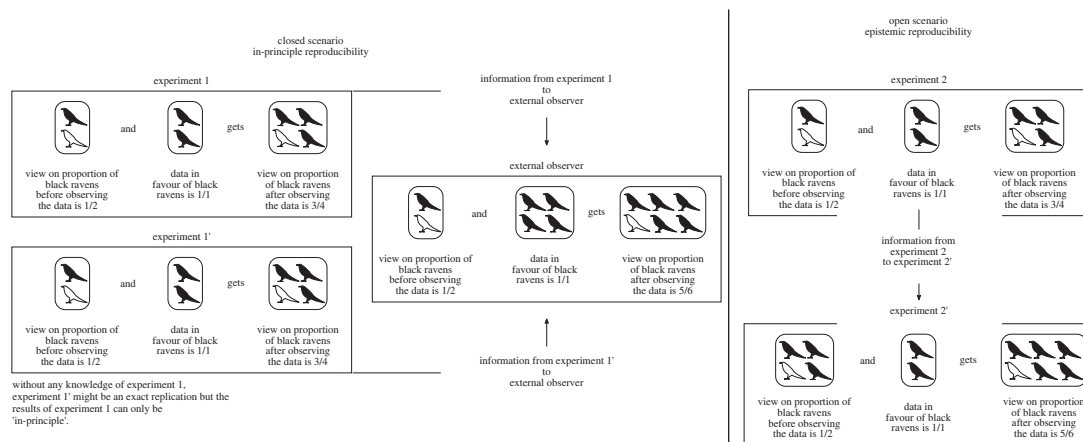


Figure 5. Epistemic versus in-principle reproducibility with an example of Bayesian information flow and learning (details are provide within appendix text).

have the same observed value b in the sample; thus, they return the same estimate \hat{p} . However, ξ_1' does not have any information pertaining to R from ξ_1 , and thus, ξ_1' is in a position neither to learn from R of ξ_1 nor to claim that it reproduced the result of ξ_1 by R' . If an external observer were to observe the experiments ξ_1 and ξ_1' , they could learn from the results of both experiments simultaneously. Starting with a prior view of equal proportion of black and white ravens, they could use the number of ravens observed in ξ_1 and ξ_1' , to conclude that R of ξ_1 is indeed reproduced by R' of ξ_1' and arrive at an updated view. When there is no information exchange with regard to R between the ξ_1 and ξ_1' , however, there is no meaningful or immediate *epistemic* interaction between ξ_1 and ξ_1' , and there is no knowledge of reproducibility unless an all-knowing third party is involved.

This closed scenario shows that if there is no openness in the sense of information flow from one experiment to the next, it is improbable (but still possible) for an experiment to reproduce the result of another experiment. In order to acknowledge this point, we say that a result can only be *in-principle* reproducible if there is no epistemic exchange between ξ_1 and ξ_1' which could speak to the reproducibility of R , with the exception of via some omniscient external observer. At times, historians of science illustrate such examples of scientific discoveries independently arrived at by different scientists unaware of each other's work.

Open scenario: There is information flow from ξ_2 to ξ_2' , with respect to R and other information relevant to obtain \hat{p} in ξ_2 . If ξ_2' incorporates this information, it is a replication. Here, ξ_2' matches the elements of ξ_2 by *social learning*. The information necessary for learning is transmitted in K and R . Starting with \hat{p} as R , ξ_2' could conclude that they have indeed reproduced it. Thus, in the open scenario, there is an *epistemic* interaction between ξ_2 and ξ_2' which contributes to the progress of science through deliberate transfer of knowledge via social learning, which gives us the notion of *epistemic reproducibility*.

As an example, we show the difference between *epistemic reproducibility* and *in-principle reproducibility* in figure 5 with an infinite population of black and white ravens and Bayesian inference. Figure 5a illustrates the closed scenario: researchers of ξ_1 assume a prior view of 1/2 on \hat{p} . After observing $n = 2$ black ravens, they update their view to $\hat{p} = 3/4$ by Bayesian inference. Researchers of ξ_1' assume a prior view of 1/2 on \hat{p} and observe identical D_o , $n = 2$ black ravens as in ξ_1 , and they update their view with same S_{post} , to reach $\hat{p} = 3/4$. However, in the absence of an external observer, these two results cannot be epistemically connected, thus reproducibility is only *in principle* in the absence of an external observer privy to both experiments. Figure 5b illustrates the open scenario: researchers of ξ_2 assume a prior view of 1/2 on \hat{p} . After observing $n = 2$ black ravens, they update their view to $\hat{p} = 3/4$ by Bayesian inference. ξ_2' is a proper replication experiment. It is informed by the result of ξ_2 as well as K , M_A , S and D_s and observes identical D_o as ξ_2 . ξ_2' , starting with a view of $\hat{p} = 3/4$ from ξ_2 , they update their view to $\hat{p}' = 5/6$. Thus, ξ_2' learns from ξ_2 , here in a Bayesian manner. The two results can be connected and thus reproducibility is epistemic.

Appendix D

Proof of result 4.2: M_A and M'_A do not have to be identical in order to reproduce a result R by R' . Under mild assumptions, the requirement for R to be reproducible by R' is that there exists a one-to-one transformation between M_A and M'_A for inferential purposes mapping to R .

We first give a proof for the statement and then follow with a specific example. We let $F_X(x)$ and $F_{X'}(x)$ be distribution functions with inverses $F_X^{-1}(x)$ and $F_{X'}^{-1}(x)$, under ξ and ξ' , respectively. By assumption, a one-to-one function, g , from $F_{X'}(x)$ to $F_X(x)$ exists. For two distribution functions whose inverses exist, the mapping of population quantiles from one to the other also exists if there is a one-to-one function between these distribution functions. All well-behaved (non-order) statistics can be represented as quantiles, so we prove result 4.2 without loss of generality by setting the quantity of inferential interest as x_q , where $F_X(x_q) = P(X \leq x_q) = q \in [0, 1]$. If using equivalent estimators of quantiles with samples from M_A and M'_A respectively, then the mapping carries over to R to R' . We have

$$x_a = F_{X'}^{-1}(a) = E_{X'}^{-1}[\mathbf{I}_{\{X' \leq x_a\}}] = g\{E_X^{-1}[\mathbf{I}_{\{X \leq x_b\}}]\} = F_X^{-1}(b) = x_b, \quad (\text{D } 1)$$

where $\mathbf{I}_{\{A\}} = 1$ if A and 0 otherwise. Equation (D 1) holds for estimators of population quantities and an estimator of x_a can be equated to an estimator of x_b via a one-to-one transformation g , by replacing the population quantities with their estimators. This result applies to non-parametric and parametric models alike, in fact to all distributions with finite means and well-defined inverses.

As an example with two parametric models from figure 1, we consider the problem of estimating the proportion of black ravens, p using ξ_{bin} as the original experiment and ξ_{negbin} as its replication. The characteristic function of ξ_{bin} and ξ_{negbin} are $(1 - p + p \text{eit})^n$ and $p^w(1 - e^{it} + p \text{eit})^{-w}$, respectively. The characteristic function for a random variable fully defines its probability model and thus, ξ_{bin} and ξ_{negbin} have distinct models. Yet p is an identifiable and estimable parameter of both experiments. The maximum likelihood estimator of p is $\hat{p}(\xi_{\text{bin}}) = \hat{p}(\xi_{\text{negbin}}) = b/n$ because ξ_{bin} and ξ_{negbin} are in a *likelihood equivalence class* with respect to parameter p . To see this, we note that the maximum likelihood estimator is obtained by setting the expression resulting from taking the derivative of the logarithm of the likelihood function (i.e. score function) with respect to p and solving for p . Under ξ_{bin} the score function is

$$\frac{d}{dp} [\log \mathbb{P}(b|p, n)] = \frac{d}{dp} (\log C_{n,b}) + \frac{d}{dp} (b \log p) + \frac{d}{dp} [w \log(1 - p)]. \quad (\text{D } 2)$$

Under ξ_{negbin} the score function is

$$\frac{d}{dp} [\log \mathbb{P}(b|p, w)] = \frac{d}{dp} (\log C_{n-1, w-1}) + \frac{d}{dp} (b \log p) + \frac{d}{dp} [w \log(1 - p)]. \quad (\text{D } 3)$$

Equations (D 2) and (D 3) differ only in their first terms, which is irrelevant to estimate p , and thus, $\hat{p}(\xi_{\text{bin}}) = \hat{p}(\xi_{\text{negbin}}) = b/n$ is the unique solution. The first terms on the right-hand side of these two equations determine the stopping rule of the experiments. In ξ_{bin} , we stop the experiment when n ravens are observed and the last raven can be black or white. In ξ_{negbin} , we stop the experiment when w white ravens are observed and the last observation must be a white raven. This difference between stopping rules means that (i) S_{pre} is different from S'_{pre} : (ii) under our choice of S_{post} and S'_{post} as the maximum likelihood estimator, the stopping rules in two models are irrelevant for estimating the proportion of black ravens in the population.

Appendix E

Proof of result 4.4: S_{post} and S'_{post} do not have to be identical in order to reproduce a result R by R' .

There are a few heuristic ways to derive well-behaved statistical estimators of parameters. Examples include: method of moments, maximum likelihood and posterior mode (Bayesian). Well-known estimators may be equal to each other in value but motivated by distinct principles. For example, for some distinct probability models in the exponential family, the method of moments and the maximum likelihood estimator return the same value. Or, by using uniform prior in Bayesian inference, the posterior mode always returns the same value as the maximum likelihood estimator. This motivates result 4.4 in the sense that S_{post} and S'_{post} do not have to be identical to reproduce R by R' .

As an example based on ξ_{bin} from figure 1, we consider the following three estimators:

- If S_{post} is the maximum likelihood estimator motivated by the likelihood principle, then we have (see appendix D)

$$\hat{p}_{MLE} = \frac{b}{n}.$$

- If S_{post} is the MME, the motivation is to set the population mean equal to the sample mean and solve for p , and we have

$$\hat{p}_{\text{MME}} = \frac{b}{n}.$$

- If S_{post} is the posterior mode under the uniform prior (a special case of conjugate prior for ξ_{bin}), we have

$$\hat{p}_{\text{MP}} = \frac{b}{n}.$$

Therefore, ξ can employ any one of these three estimators as S_{post} , and ξ' can employ another as S'_{post} and still reproduce R by R' , as if they have used the same statistical method. For other modes of statistical inference such as hypothesis tests and prediction, we can find examples of numerically equivalent methods that are not identical in motivation (e.g. [53]).

Appendix F

Proof of result 4.5: D_s and D'_s do not have to be identical in order to reproduce a result R by R' .

The data structures of probability models that correspond to ξ_{bin} , ξ_{negbin} , ξ_{hyper} , ξ_{poi} , ξ_{exp} and ξ_{nor} are all distinct. In ξ_{bin} and ξ_{negbin} , the data structures are a sample of size n ravens and a sample of size w white ravens, respectively, from an infinite population in which p is constant. Stopping rules of the sampling in these experiments are different from each other: the last raven must be white in ξ_{negbin} but not in ξ_{bin} . In ξ_{hyper} , the stopping rule is the same as ξ_{bin} , but the parameter p changes with each sample obtained due to finite population assumption in ξ_{hyper} .

In ξ_{poi} and ξ_{exp} , $np \rightarrow \lambda$ is the rate of black ravens appearing in the process. The observable in ξ_{poi} is the random variable b_t , the count of black ravens at time t , and we denote the count of black ravens at time $t + \delta$ by $b_{t+\delta}$. The observable in ξ_{exp} is the random waiting time δ to observe another black raven assuming a black raven is observed at time t . The equivalence between the parameters of ξ_{poi} and ξ_{exp} is given by

$$P(T \leq t) = 1 - P(b_{t+\delta} - b_t = 0), \quad (\text{F } 1)$$

where the cumulative distribution function of the time variable in M_A in ξ_{exp} is related to the counts in M_A in ξ_{poi} by probability of no event in time period δ . Equation (F 1) implies that no black raven is observed in δ . By Poisson probability mass function, we have this probability as $P(b_{t+\delta} - b_t = 0) = e^{-\delta}$, and we have $P(T \leq t) = 1 - e^{-\delta}$. This identifies T as an exponential random variable in ξ_{exp} implying that the data structures in ξ_{poi} and ξ_{exp} are distinct. Yet, irrespective of all these differences in data structures, ξ_{bin} and ξ_{negbin} estimate the same parameter, p . Further, ξ_{poi} and ξ_{exp} also estimate the same parameter, λ . Hence, R can be reproduced by R' without the data structures being identical in ξ and ξ' .

Appendix G

Proof of result 5.1: Assume a sequence $\xi^{(1)}, \xi^{(2)}, \dots, \xi^{(j)}$ of idealized experiments in which a result R is of interest. Then, the estimated reproducibility rate of R in this sequence converges to the mean reproducibility rate of R in J replication experiments.

Conditional on all other elements of an idealized experiment, definition 2.5 and consequently equation (2.1) assume that data are generated independently in each replication which implies that $R^{(i)}$ are independent and identically distributed random variables. Result 5.1 is straightforward for independent and identically distributed random variables. Unconditionally on the elements, however, $R^{(1)}, R^{(2)}, \dots$ in the sequence $\xi^{(1)}, \xi^{(2)}, \dots$ are *not* independent and identically distributed, implying that $R^{(i)}$ are not drawn from the same sampling distribution of results. An easy way to see this is to pick $\xi^{(i)}$ and $\xi^{(j)}$ distinct at least with respect to one element. In exact replications, $R^{(i)}$ and $R^{(j)}$ will converge to their unique true reproducibility rate $\phi^{(i)}$ and $\phi^{(j)}$ by equation (2.1). However, equation (2.1) can be generalized to obtain result 5.1 as follows using a theorem due to Kolmogorov (see [54]).

We let $\phi_N^{(1)}, \phi_N^{(2)}, \dots$ be estimates of reproducibility rates, with means $\phi^{(1)}, \phi^{(2)}, \dots$ and variances $N^{-1}\phi^{(1)}(1-\phi^{(1)}), N^{-1}\phi^{(2)}(1-\phi^{(2)}), \dots$, respectively. We assume that the series $\sum_{i=1}^{\infty} i^{-1}\phi^{(i)}(1-\phi^{(i)})$ converges. Then,

$$N^{-1} \sum_{i=1}^N \phi_N^{(i)} \rightarrow N^{-1} \sum_{i=1}^N \phi^{(i)}, \text{ almost surely.} \quad (\text{G } 1)$$

Expression (G 1) states that the estimated reproducibility rate of results from non-exact replication experiments meaningfully converges to the mean true reproducibility rate of the idealized experiments performed. The case of exact replications given by equation (2.1) is a special case of the equation (G 1), where all non-exact replications are identical to each other (and thus exact) with respect to the result obtained in an original idealized experiment. That is, if equation (G 1) is applied to $\xi \equiv \xi^{(1)} \equiv \xi^{(2)} \equiv \dots \equiv \xi^{(N)}$, where the true reproducibility rate for R_0 obtained from ξ is ϕ , and we obtain

$$N^{-1} \sum_{i=1}^N \phi_N^{(i)} \rightarrow N^{-1} \sum_{i=1}^N \phi^{(i)} = N^{-1} \sum_{i=1}^N \phi = \phi, \text{ almost surely.} \quad (\text{G } 2)$$

Appendix H

Reproducibility rate of R as a model selection problem, in the context of linear regression models.

In addition to the simulation example given in figure 3, here we present a second simulation example to illustrate the convergence of reproducibility rates from exact and non-exact replication experiments to their true value. Our example involves the model selection problem in the context of linear regression models. Briefly, we assume the linear regression model

$$\mathbf{y} = \mathbf{X}\beta + \epsilon,$$

where \mathbf{y} is $n \times 1$ vector of responses, \mathbf{X} is $n \times k$ matrix of fixed observables with first column entries equal to 1, β is $k \times 1$ vector of parameters and ϵ is $n \times 1$ vector of independent and identically distributed normal errors with mean 0 and unknown variance. The statistical problem is as follows: Given D with D_v independent and identically distributed and D_s constituting $n \times 1$ responses and $n \times k$ observables, select the best linear regression model among three models with respect to a model selection criterion (S_{post}). The saturated model is given by

$$\mathbf{X} = (\mathbf{1}, \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3), \text{ with } \beta = (\beta_0, \beta_1, \beta_2, \beta_3),$$

where $\mathbf{x}_1, \mathbf{x}_2$ and \mathbf{x}_3 are $n \times 1$ vectors of first, second and third predictors, respectively, and β_1, β_2 and β_3 are their respective regression coefficients. The set of three models considered in the model selection problem are as follows:

1. $\mathbf{X} = (\mathbf{1}, \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)$, with $\beta = (\beta_0, \beta_1, \beta_2, \beta_3)$,
2. $\mathbf{X} = (\mathbf{1}, \mathbf{x}_1, \mathbf{x}_2)$, with $\beta = (\beta_0, \beta_1, \beta_2)$,
3. $\mathbf{X} = (\mathbf{1}, \mathbf{x}_1, \mathbf{x}_3)$, with $\beta = (\beta_0, \beta_1, \beta_3)$.

In all cases, the true model generating the data is model 3. For each ξ (and their exact replications), we vary four elements $M_A, R, S_{\text{post}}, D_s$ in a $2 \times 2 \times (2 \times 2 + 1)$ simulation study design:

1. M_A : Model as determined by the signal-to-noise ratio in the true model generating the data. Two values are: signal:noise = 1:1, which is equivalent to the statistical condition $E(Y): \sigma = 1:1$, and signal:noise = 1:4, which is equivalent to the statistical condition $E(Y): \sigma = 1:1$ in figure 6.
2. R : Result of the original experiment. Two values are: true and false.
3. S_{post} : Model selection method. Two values are: Akaike's information criterion (AIC) and Bayesian information criterion (BIC).
4. D_s : Data structure. Two values are as follows: sample sizes $n = 10$ and $n = 100$.
5. Additionally: Uniformly randomly chosen non-exact replications at each step of the sequence from the set of all idealized experiments.

We performed 100 runs of a sequence of 1000 exact replication experiments for each of the 16 experimental conditions, plus 100 runs of a sequence of 1000 non-exact replication experiments where $(M_A, R, S_{\text{post}}, D_s)$ is chosen uniformly randomly from 16 conditions. Four of the experimental conditions (M_A, R values) are

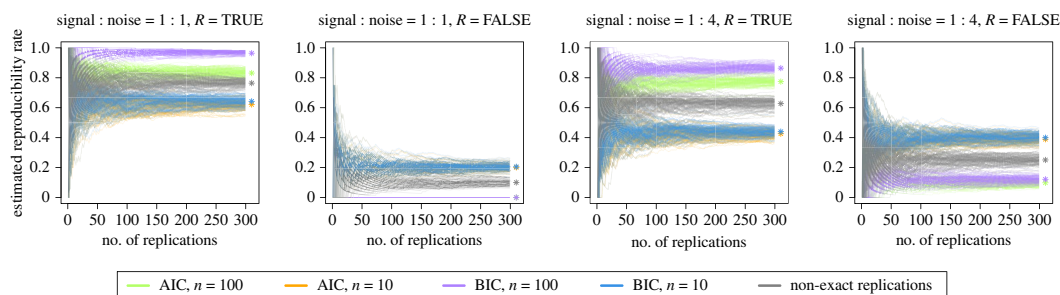


Figure 6. A simulation example to illustrate the convergence of reproducibility rates from exact and non-exact replication experiments to their true value. See text within the appendix for description of this figure.

shown in figure 6: (a) signal : noise = 1 : 1 and TRUE result; (b) signal : noise = 1 : 1 and FALSE result; (c) signal : noise = 1 : 4 and TRUE result; (d) signal : noise = 1 : 4 and FALSE result. Five experimental conditions (S_{post} , D_s values + non-exact replications) are shown in coloured lines in figure 6. Green: AIC, $n = 100$; orange: AIC, $n = 10$; purple: BIC, $n = 100$; blue: BIC, $n = 10$; grey: non-exact replications. The result of interest is the reproducibility rate of the result of the original experiment. The plots in figure 6 display runs of replication experiments indexed by colour. Each run converges to a point (indicated by star) representing the true reproducibility rate of a given run. As a whole, these plots illustrate how true reproducibility rate changes depending on the elements of ξ and the effect of divergence of ξ' from ξ . We emphasize that all parameters of the simulation example in figure 6 are chosen so that one can discern the effect of varying models, methods, and data structures.

We interpret the results as follows:

1. The reproducibility rates for false results and for true results sum to 1, which is a verification of simulation experiments.
2. By the true rates of reproducibility marked by stars, we observe that they depend on the true data-generating mechanism, and the elements of the original experiment, S_{post} and D_s . For example, as the noise increases, the true reproducibility rate gets smaller, and the variance of the estimated reproducibility rate increases. So for larger noise, replication results are expected to be highly variable. True reproducibility rates of true results also change with sample size and method.
3. Reproducibility rate increases with the sample size for true results, whereas it decreases for false results such that low sample size makes false results more reproducible in our simulations.
4. Even when the true reproducibility rate is high, we might see a lot of variation in observed reproducibility rate after a small number of replications even when they are exact replications. Non-exact replications yield highly variable observed reproducibility rates that do not converge to the true reproducibility rate of the original result.

This simulation experiment complements the one presented in the main text (figure 3) by providing a different illustration from our toy example. The context of linear regression models is readily relevant to many practising scientists. Moreover, this simulation extends the results to new contexts by observing the outcome of interest under different levels of system noise and both true and false original results. Ultimately both simulations show considerable variability in true reproducibility rates as a function of the elements of and relationship between original and replication experiments.

Appendix I

True results are not necessarily reproducible and perfectly reproducible results may not be true.

Reproducibility is a function of the true unknown data-generating model and the elements of ξ . Devezet *et al.* [10] provide some account. We give a brief overview with a proof by counterexample. Conditional on R from ξ , we let $\xi^{(1)}, \xi^{(2)}, \dots$ be exact replications of ξ and $\mathbf{I}_{\{b^*\}}$ be the indicator function that equals 1 if the first raven in the sample is black, and 0 otherwise. To prove the first part of the statement, we choose the estimator

$$\hat{p} = \frac{b + \mathbf{I}_{\{b^*\}}}{n + \mathbf{I}_{\{b^*\}}}.$$

The estimator \hat{p} is valid on $[0, 1]$ by: if $b = n$, then the first raven sampled must be black and $\hat{p} = 1$, else if $b = 0$, then the first raven must be white and $\hat{p} = 0$ such that $\hat{p} \in [0, 1]$. However, \hat{p} is unbiased for p only

with probability $(1 - p)$. The reason is that the probability that first raven is white raven is $(1 - p)$, and if it is a white raven, we get $\hat{p} = b/n$ giving $E(\hat{p}) = E(b/n) = (1/n)(np) = p$. In contrast, \hat{p} is biased for p with probability $(1 - p)$. The reason is that the probability that first raven is black raven is p , and if it is a black raven, we obtain $E(\hat{p}) \neq p$. This does not only show that the true results are not always reproducible but also shows that the reproducibility rate can be a function of the true parameter.

To prove the second part of the statement, choose the estimator $\hat{p} = c$, where c is a constant in $[0, 1]$. $E(\hat{p}) = c$. This expectation is only equal to p when $p = c$. However, the result using this \hat{p} is reproducible with probability 1, thereby completing the proof.

References

- Open Science Collaboration. 2015 Estimating the reproducibility of psychological science. *Science* **349**, aac4716–1–aac4716–8.
- Leonelli S. 2018 Rethinking reproducibility as a criterion for research quality. In *Including a symposium on Mary Morgan: curiosity, imagination, and surprise*, vol. 36B, pp. 129–146. Bingley, UK: Emerald Publishing Limited.
- Radder H. 1992 Experimental reproducibility and the experimenters' regress. In *PSA: Proc. of the Biennial Meeting of the Philosophy of Science Association*, vol. 1992, pp. 63–73. Philosophy of Science Association.
- Radder H. 1996 *In and about the world: philosophical studies of science and technology*. Albany, NY: SUNY Press.
- Fidler F, Wilcox J. 2018 Reproducibility of scientific results. In *The Stanford encyclopedia of philosophy* (ed. EN Zalta). Stanford, CA: Metaphysics Research Lab, Stanford University.
- Penders B, Holbrook JB, de Rijke S. 2019 Rinse and repeat: understanding the value of replication across different ways of knowing. *Publications* **7**, 1–15. (doi:10.3390/publications7030052)
- Camerer CF *et al.* 2016 Evaluating replicability of laboratory experiments in economics. *Science* **351**, 1433–1436. (doi:10.1126/science.aaf0918)
- Stodden V. 2011 Trust your science? Open your data and code. *Amstat News* **July**, 21–22.
- Baumgartner B, Devezer B, Buzbas EO, Nardin LG. 2018 Openness and reproducibility: insights from a model-centric approach. *arXiv preprint*. (<https://doi.org/10.48550/arXiv.1811.04525>)
- Devezer B, Navarro DJ, Vandekerckhove J, Ozge Buzbas E. 2021 The case for formal methodology in scientific reform. *R. Soc. Open Sci.* **8**, 200805. (doi:10.1098/rsos.200805)
- Collins FS, Tabak LA. 2014 Policy: NIH plans to enhance reproducibility. *Nat. News* **505**, 612–613. (doi:10.1038/505612a)
- Iqbal SA, Wallach JD, Khoury MJ, Schully SD, Ioannidis JP. 2016 Reproducible research practices and transparency across the biomedical literature. *PLoS Biol.* **14**, 1–13. (doi:10.1371/journal.pbio.1002333)
- National Academies of Sciences, Engineering, and Medicine. 2017 *Fostering integrity in research*. Washington, DC: National Academies Press.
- Nosek BA *et al.* 2015 Promoting an open research culture. *Science* **348**, 1422–1425. (doi:10.1126/science.aab2374)
- Nosek BA *et al.* 2022 Replicability, robustness, and reproducibility in psychological science. *Annu. Rev. Psychol.* **73**, 719–748. (doi:10.1146/annurev-psych-020821-114157)
- Kerr NL. 1998 Harking: hypothesizing after the results are known. *Pers. Soc. Psychol. Rev.* **2**, 196–217. (doi:10.1207/s15327957pspr0203_4)
- Munaò MR *et al.* 2017 A manifesto for reproducible science. *Nat. Human Behav.* **1**, 1–9.
- Bruns SB, Ioannidis JPA. 2016 *p*-curve and *p*-hacking in observational research. *PLoS ONE* **11**, 1–13. (doi:10.1371/journal.pone.0149144)
- 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. Columbia, SC: Department of Statistics, Columbia University **348**, 1–7. See http://www.stat.columbia.edu/~gelman/research/unpublished/p_hacking.pdf.
- Feest U. 2019 Why replication is overrated. *Philos. Sci.* **86**, 895–905. (doi:10.1086/705451)
- Devezer B, Nardin LG, Baumgartner B, Buzbas EO. 2019 Scientific discovery in a model-centric framework: reproducibility, innovation, and epistemic diversity. *PLoS ONE* **14**, 1–23. (doi:10.1371/journal.pone.0216125)
- Wagenmakers E-J, Wetzels R, Borsboom D, van der Maas HL, Kievit RA. 2012 An agenda for purely confirmatory research. *Perspect. Psychol. Sci.* **7**, 632–638. (doi:10.1177/1745691612463078)
- Lindley DV. 2000 The philosophy of statistics. *J. R. Stat. Soc.: Ser. D (The Statistician)* **49**, 293–337.
- Feyerabend P. 2010 *Against method*, 4th edn. London, UK: Verso Books.
- Murre JM. 2021 The Godden and Baddeley (1975) experiment on context-dependent memory on land and underwater: a replication. *R. Soc. Open Sci.* **8**, 200724. (doi:10.1098/rsos.200724)
- Godden DR, Baddeley AD. 1975 Context-dependent memory in two natural environments: on land and underwater. *Br. J. Psychol.* **66**, 325–331. (doi:10.1111/j.2044-8295.1975.tb01468.x)
- Hyman I. 2021. [@Ira_Hyman]. (2021, November 4). I have seen recent tweets on a failure to replicate context dependent memory – the Godden & Baddeley 1975 scuba diving study. Should you drop this from teaching? No. Place dependent memory effects are huge, reliable. Some studies will fail to create the effect. Thread. 1/26 [Tweet]. Twitter. See https://twitter.com/ira_hyman/status/1456336780412653584.
- Botvinik-Nezer R *et al.* 2020 Variability in the analysis of a single neuroimaging dataset by many teams. *Nature* **582**, 84–88. (doi:10.1038/s41586-020-2314-9)
- Silberzahn R *et al.* 2018 Many analysts, one data set: making transparent how variations in analytic choices affect results. *Adv. Methods Pract. Psychol. Sci.* **1**, 337–356. (doi:10.1177/2515245917747646)
- Cooper RP, Guest O. 2014 Implementations are not specifications: specification, replication and experimentation in computational cognitive modeling. *Cogn. Syst. Res.* **27**, 42–49. (doi:10.1016/j.cogsys.2013.05.001)
- Guest O, Martin AE. 2021 How computational modeling can force theory building in psychological science. *Perspect. Psychol. Sci.* **16**, 789–802. (doi:10.1177/1745691620970585)
- Loken E, Gelman A. 2017 Measurement error and the replication crisis. *Science* **355**, 584–585. (doi:10.1126/science.aal3618)
- Stanley DJ, Spence JR. 2014 Expectations for replications: are yours realistic? *Perspect. Psychol. Sci.* **9**, 305–318. (doi:10.1177/1745691614528518)
- Grujters SL. 2022 Making inferential leaps: manipulation checks and the road towards strong inference. *J. Exp. Soc. Psychol.* **98**, 104251. (doi:10.1016/j.jesp.2021.104251)
- Hardwicke TE *et al.* 2018 Data availability, reusability, and analytic reproducibility: evaluating the impact of a mandatory open data policy at the journal cognition. *R. Soc. Open Sci.* **5**, 1–18. (doi:10.1098/rsos.180448)
- Molloy JC. 2011 The open knowledge foundation: open data means better science. *PLoS Biol.* **9**, 1–4. (doi:10.1371/journal.pbio.1001195)
- Borgman CL. 2012 The conundrum of sharing research data. *J. Am. Soc. Inf. Sci. Technol.* **63**, 1059–1078. (doi:10.1002/asi.22634)
- Janssen M, Charalabidis Y, Zuidewijk A. 2012 Benefits, adoption barriers and myths of open data and open government. *Inf. Syst. Manage.* **29**, 258–268. (doi:10.1080/10580530.2012.716740)

39. Box GE. 1976 Science and statistics. *J. Am. Stat. Assoc.* **71**, 791–799. (doi:10.1080/01621459.1976.10480949)
40. Dennis B, Ponciano JM, Taper ML, Lele SR. 2019 Errors in statistical inference under model misspecification: evidence, hypothesis testing, and AIC. *Front. Ecol. Evol.* **7**, 1–28. (doi:10.3389/fevo.2019.00372)
41. Bak-Coleman J, Mann RP, West J, Bergstrom CT. 2022 Replication does not measure scientific productivity. *SocArXiv* preprint. (doi:10.31235/osf.io/rkyf7)
42. Chang H. 2004 *Inventing temperature: measurement and scientific progress*. Oxford, UK: Oxford University Press.
43. Kuhn TS. 1996 *The structure of scientific revolutions*, 3rd edn. Chicago, IL: The University of Chicago Press.
44. Schauer J, Hedges L. 2021 Reconsidering statistical methods for assessing replication. *Psychol. Methods* **26**, 127. (doi:10.1037/met0000302)
45. Nosek BA, Ebersole CR, DeHaven AC, Mellor DT. 2018 The preregistration revolution. *Proc. Natl Acad. Sci. USA* **115**, 2600–2606. (doi:10.1073/pnas.1708274114)
46. van't Veer AE, Giner-Sorolla R. 2016 Pre-registration in social psychology—a discussion and suggested template. *J. Exp. Soc. Psychol.* **67**, 2–12.
47. Steinle F. 1997 Entering new fields: exploratory uses of experimentation. *Philos. Sci.* **64**, S65–S74. (doi:10.1086/392587)
48. Baribault B, Donkin C, Little DR, Trueblood JS, Oravecz Z, van Ravenzwaaij D, White CN, De Boeck P, Vandekerckhove J. 2018 Metastudies for robust tests of theory. *Proc. Natl Acad. Sci. USA* **115**, 2607–2612. (doi:10.1073/pnas.1708285114)
49. Steegen S, Tuerlinckx F, Gelman A, Vanpaemel W. 2016 Increasing transparency through a multiverse analysis. *Perspect. Psychol. Sci.* **11**, 702–712. (doi:10.1177/1745691616658637)
50. Fletcher SC. 2021 How (not) to measure replication. *Eur. J. Philos. Sci.* **11**, 57. (doi:10.1007/s13194-021-00377-2)
51. Buzbas E, Devezer B. 2022 Bridging the gap between formal theory and scientific reform practices. *bioRxiv* preprint. (doi:10.1101/2022.12.07.519533)
52. Serfling RJ. 1980 *Approximation theorems of mathematical statistics*. New York, NY: John Wiley and Sons.
53. Shively T, Walker S. 2013 On the equivalence between Bayesian and classical hypothesis testing. (<http://arxiv.org/abs/1312.0302>)
54. Rao CR. 1973 *Linear statistical inference and its applications*, vol. 2. New York, NY: Wiley.