

Recommendations for Increasing the Transparency of Analysis of Pre-Existing Datasets

Sara J. Weston¹, Stuart J. Ritchie², Julia M. Rohrer³, Andrew K. Przybylski^{4,5}

¹Department of Psychology, University of Oregon, Eugene, OR

²Social, Genetic and Developmental Psychiatry Centre, King's College London, London, UK

³Department of Psychology, University of Leipzig, Leipzig, Germany; International Max Planck Research School on the Life Course, Max Planck Institute for Human Development, Berlin, Germany; German Institute for Economic Research (DIW Berlin), Berlin, Germany.

⁴Oxford Internet Institute, University of Oxford, Oxford, United Kingdom.

⁵Department of Experimental Psychology, University of Oxford, Oxford, United Kingdom.

Abstract

Secondary data analysis, or the analysis of pre-existing data, can be a powerful tool for the resourceful researcher. Never has this been more true than now, when technological advances allow for easier sharing of data across labs and continents and the mining of large sources of “pre-existing data”. However, secondary data analysis is often ignored as a methodological tool, either when developing new open science practices or improving analytic methods for robust data analysis. In this paper, we hope to provide researchers with the knowledge necessary to incorporate secondary data analysis into their toolbox. Specifically, we define secondary data analysis as a tool and in relation to other common forms of analysis (including exploratory and confirmatory, observational and experimental). We highlight the advantages and disadvantages of this tool. We describe how engagement in transparency can improve and alter our interpretations of results from secondary data analysis and provide resources for robust data analysis. We close by suggesting ways in which subfields and institutions could address and improve the use of secondary data analysis.

Never before have data been so widely available to researchers. Online storage platforms for academic scientists, such as Harvard’s Dataverse and the Open Science Framework, make sharing data across labs, countries, and continents instantaneous. Aimed at a wider audience of data users, government-funded data collection initiatives organize and track individuals at an enormous scale. Coupled with the rise of social media and smartphone technology, behavioral scientists have a wide range of trace data to analyze and combine with a rich array of datasets. Despite this wealth of data, conversations regarding data analysis and modeling in psychology often assume researchers are collecting new data for each analysis. Certainly, a great many principles of this “primary” data analysis are still relevant, applicable, and important in the use of pre-existing data. However, pre-existing data brings with it new concerns—for example, potential bias and lack of experimental control—that warrant careful consideration. Moreover, the benefits of using pre-existing data are often overlooked. In this paper, we describe the analysis of pre-existing data, often called “secondary” data analysis, and outline its value to psychological researchers. We also discuss the potential pitfalls of secondary data analysis, especially in terms of recent advances in open science and transparency. We end with recommendations for increasing the transparency of secondary data analysis and improving the robustness of the results obtained from these methods, including some ideas regarding preregistration. We have written this paper for scientists who are interested in adding secondary data analysis to their research toolbox, and for anyone who wish to use pre-existing data fruitfully and responsibly.

What is secondary data analysis?

We consider *pre-existing data* to be any data that exist before researchers formulate their research hypothesis. Pre-existing data can take many forms. We focus on two in this current manuscript: large-scale survey studies and single-lab studies.

Large-scale survey studies routinely assess a broad number of variables, often from a substantial number of respondents. Typical designs include repeated assessments of the same set of individuals over time, often years or decades in the case of panel or cohort studies, or repeated cross-sectional assessments of representative samples. Such large-scale survey studies are often formed to track changes in the attitudes, health, or economics of a population over time; consequently, they tend to be larger in size, in terms of the number of participants sampled, the number and scope of questions assessed, and the research team. Many panel studies—such as the German Socioeconomic Panel Study (Wagner, Frick & Schupp, 2007) the British Household Panel Study (University of Essex, 2010), and the National Longitudinal Study of Youth 1979 (Bureau of Labor Statistics, 2017)¹—are funded by governments or other large organisations, and have their data made publically available, or available upon registration.

Pre-existing data do not have to be collected on a large scale. When running studies, research labs often choose to collect data that is not directly relevant to the primary research question. Alternatively, after analysis or publication of a study, researchers may think of a different question that the previously-collected data may be able to answer. In both of these cases, we consider these data collected from *smaller-scale lab studies* to be pre-existing data. In this way, the process of generating and sharing data for use by others need not be left to research councils and national governments. Given the potentially limited sample size of these smaller-scale investigations, considerations of statistical power cannot be ignored when

analyzing data of this nature. Single-lab studies may also resemble panel studies, in that researchers may track the same participants over time and repeatedly measure a variety of constructs.

Pre-existing data can take other forms. One of the fastest growing areas of research is “big data,” or data collected through the use of modern technologies including the Internet and smartphones (Hashem et al., 2015; Kosinski, Stillwell & Graepel, 2013). Often these data are collected without a primary research question in mind and may be mined by researchers. We believe the claims regarding and recommendations for using pre-existing data extend to analyses using big data. We consider *secondary data analysis* to be the analysis of any pre-existing data².

Psychologists often think about research in terms of two “modes”: exploratory (i.e. theory-building) and confirmatory (i.e. theory-testing) analyses (Wagenmakers et al., 2012). The former is a common focus for secondary data analysis and is one of its great strengths. Since pre-existing datasets often contain many—even many thousands—of variables, researchers have the flexibility to explore many relationships between constructs. Researchers may run exploratory analyses in pre-existing datasets without wasting valuable time or financial resources. If they find evidence of such a relation, they can choose to invest in another study to confirm it; if they find little evidence, they may decide it would be a waste to collect new data.

On the other hand, it is also possible to use pre-existing data to test theories in a confirmatory fashion. However, this comes with an important caveat: many commonly applied statistical tests were developed under specific assumptions, for example, that data are collected to test one precise hypothesis. Consequently, researchers aiming to use secondary data analysis to provide evidence which might help confirm a theory must take extra steps to ensure the

robustness of their results. We provide a selection of robustness-increasing ideas, which, it should be noted, are not mutually exclusive, in Table 1.

--Insert Table 1 here--

In the same vein, secondary data analysis can be both correlational and experimental in nature. It is true that correlational work makes up the bulk of secondary data analysis, given that much of this work uses panel studies and other survey-type data (see Rohrer, 2017, for discussion of causality in psychological research). However, if we consider the case of single-lab studies, experimental work might also fall under this umbrella. For example, a study designed to assess the effectiveness of an intervention on academic performance might be re-analyzed for effects on additional secondary outcomes, such as happiness or sleep quality, at a later time or by another group of researchers. Quasi-experiments, based on exogenous (often historical) factors that can be harnessed using methods developed in econometrics, also bring an experimental aspect to secondary analysis. For example, the Lanham Act of 1940 provided free, universal child care in the United States during World War II. Using US Census data, researchers were able to track cohort outcomes to estimate the effect of this policy and found a strong and persistent increase in well-being (Herbst, 2017). Methods such as instrumental variables and regression discontinuity analysis can, in cases where their assumptions are met, allow causal inferences from correlational data (Kim & Steiner, 2016), and genetic versions of these techniques, such as Mendelian Randomization, are bringing a new causal aspect to secondary-data studies in biomedicine and beyond (Pingault et al., 2018).

Advantages and disadvantages of secondary data analysis

Pre-existing data, if properly analyzed, offers great advantages: it can help situate effects in real-world behavior and outcomes and in diverse samples—or at least samples more diverse than psychology undergraduate participants (Machery, 2010)—with increased generalizability. It can, in the case of meta-analysis, be used to refine estimates found in prior work. It can be used to investigate hard-to-detect effects thanks to sample sizes that often exceed what is feasible for laboratory studies and allow high-powered statistical tests. It often enables cross-country and cross-cultural research of considerable scope.

Panel studies repeatedly assess participants over years, even decades, allowing for complex longitudinal modeling. Many panel studies are conducted by teams representing a variety of disciplines, including psychology, economics, epidemiology, sociology, and demography; often these datasets contain unique combinations of explanatory and criterion variables. Researchers sampling from these datasets have the opportunity to pair constructs from disparate fields to generate novel research questions. In addition, psychology researchers can benefit from the influence of these other fields. For example, demographers may work to ensure sampling of various geographic locations or subpopulations, allowing for more accurate representations of a country. Panel studies often receive the funding necessary to assess biomarkers of health, giving researchers the data necessary to study small-sized yet potentially meaningful relations between psychosocial and biological variables, such as brain MRI measures or combinations of genetic variants (the largest recent example of a biomedical panel study is the UK Biobank; Collins, 2012).

By including many variables in a single dataset, researchers can make room for creativity and exploration. Collaboration with other lab members or other labs is an excellent learning

opportunity for early career researchers, as they navigate different interests, limited resources, and new technologies. The resulting dataset provides researchers with a resource they can return to again and again for exploration, exercise, and collaboration with others. Even those datasets originally collected for a single study can serve as teaching tools, opportunities to explore an idea, and prototypes for designing new studies. Of course, at all times, researchers ought to take care that (1) the data situation and the analysis remain transparent and that (2) measures are taken to ensure robust inference (see sections below).

On top of all these potential uses, secondary data analysis is an efficient way to conduct research: pre-existing data is often free or costs far less than paying the same number of participants, and researchers spend no time at all collecting data. This makes pre-existing data an especially attractive option for researchers with limited funds or time, including graduate students, post-doctoral fellows, researchers at teaching-oriented universities, and mentors of undergraduate theses. From the perspective of science as an endeavor constrained by limited resources, *not* using pre-existing data when it is available and suitable to answer a research question could be considered inefficient and wasteful.

However, secondary data analysis is not without disadvantages. In cases when applying secondary data analysis to data collected by someone else, a researcher relinquishes control over many important aspects of a study, including the specific research questions they can answer. It may appear obvious, but if the researcher is interested in the relations between A and B, then both A and B must be measured, with a certain degree of internal reliability and external validity. Unfortunately, particularly in the context of large-scale survey studies, this criterion may not always be met. Due to the breadth of such studies, data collectors may opt for short, coarse, and

potentially unreliable measures of those constructs to save time. For example, some of the cognitive tests in the initial sweep of the UK Biobank study, likely due to their being bespoke tests with very short durations, had very poor reliability (Lyall et al., 2016). These issues may be lessened when analyzing one's own data, which is a frequent use of secondary data analysis; however, researchers will still grapple with data that was designed to answer a question different from the specific one they currently study or was not designed with any specific questions in mind. Certain constructs might not have been assessed, or the ordering of an experiment might prohibit the correct temporal analysis. Those interested in longitudinal work might also find that the infrequency of measurement occasions or the length of time between them does not fit their research question. Furthermore, conclusions are necessarily restricted to the populations included in the study. In short, researchers must weigh the convenience and power of pre-existing data with the limitations a dataset imposes on the analysis and research question. Like any other research tool, secondary data analysis is best when used in conjunction with other methods (e.g., primary data collection, including experimentation and survey research; see Munafò & Davey Smith, 2018, for discussion of “triangulation” of research findings across multiple lines of evidence).

Despite its potential, secondary data analysis has been eschewed by some under the notion that it leads to “research parasites”—researchers who do not produce new data but simply live off the data collected by others (cf. Longo & Drazen, 2016). This concern appears symptomatic of misaligned incentives in the field: researchers are not rewarded for collecting high quality data, which could defuse concerns that others “cash in” on one's data collection labor; instead, they are rewarded for presenting striking results. Whereas many reforms are

currently aimed at incentivizing better analyses and transparency (e.g., badges for open practices, Blohowiak et al., 2018; see below), we should consider building incentives for those who collect high-quality data and share it with others. For example, a dataset archived in a public repository such as Dataverse could be equivalent to a publication on a curriculum vitae; and if other researchers use the dataset in a productive manner, this downstream impact should be credited. We might consider developing indicators of quality measurement or repeated measures or large samples, so the evaluation of job or tenure candidates can include attention to these aspects of data collection. Fully acknowledging the collection of high-quality data as an integral contribution to science might require further development for data-sharing norms; publicly-available, high-quality data are of limited use without documentation that enables others to use the data (see Scott & Kline, 2018, for a discussion).

Secondary data analysis through the lens of Open Science

The field of psychology broadly has entered a phase of reflection and reform, largely characterized by an inability to replicate and reproduce many key findings (Pashler & Wagenmakers, 2012). Large-scale collaborative efforts to evaluate the replicability of psychological effects have focused almost exclusively on studies using primary data collection and experimental methods (e.g., Nosek et al., 2015). This is to be expected; replications of such studies are easier to carry out since they typically have smaller sample sizes and more controlled environments than, for example, longitudinal cohort research. Researchers replicating lab-based research can more easily achieve high power and directly copy the testing conditions in the original experiment. We applaud these efforts, which have shed a great deal of light on which

psychological findings can be relied upon and under which circumstances. But a consequence of the focus on experimental studies is uncertainty regarding the replicability of secondary research.

One hint as to the robustness of secondary data analysis is its computational *reproducibility*. While an effect's *replicability* is the extent to which a researcher can find the same effect with different data, its *reproducibility* is the extent to which a researcher can find the same effect with the same data. Reproducibility is a key feature of transparent and robust research, as it results from well-documented analyses. To our knowledge, no one has tried explicitly to estimate the reproducibility of psychological effects found through secondary data analysis. However, this has been attempted in the field of economics, where secondary data analysis is the norm (Chang & Li, 2018). Of more than 60 papers, fewer than half were reproducible, and the study researchers required assistance from the original authors in many cases. Economics journals typically require the submission of code along with a manuscript, a practice that has not yet become mainstream in psychology. This leads us to predict that the reproducibility of psychological findings using secondary data analysis will be lower than that in economics research.

To address issues of reproducibility and replicability, many scientists have advocated for the broad adoption of open science values and practices (e.g., Klein et al., 2014; Nosek et al., 2015), most often implemented through disclosure and transparency in various forms. For example, one of the practical reforms of open science is the implementation of badges. These visual icons are attached to a published article along with links to online resources to signal that open science practices—the current set of badges are for open data, open materials, and preregistration (Kidwell et al., 2016)—have been used in the reported studies. These badges have

been adopted by a number of psychology journals, including *Psychological Science* (Eich, 2014). More generally, psychologists have outlined practices for all members of the scientific community, including researchers, teachers, and journal editors, to adopt in service of increasing the quality of research (Asendorpf et al., 2013; Funder et al., 2014; Lakens & Evers, 2014; van Assen et al., 2014).

Whereas the adoption of open science practices appears to have increased the transparency of psychological science generally (Kidwell et al., 2016), the focus on laboratory-based methodologies has largely neglected the challenges faced by researchers using pre-existing data. Most, if not all, journal badges are inaccessible if pre-existing data are used (or introduce new ethical complications; Finkel, Eastwick & Reis, 2015). The Open Data badge is often unavailable because most panel studies require registering with study coordinators to access data, and data sharing agreements prohibit sharing data among unregistered researchers. The Open Materials badge often cannot be awarded since many studies, especially those initiated decades ago, make use of copyrighted measures that are not permitted to be shared online. Finally, the Preregistration badge hinges upon posting analytic plans before data collection. Even if researchers do not analyze the data prior to registering an analytic plan, they cannot definitively prove (for example, with time-stamped variables) that they have not “peeked” at the data (run a few indicative tests) before making their hypotheses, nor can they prove they have not read other studies which use the data to address similar questions. To some, this prohibits the use of preregistration for secondary data analysis, although, as we make the case below, this need not be true.

The implicit (and sometimes explicit) exclusion of secondary data analysis from open science practices is unfortunate: like all scientific endeavors, secondary data analysis in practice often comes with many pitfalls, and could be further improved if these were addressed. Aside from issues like the lack of experimental control and the resulting restrictions on causal interpretation (see above), secondary data analysis comes with a number of problems familiar to followers of the “replication crisis” (Pashler & Harris, 2012). For instance, given the proliferation of variables in these datasets, it is all too easy to “*p*-hack” one’s way to statistically significant, eye-catching results (Simonsohn, Nelson, & Simmons, 2014). This can be done in a variety of ways: for example, outcome-switching, a practice common in clinical trials (Chan et al., 2004), is also prevalent (in our experience) in secondary analysis; tests of interactions between potential predictors can be added on a whim; and, researchers can simply plug in one covariate after another until a significant result is obtained.

Another common, problematic practice is subgroup analyses. Sometimes this is obvious—for example, examining specific ethnic groups separately—but subgroup analyses can be less conspicuous. For example, researchers may choose to analyze data from a single wave of a longitudinal panel study. In the case of repeated measures, researchers can examine multiple cross-sectional relationships and present only the significant results. For longitudinal data, researchers might, deliberately or otherwise, ignore variables collected at another wave which have important statistical or theoretical links to the constructs of interest.

Certainly these kinds of practices are possible in most studies (Simmons, Nelson, & Simonsohn, 2011), but in large, pre-existing datasets, the temptation to “try it” with another variable or subgroup—selected *post hoc*—is often strong, and the large sample sizes involved

mean that perseverance is likely to be “rewarded” with a p -value falling below the alpha level for significance. In this way, researchers are more likely to present models that fit random variation in their data—especially so as models increase in complexity—instead of revealing reliable, generalizable associations. That is, they are more likely to overfit their models to the data and reduce the potential for replication of their results.

Unique to secondary data analysis is the problem of familiarity with the data. A key reason for using a pre-existing dataset is that it may be the sole source of data appropriate for evaluating a particular research question. For example, questions about lifespan development require decades of time (for example, the unique lifespan data from the Lothian Birth Cohorts described by Deary et al., 2012). Biomarker and genetic data often require a very large team of research assistants, medical professionals, and data scientists (found in large quantities in few places other than the UK Biobank sample; Collins, 2012). Given the limited numbers of datasets available to answer such questions, along with the huge number of variables available in existing ones, it is expected that researchers will return to the same dataset multiple times to investigate different (but similar) research questions. Unfortunately, this practice introduces biases, as researchers become aware of relations in the data. Consequently, researchers can design complex models that fit the data with very few changes or propose very specific hypotheses that are substantiated with few caveats. These are not truly “predictions”, because the researchers already had some knowledge of how the variables relate. As pointed out by Gelman & Loken (2013), the problem is not necessarily the number of ways a researcher analyzes their data, but the number of potential ways they *could* do so. When researchers make analytic decisions based on their data

rather than their theory, the results must be interpreted in the context of multiple potential comparisons.

The proliferation of published research using these datasets means even a researcher who has never worked with a particular dataset before will likely have some knowledge of the patterns within it. Many pre-existing datasets have been repeatedly mined in this way, usually by scholars in the same subfields. For example, the Health and Retirement Study (HRS; Juster & Suzman, 1995) has been used by personality psychologists studying health to examine smoking (Weston & Jackson, 2015), longevity (Hill et al., 2011), and biomarkers of health (Lucetti et al., 2014). It is to be expected that these researchers will read each other's scholarly work, as these studies provide substantial information for generating and testing health-psychological theories. However, in the process of developing well-grounded hypotheses, these researchers also become aware of relationships in the HRS dataset, regardless of whether they had previously analyzed these data, and are potentially biased by what they have learned. This “curse of knowledge” does not preclude researchers from analyzing a dataset they have read about. But prior knowledge of these data sets can bias the choices researchers make in which research questions to ask, how to wrangle variables, and how to fit models.

Because of the opportunity to capitalize on “researcher degrees of freedom” and the increased likelihood of results-biased decision-making, research employing secondary data analysis must be held to as a high—if not a higher—standard as research using primary data collection. In what follows, we make recommendations for increasing the robustness of secondary data analysis on two fronts: by increasing the transparency of secondary data analysis, and by estimating the robustness of results estimated in such research. We make these

suggestions to researchers who value open science and wish to produce research that will stand the tests of time and replication. However, it is our hope that these recommendations will inspire journal editors, grant reviewers, tenure committees and all those who have the power formally to change incentives in the field.

Recommendations for transparent secondary data analysis

Increasing transparency is a cornerstone of the current open science reform movement. The object of attempts to increase transparency is to live up to the ideal summarised in the motto of the UK's Royal Society: *Nullius in verba*, or “take nobody's word for it”. A scientist need not be taken at her word when all her materials, methods, and actions are available for anyone to see. The badges describe above all traffic in transparency: Open Data and Open Materials are the ingredients of a study, and Preregistration clarifies which analytic decisions were determined before the data were analyzed and which were not. This last point is key. If data-analytic decisions are based on the collected data itself, then traditionally used statistical tests can no longer successfully control error rates.

The tendency to make decisions based on data rather than theory becomes more likely, maybe even certain, in the case of pre-existing data, especially if a researcher has used or even read about the data in the past. Take, for example, the proliferation of papers using the Health and Retirement Study (HRS) described above. During a thorough literature review, a personality and health researcher reads frequently about this dataset and become aware that the traits extraversion and conscientiousness are highly correlated in the HRS. She therefore chooses to use the latter as a covariate when examining the relationship of the former to health. This alone is

not problematic; the problem is that readers of her study have no way to know that her decision was based on her prior knowledge. Transparency clarifies for readers of science which decisions were theory-based and independent of the data and which were not, which allows them to interpret results. More specifically, readers (and the researchers conducting the study) should have less confidence in the results from analyses that were designed, in part, by prior knowledge of the data. We recommend several ways in which researchers can transparently document a secondary data analysis:

First, researchers can *provide links to codebooks and data access instructions*. If the pre-existing dataset is a panel study or available for purchase, there are likely to be publically available codebooks or data access websites. These can helpfully supplement the Method section of papers reporting analyses of the dataset as well as STROBE-guideline-based workflows (von Elm et al., 2007). If a researcher owns the dataset, they can create their own codebook with relevant information (e.g., Vardigan, Heus, & Thomas, 2008; for an example, see Condon & Revelle, 2016). If the data are not their own, they can describe how they were able to access them. Panel studies and data banks often email researchers when they have provided access to some or all of the data. Copies of this correspondence, or any data access statements, can be made available to readers as supplementary material. This often contains a date, which can be important time-stamp information if the choice was made to register analyses prior to accessing the data (see below).

Second, researchers should *communicate how the data have been used in prior research*, both by the researchers themselves and by others in the published literature. Enumerating prior experience with a dataset is not meant to prohibit researchers from using that dataset again. The

process of reconstructing this context is next to impossible if it is not done by researchers in an incremental way. Doing so openly simply makes clear both to the readers and the researchers themselves how much prior knowledge went into generating hypotheses or designing models. Regarding research by others, the researcher might simply document the instances they have come across during their literature review. A thorough description of the prior literature is likely central to developing a research hypothesis and writing the Introduction section to an article, so we recommend integrating this description into the eventual literature review. Spending this time prior to the analyses will likely save time when writing up the results. Regarding previous research by that same researcher, providing citations to past publications that are pertinent to the research question is one simple way to reduce the chance that another researcher will inadvertently duplicate previously published work.

A note of caution: It is quite likely that a researcher's analysis history with a particular pre-existing dataset is not limited to simply what has been published. Researchers should disclose any analysis that is relevant to the project. Specifically, this includes any statistic or visualization that includes at least one variable in the project. We believe that this process will become easier as preregistrations and preprints are more widely and consistently used. Ideally, it will become relatively easy to link to prior projects that carefully document both published and unpublished analyses through a platform like the Open Science Framework (OSF). Today, however, this will not be an easy task for most researchers. Because preregistrations and preprints have only recently been adopted in psychological science, this task may actually prove impossible for some. There are no easy solutions for ensuring and checking that researchers have disclosed all knowledge of a dataset. This unfortunately creates space for motivated naivete and

strategic laziness. We must therefore acknowledge that this recommendation—disclosing all prior knowledge—addresses only part of the problem. We hope that work continues on this front.

Third, researchers can *document the data wrangling and analysis pipeline*. Sharing of the analytic script is not always considered part of sharing materials, depending on the journal, but it is especially important for researchers using pre-existing data. A key component of research is documenting by way of code and precise instruction the steps required to access, merge, and prepare the data prior to formal statistical analysis. These procedures are often extremely complex and important details are often left out of academic publications and are dependent on the time and exact version of data as were accessed. As models increase in their complexity, it often becomes more difficult to describe to readers how data were modified and analyzed, especially given the space constraints of many journals' Method sections, for example the economics study described above. Sharing the analysis script instantly deals with this problem.

Fourth, we recommend that *secondary data analysis be preregistered*. Preregistration of secondary data analysis should be similar to preregistration of primary data collection, in that preregistration should occur before the analyses are conducted. Preregistration forms should enumerate any planned analyses and all analytic decisions related to those analyses, for example, the numeric definition of outliers and the procedure for handling them, or how a particular measure will be scored. Researchers can also preregister analyses for upcoming waves of publically available datasets. We note that this system could be expanded to exploratory data analysis as well: the preregistration could simply note a plan to explore relationships between specified variables. At the time of writing, an interactive form is being developed for the preregistration of analysis using pre-existing data. There is also a template OSF project

(osf.io/x4gzt) which guides researchers through the information relevant for preregistering analyses.

Applying the term pre-registration to the analysis of pre-existing, potentially accessible, secondary data is somewhat controversial. Some have argued that the term should be reserved for registration of studies prior to data collection (e.g., Chambers et al., 2015). One of the arguments is that there is no way definitively to prove that a researcher has not looked at the data prior to analysis.

Pre-registrations are very much an imperfect business at present. Many preregistration protocols are too vague to safeguard against *p*-hacking (Wicherts et al., 2016); the published manuscript might not follow the pre-specified analysis plan, nor is it clear how or whether journals should evaluate fidelity to an analysis plan (Tucker, 2014). Reviewers and readers must still compare the preregistration to a final study to evaluate adherence, and adherence itself tells us little about other important aspects of study quality (e.g., the validity of the design). Hence, one could argue that a preregistration per se does not imply much, which is why the label should not be interpreted as a signal of superior quality. We should note that preregistration is not a box-ticking exercise. In evaluating a manuscript, attention should be paid not just to whether a preregistration exists, but to the content and quality of that pre-registration.

Given that preregistration seems to have acquired such a special meaning, and given its prominent role in the “Open Science Trifecta” (Open Data, Open Materials, Preregistration), reserving the use of the term only for primary data analysis might make researchers who rely on secondary data assume incorrectly that open science is not of relevance or availability to them. A simplistic “preregistration or it didn’t happen” mindset might even lead researchers to conclude

that secondary data analysis is second class research because it cannot be fully preregistered, furthering the chasm between different research traditions.

Hence, we argue that the term preregistration *can* be applied to secondary data analysis, mostly for pragmatic reasons, and at the same time, we would like to encourage more discussion about what preregistrations can and cannot achieve, both in the context of primary and secondary data analysis. For example, preregistrations are *always* trust-based regardless of whether the data already existed, since it is possible to “pre”register a study that has already been conducted and analyzed; and since there is currently no mechanism in place that prevents researchers from filing multiple pre-registrations (potentially on different platforms) with slightly different analysis plans and later selectively reporting the one that “worked.” Scientists who yearn for a bulletproof approach that cannot be gamed by insincere authors might prefer Registered Reports, which (1) make it very hard to “pre”-register data analyses that have already been performed because reviewers’ feedback during the initial stage can lead to substantial changes in these analyses and (2) partially remove the incentive to produce a certain result thanks to in-principle-acceptance prior to data collection/analysis. For the former reason, Registered Reports, unlike weaker preregistration of analysis plans, might preclude secondary data analysis when researchers cannot supply evidence that they had no prior access to the data, although this too can be a point of discussion.

Improved inference based on secondary data analysis

Researchers often face a large number of decisions while analyzing their data (e.g., whether and how to transform variables, which covariates to include, which estimator to use),

and they might often genuinely be unsure about the best way statistically to approach their research question. This is even more of an issue with data sets that are rich in variables. Thus, in the context of secondary data analysis, the robustness of one's finding becomes a central concern: would conclusions substantially change if a different plausible model specification were used?³

Based on empirical testing, Young and Holsteen (2017) described three different patterns of model robustness: the result holds no matter how the model is specified (i.e., the finding is robust); the result depends on some specific model ingredients such as a particular covariate (i.e., there is systematic variability); or the result depends on a very specific combination of parameters and only arise in one (or a few) of many possible models ("knife edge" specification). A robust finding would increase confidence that conclusions are not based on a fluke. Systematic variability would call for follow-up analyses to better understand the role of the particular model ingredient—for example, perhaps the inclusion of a covariate introduces (or removes) a spurious contamination. Knife edge specifications call for prudence: if only one in a multitude of plausible models supports a particular finding, it could plausibly just be a fluke in the data.

The simplest way to probe the robustness of one's finding is to perform so-called robustness checks (also known as sensitivity analyses), which are a staple in economics research but appear to be less common in psychology (see e.g. Duncan, Engel, Claessens, & Dowsett, 2014, for a comparison of journals in economics versus those in developmental psychology)⁴. In the most standard form, the model is re-run with one element of the specification changed. Robustness checks might range from a simple footnote (e.g., "results remained unaffected when

age was included as a covariate”) to a supplementary website presenting results from dozens of models in such a way that they can be compared by the reader (e.g., Arslan et al., 2017).

The fundamental idea of robustness checks can be expanded by considering all possible combinations of all plausible model ingredients. Naturally, this quickly leads to rapid growth of the number of possible models: for example, even if there are only three simple dichotomous decisions to be made (control for gender yes/no; control for age yes/no; remove outliers yes/no), already 8 different model specifications ($2 \times 2 \times 2$) result. Several researchers have advocated that all these models should be run and reported. For example, Steegen, Tuerlinckx, Gelman, and Vanpaemel (2016) label this approach a “multiverse analysis”; Young (2018) described this as the “computational solution” to model uncertainty; and Simonsohn, Simmons, and Nelson (2015) developed “specification curves”, which allow researchers to calculate a *p*-value across all specifications.

Though originally developed for small-scale laboratory data there have been at least two successful applications of specification curve analysis in the studies of birth-order effects on personality (Rohrer, Egloff & Schmukle, 2017) and impact of digital technology use on psychological well-being (Orben & Przybylski, 2019). Both of these investigations built on large-scale social datasets. Results derived from the work focused on birth-order effects suggested that some published results in this literature might depend on knife-edge specifications. Considering the effects of gaming and social media, the thorough analytic approach applied to large-scale datasets (i) provided a robust estimate of the modest impact on young people and (ii) put the effect sizes within clear everyday contexts using the inherent

richness of the available data (for example, by comparing the associations between technology and well-being with the associations between potato consumption and well-being).

Robustness checks and their expansions are still no guarantee that the data have not been overfit. Among economists, one sometimes hears jokes about how wondrous it is that robustness checks always only work to confirm the finding; the danger of selectively only including model ingredients that support one's preferred conclusion is certainly higher than zero. Hence, to further strengthen these approaches, they can also be preregistered.

Beyond robustness checks, there are additional approaches that can be used to ensure that biases do not affect secondary data analyses and to avoid overfitting. We have included some of these approaches for reference in Table 1. We note that these recommendations are not specifically for secondary data analysis and are used with great success in analyzing primary data sources.

Into the future

We have recommended methods to ensure that secondary data analysis is transparent and that results from secondary data analysis are robust. We finish with three calls to action.

First, for researchers who run laboratory experiments with the potential for further analyses, we encourage them to consider making their datasets available for others to analyse in a secondary-analysis fashion. Such datasets—whether made completely open or accessed with permission—constitute valuable resources for future research. As discussed above, we believe the production and curation of such datasets should be considered a research output with value akin to the writing of papers or the development of statistical software packages.

Second, we turn to fields where secondary data analysis is frequently used, such as the subfields of personality, individual differences, and development. The use of secondary data analysis, specifically the use of a few large surveys, can create the illusion of replication or convergence across an area of research. We say “illusion,” because a large proportion of published results within a field may be found using the same panel study or dataset. This results in a relatively large number of publications reporting similar effects, growing a literature of evidence for a general relationship or idea, but failing to expand the number of truly independent tests. For example, the German Socio-economic Panel Study (SOEP) has repeatedly been used to track personality development, including twice in the same issue of the same journal (Lucas & Donnellan, 2011; Specht, Egloff & Schmukle, 2011). This is not necessarily problematic—and could even be beneficial to probe the robustness of different analytical approaches—if it is clear to readers that many studies use the same data for similar or sometimes identical questions. However, without better indexing for data (e.g., clear tags referring to the data source), the extent of use of a particular dataset is difficult for readers to evaluate.

For example, how much of the evidence for the link between trait conscientiousness and health is based on data from the Health and Retirement Study? Dependence between published findings limits our certainty in an effect. If a published literature is largely supported by one dataset, or even a small number of datasets, we should be less certain that the effect in question is generalizable to people outside of that sample. More importantly, multiple related findings from a single dataset may not suggest multiple independent effects, but rather one effect with shared variance across a number of indicators. In the conscientiousness-health example, conscientiousness has been found to be associated with mortality (Hill et al., 2011), incidence of

chronic conditions (Weston, Hill & Jackson, 2015), health behaviors (Hakulinen et al., 2015; Roberts, Jackson & Edmonds, 2009), and sleep (Hintsanen et al., 2014), but each of these studies uses the HRS dataset. We cannot count each of these publications as reporting an independent result: they all rely on the same sample of individuals as well as on the same operationalization of conscientiousness. Without independent verification of each of these associations in new datasets, this evidence merely suggests that conscientiousness, as measured in this particular study, may be related to some (likely overlapping) aspects of health in this sample. We call for introspection and systematic review of key findings in our subfields, especially for those findings about which we feel certainty. Do we find these effects in multiple, independent datasets? Or do we find them in only one or two datasets, over and over again?

We also call for systematic reproducibility checks on studies using secondary data analysis. How many of our findings have been preregistered, or make code available, or in any way ensure others can readily reproduce the effects? We especially appeal to academic journals, which could serve the field by hiring statistical editors or reviewers whose jobs are to reproduce analyses and results based on the code and data provided or specified. At the time of writing, this has already been implemented by six academic journals, including *Meta-Psychology*, and so a viable model exists (see the Data Reproducibility Policies section of the wiki page at osf.io/kgnvva).

Finally, we turn to those interested in developing technologies for the advancement of open science and call for the development of tools specific to secondary data analysis. Certainly there are ways to adapt existing tools for our work (e.g., preregistering secondary data analysis). But there are specific challenges in the case of pre-existing data that can be addressed. Others

have proposed data “check-out” systems as a form of preregistration and monitoring of data use (Scott & Kline, 2018). We also suggest the development of tools for tracking and reporting prior knowledge of a dataset.

Conclusion

Our purpose here was to urge psychologists to consider the use of secondary data analysis as a power and economic tool for exploring important research questions. We urge those with limited resources especially to consider the ways in which pre-existing data may supplement or build research and teaching programs. We note that secondary data has many strengths, but it also falls prey to several of the limitations and questionable research practices that continue to haunt psychology during its “replication crisis”. Happily, however, many of the specific reforms that have begun to improve the credibility of primary research either have analogies to, or can be directly implemented in, secondary analysis. Preregistration can be implemented in terms of registering analysis before data are accessed. Inferences can be improved using strategies such as cross-validation, hold-out samples, and multi-cohort analyses. The robustness of results can be tested to destruction in a multiverse analysis. These and the other techniques we have outlined in this paper form a manifesto for the improvement of secondary data analysis, to ensure that this critically-important type of research is carried along with the Open Science revolution.

Author contributions

All authors contributed to the writing and editing of this manuscript.

Acknowledgements

None of this work would have been possible without the contributions of those who participated in a discussion of open science strategies for secondary data analysis, which took place at the 2nd Annual Meeting of the Society for the Improvement of Psychological Science (SIPS) in Charlottesville, VA in 2017. Those participants were Cassandra Brandes, Bill Chopik, Lisa Hoplock, Patrick Forscher, Nick Fox, Rich Lucas, Kathryn Mills, John Protzko, Kathleen Reardon, Kate Rogers and Kaitlyn Werner.

Since the initial drafting of this manuscript, another meeting of SIPS took place in Grand Rapids, MI in 2018. During this meeting, another group of researchers built on the principles and recommendations of this manuscript and developed a template for registering secondary data analyses. Those researchers included Olmo van den Akker, Marjan Bakkar, Brian Brown, Lorne Campbell, Bill Chopik, Oliver Clark, Rodica Damien, Pam Davis-Kean, Charlie Ebersole, Andrew Hall, Matthew Kay, Jessica Kosie, Elliot Kruse, Jerome Olsen, Stuart Ritchie, Courtney Soderberg, K.D. Valentine, Anna Van't Veer, and Sara Weston.

References

- Arslan, R. C. (2017, September 14). Overfitting vs. open data. [Blog post]. Retrieved from <http://www.the100.ci/2017/09/14/overfitting-vs-open-data/>
- Arslan, R. C., Willführ, K. P., Frans, E., Verweij, K. J. H., Bürkner, P., Myrskylä, M., ... Penke, L. (2017, August 4). Paternal age and offspring fitness: online supplementary website (Version v2.0.1). Zenodo. <http://doi.org/10.5281/zenodo.838961>
- Benjamin, D. J., Berger, J. O., Johannesson, M., Nosek, B. A., Wagenmakers, E. J., Berk, R., ... & Cesarini, D. (2018). Redefine statistical significance. *Nature Human Behaviour*, 2(1), 6.
- Blohowiak, B. B., Cohoon, J., de-Wit, L., Eich, E., Farach, F. J., Hasselman, F., ... DeHaven, A. C. (2018, November 14). Badges to Acknowledge Open Practices. Retrieved from osf.io/tvyxz
- Chang, A. C., & Li, P. (2018). Is Economics Research Replicable? Sixty Published Papers from Thirteen Journals Say “Often Not”. *Critical Finance Review*, 7.
- Collins, R. (2012). What makes UK Biobank special?. *The Lancet*, 379(9822), 1173-1174.
- Credé, M., & Phillips, L. A. (2017). Revisiting the Power Pose Effect: How Robust Are the Results Reported by Carney, Cuddy, and Yap (2010) to Data Analytic Decisions?. *Social Psychological and Personality Science*, 8(5), 493-499.
- Hakulinen, C., Hintsanen, M., Munafò, M. R., Virtanen, M., Kivimäki, M., Batty, G. D., & Jokela, M. (2015). Personality and smoking: Individual-participant meta-analysis of nine cohort studies. *Addiction*, 110(11), 1844-1852.

- Hashem, I. A. T., Yaqoob, I., Anuar, N. B., Mokhtar, S., Gani, A., & Khan, S. U. (2015). The rise of “big data” on cloud computing: Review and open research issues. *Information systems*, 47, 98-115.
- Herbst, C. M. (2017). Universal child care, maternal employment, and children’s long-run outcomes: Evidence from the US Lanham Act of 1940. *Journal of Labor Economics*, 35(2), 519-564.
- Hill, P. L., Turiano, N. A., Hurd, M. D., Mroczek, D. K., & Roberts, B. W. (2011). Conscientiousness and longevity: an examination of possible mediators. *Health Psychology*, 30(5), 536.
- Hintsanen, M., Puttonen, S., Smith, K., Törnroos, M., Jokela, M., Pulkki-Råback, L., ... & Venn, A. (2014). Five-factor personality traits and sleep: Evidence from two population-based cohort studies. *Health Psychology*, 33(10), 1214.
- Hofer S. M. Piccinin A. M . 2009. Integrative data analysis through coordination of measurement and analysis protocol across independent longitudinal studies. *Psychological Methods* , 14, 150–164. doi:10.1037/a0015566
- Kim, Y., & Steiner, P. (2016). Quasi-experimental designs for causal inference. *Educational Psychology*, 51(3-4), 395-405. doi: 10.1080/00461520.2016.1207177
- Klein, R. A., Ratliff, K. A., Vianello, M., Adams, R. B., Jr., Bahník, Š., Bernstein, M. J., ... Nosek, B. A. (2014). Investigating variation in replicability: A “many labs” replication project. *Social Psychology*, 45(3), 142-152.
- <http://dx.doi.org/10.1027/1864-9335/a000178>

- Kosinski, M., Stillwell, D., & Graepel, T. (2013). Private traits and attributes are predictable from digital records of human behavior. *Proceedings of the National Academy of Sciences*, 201218772.
- Lakens, D. (2018, December 1). Justify your alpha by decreasing alpha levels as a function of the sample size. [Blog post]. Retrived from <http://daniellakens.blogspot.com/2018/12/testing-whether-observed-data-should.html/>
- Longo, D. L., & Drazen, J. M. (2016). Data sharing. DOI: 10.1056/NEJMe1516564
- Lucas, R. E., & Donnellan, M. B. (2011). Personality development across the life span: Longitudinal analyses with a national sample from Germany. *Journal of personality and social psychology*, 101(4), 847.
- Lyall, D. M., Cullen, B., Allerhand, M., Smith, D. J., Mackay, D., Evans, J., ... & Pell, J. P. (2016). Cognitive test scores in UK Biobank: data reduction in 480,416 participants and longitudinal stability in 20,346 participants. *PLOS ONE*, 11(4), e0154222.
- MacCoun, R., & Perlmutter, S. (2015). Blind analysis: hide results to seek the truth. *Nature News*, 526(7572), 187.
- Munafò, M. R., & Davey Smith, G. (2018). Robust research needs many lines of evidence. *Nature*, 553(7689), 399-401.
- Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D., Breckler, S. J., et al. (2015). Promoting an open research culture. *Science*, 348(6242), 1422–1425. <http://doi.org/10.1126/science.aab2374>
- Orben, A., & Przybylski, A. K. (2019). The association between adolescent well-being and digital technology use. *Nature Human Behaviour*, 1.

- Pashler, H., & Wagenmakers, E. J. (2012). Editors' introduction to the special section on replicability in psychological science: A crisis of confidence? *Perspectives on Psychological Science*, 7(6), 528-530. doi: 10.1177/1745691612465253
- Pingault, J. B., O'Reilly, P. F., Schoeler, T., Ploubidis, G. B., Rijsdijk, F., & Dudbridge, F. (2018). Using genetic data to strengthen causal inference in observational research. *Nature Reviews Genetics*, 19, 566-580. doi: 10.1038/s41576-018-0020-3
- Roberts, B. W., Smith, J., Jackson, J. J., & Edmonds, G. (2009). Compensatory conscientiousness and health in older couples. *Psychological Science*, 20(5), 553-559.
- Rohrer, J. M., Egloff, B., & Schmukle, S. C. (2017). Probing birth-order effects on narrow traits using specification-curve analysis. *Psychological Science*, 28(12), 1821-1832.
- Scott, K. M., & Kline, M. (2018, July 30). Enabling Confirmatory Secondary Data Analysis by Logging Data 'Checkout'. <https://doi.org/10.31234/osf.io/87wjc>
- Simonsohn, U., Simmons, J. P., & Nelson, L. D. (2015). Specification curve: Descriptive and inferential statistics on all reasonable specifications.
- Specht, J., Egloff, B., & Schmukle, S. C. (2011). Stability and change of personality across the life course: The impact of age and major life events on mean-level and rank-order stability of the Big Five. *Journal of personality and social psychology*, 101(4), 862.
- Steege, S., Tuerlinckx, F., Gelman, A., & Vanpaemel, W. (2016). Increasing transparency through a multiverse analysis. *Perspectives on Psychological Science*, 11(5), 702-712.
- Tucker, J. (2014, September, 18). Experiments, preregistration, and journals. [Blog post]. Retrieved from <https://blog.oup.com/2014/09/pro-con-research-preregistration/>

- van't Veer, A. E., & Giner-Sorolla, R. (2016). Pre-registration in social psychology—A discussion and suggested template. *Journal of Experimental Social Psychology*, 67, 2-12.
- Varma, S., & Simon, R. (2006). Bias in error estimation when using cross-validation for model selection. *BMC bioinformatics*, 7(1), 91.
- Vartanian, T. P. (2010). Secondary data analysis. Oxford University Press.
- Wagenmakers, E. J., Wetzels, R., Borsboom, D., van der Maas, H. L., & Kievit, R. A. (2012). An agenda for purely confirmatory research. *Perspectives on Psychological Science*, 7(6), 632-638.
- Weston, S. J., Hill, P. L., & Jackson, J. J. (2015). Personality traits predict the onset of disease. *Social Psychological and Personality Science*, 6(3), 309-317.
- Wicherts, J. M., Veldkamp, C. L., Augusteyn, H. E., Bakker, M., Van Aert, R., & Van Assen, M. A. (2016). Degrees of freedom in planning, running, analyzing, and reporting psychological studies: A checklist to avoid p-hacking. *Frontiers in Psychology*, 7, 1832.
- Yarkoni, T., & Westfall, J. (2017). Choosing prediction over explanation in psychology: Lessons from machine learning. *Perspectives on Psychological Science*, 12(6), 1100-1122.
- Young, C., & Holsteen, K. (2017). Model uncertainty and robustness: a computational framework for multimodel analysis. *Sociological Methods & Research*, 46(1), 3-40.
- Young, C. (2018). Model Uncertainty and the Crisis in Science. *Socius*, 4, 2378023117737206.

Footnotes

¹ We recommend the reader browse the “Cohort Profile” section in each issue of the International Journal of Epidemiology for details on a huge number of other such datasets.

² The term “secondary data” is sometimes used in a somewhat different context: to refer to data that is collected by one researcher (or team of researchers) and analyzed by a second (or team; e.g., Vartanian, 2010). We choose not to use this definition, since pre-existing data can be collected by the same team that wishes to analyze it. However, we retain the use of term “secondary data analysis” to connect our work with others who have sought to improve the robustness of research using secondary data and have curated lists of available data sets.

³ Of course, robustness can also be a central concern in primary data analysis, as illustrated in Credé & Phillips (2017).

⁴ It should be noted that these authors use a slightly different definition of robustness check that includes the replication of a finding on a new data set.

Table 1. Approaches for improving inference based on (secondary data) analysis.

Method	Description
Data Blind Analysis	To avoid their preconceptions affecting their data analyses, particle physicists and cosmologists use blind analysis (MacCoun & Perlmutter, 2015): aspects of the data are altered (e.g., random noise is added to data points, variable labels are shuffled), all analytical decisions are made on this altered data set, and finally the analysis is run on the “real”, original data. Such an approach could also be used by psychologists analyzing secondary (though also, for that matter, primary) data.
Cross-validation	In the context of machine learning, cross-validation is a standard approach to avoid the statistical model being overfit to the data at hand. The data set is repeatedly split into training and test subsets; the training data serves to estimate the model parameters, whereas the test data is used to evaluate the performance of the model (see Yarkoni & Westfall, 2017, for an introduction). If there are additional analytic flexibilities in model specification (e.g., decisions about which variables to include), this can be expanded to nested resampling (Varma & Simon, 2006), in which analytic decisions are based on a separate part of the data, and the model is

	then estimated and evaluated (using cross-validation) on the remaining part of the data.
Hold-out data	<p>The very nature of secondary data opens the door to one highly effective mechanism to avoid overfitting: data curators could hold back parts of the data. Researchers could use the data available to them to specify and estimate their models, and the holdout data, provided after the completion of this initial analysis, could be used to give an unbiased estimate of model performance (suggested by Arslan, 2017). For example, in the Fragile Families Challenge, researchers received access to parts of a longitudinal dataset with more than 10,000 variables and were challenged to predict parts of the data they had not seen (http://www.fragilefamilieschallenge.org/). To our knowledge, no major data holder or curator has yet implemented systematic holdouts, but these might be a promising future avenue.</p>
Adjusted alpha level	<p>Another approach to limit false-positive findings is setting the alpha-level. For example, researchers might want to use a more conservative level of .005 instead of .05 (Benjamin et al., 2018), or decrease their alpha as a function of sample size to balance error rates (Lakens, 2018). Note that this suggestion, as with many of the others, is by no means limited to secondary data analysis.</p>

Coordinated analysis	<p>The existence of multiple, independent, large-scale survey studies also allows for evaluation of generalizability in the context of secondary data analyses. In this kind of multi-cohort “coordinated” analysis (as suggested by Hofer & Piccinin, 2009), researchers can test the same (or similar) analytic models in different samples, representing, for example, different geographic locations or cohorts, or different measurement instruments. Results can be pooled to better estimate an effect size and evaluate heterogeneity across differences in populations and methods.</p>
Exploratory data analysis	<p>All the above recommendations have been applicable to confirmatory data analysis, but it is also important to consider exploratory methods. It has been argued that a major flaw of the way research is currently reported is that exploratory research is often written up as if it were confirmatory all along (Wagenmakers et al., 2012). Clearly identifying exploratory analyses in a manuscript helps readers better assess the robustness of a particular result, and opens the door for high-quality confirmatory follow-up research. We recommend that researchers omit <i>p</i>-values and other tests of significance from exploratory analyses, as these cannot properly be interpreted without a confirmatory framework.</p>

