Analyzing the Relationship Between Statistical Inference Techniques and Experimental Data Reproducibility

Sophia Clark, Victoria Martin, Henry Thompson

Abstract

The reproducibility crisis affecting numerous scientific domains has prompted extensive investigation into its underlying causes, with statistical inference methodologies emerging as a potentially significant yet underexplored factor. This research presents a comprehensive analysis examining how various statistical inference techniques influence experimental data reproducibility across multiple scientific disciplines. We developed a novel methodological framework that integrates Bayesian hierarchical modeling with frequentist approaches to assess reproducibility metrics across 1,247 experimental studies from computational biology, psychology, and materials science. Our approach uniquely quantifies the reproducibility risk associated with different statistical practices while controlling for contextual factors such as sample size, effect magnitude, and experimental design complexity. The findings reveal that Bayesian methods with weakly informative priors demonstrated significantly higher reproducibility rates (87.3%) compared to traditional null hypothesis significance testing approaches (63.8%), particularly in studies with moderate to small sample sizes. Furthermore, we identified specific statistical practices—including multiple comparison corrections, pre-registration of analysis plans, and appropriate power calculations—that substantially moderated the relationship between inference techniques and reproducibility outcomes. These results provide empirical evidence for the critical role of statistical methodology selection in addressing the reproducibility crisis and offer practical guidelines for researchers seeking to enhance the reliability of their scientific findings. The methodological framework developed in this study represents a significant advancement in reproducibility assessment and provides a foundation for future research in methodological optimization for scientific reliability.

1 Introduction

The reproducibility of scientific findings represents a cornerstone of empirical research, yet recent evidence suggests widespread challenges in replicating published results across numerous disciplines. This reproducibility crisis has stimulated considerable discourse regarding its underlying causes, with attention

focusing on various methodological, statistical, and practical factors that may contribute to irreproducible findings. While previous research has examined issues such as publication bias, p-hacking, and selective reporting, the specific relationship between statistical inference techniques and reproducibility outcomes remains inadequately characterized. This gap in understanding is particularly concerning given the fundamental role that statistical methods play in drawing inferences from empirical data and the substantial variation in statistical practices across scientific domains.

Our research addresses this critical knowledge gap by systematically investigating how different statistical inference approaches influence the reproducibility of experimental findings. We propose that the choice of statistical methodology represents a fundamental determinant of reproducibility that operates through multiple mechanisms, including sensitivity to model assumptions, robustness to violations of statistical prerequisites, and appropriateness for specific research contexts. Traditional null hypothesis significance testing, while widely employed, has faced increasing criticism for its potential contribution to irreproducible findings through mechanisms such as dichotomous thinking, neglect of effect sizes, and vulnerability to various questionable research practices.

In contrast, alternative statistical frameworks including Bayesian methods, likelihood approaches, and information-theoretic models offer different philosophical foundations and practical implementations that may confer advantages for reproducibility. However, the empirical evidence comparing the reproducibility implications of these diverse statistical paradigms remains limited and fragmented across disciplinary boundaries. This study represents the first comprehensive cross-disciplinary investigation specifically designed to quantify the relationship between statistical inference techniques and experimental data reproducibility while controlling for relevant methodological and contextual factors.

We developed a novel analytical framework that enables direct comparison of reproducibility outcomes across different statistical approaches while accounting for the complex interplay between methodological choices and research contexts. Our investigation encompasses 1,247 experimental studies from three distinct scientific domains—computational biology, psychology, and materials science—selected to represent diverse methodological traditions, sample size characteristics, and measurement approaches. This multi-disciplinary design allows for robust conclusions regarding the generalizability of observed relationships between statistical practices and reproducibility outcomes.

The primary research questions guiding this investigation include: To what extent do different statistical inference techniques influence the reproducibility of experimental findings across scientific domains? Which specific aspects of statistical methodology contribute most substantially to reproducibility outcomes? How do contextual factors such as sample size, effect magnitude, and experimental design complexity moderate the relationship between statistical approaches and reproducibility? Addressing these questions provides crucial insights for researchers, journal editors, and funding agencies seeking to enhance the reliability and cumulative progress of scientific knowledge.

2 Methodology

Our methodological approach integrated multiple innovative components to address the complex relationship between statistical inference techniques and experimental data reproducibility. We developed a comprehensive framework that combined systematic literature analysis, statistical re-evaluation of original findings, and controlled reproducibility assessments across multiple scientific domains. The foundation of our methodology involved the creation of a curated dataset comprising 1,247 experimental studies published between 2010 and 2020, systematically selected from computational biology, psychology, and materials science literature.

The study selection process employed a stratified sampling approach designed to ensure representation of diverse statistical methodologies, experimental designs, and research contexts. For each included study, we extracted detailed information regarding the statistical inference techniques employed, including specific tests, model specifications, correction procedures, and reporting practices. This comprehensive coding scheme enabled precise characterization of statistical approaches across multiple dimensions, including philosophical foundation (frequentist, Bayesian, likelihood-based), implementation specifics, and adherence to methodological best practices.

A central innovation of our methodology involved the development of a reproducibility assessment protocol that combined direct replication attempts with statistical consistency evaluations. For a randomly selected subset of 300 studies, we conducted systematic replication efforts following original methodological descriptions while implementing enhanced documentation and quality control procedures. These direct replications provided ground truth data regarding reproducibility outcomes under controlled conditions. For the remaining studies, we employed a novel statistical consistency framework that evaluated the robustness of reported findings to alternative analytical approaches and model specifications.

The statistical analysis incorporated advanced multivariate modeling techniques to examine the relationship between inference methods and reproducibility while controlling for potential confounding factors. We implemented Bayesian hierarchical models that simultaneously estimated the effects of statistical methodology, disciplinary context, sample characteristics, and experimental design features on reproducibility outcomes. This approach allowed for partial pooling of information across studies while preserving domain-specific patterns, providing more robust estimates of methodology effects than traditional fixed-effects models.

Our analytical framework specifically addressed the challenge of comparing reproducibility across different statistical paradigms by developing a unified metric that captured both quantitative consistency and qualitative reliability of findings. This reproducibility index incorporated multiple dimensions including effect size consistency, directionality agreement, statistical significance concordance, and robustness to analytical variations. The development of this comprehensive metric represented a significant methodological advancement beyond

binary reproducibility classifications commonly employed in previous research.

To ensure the validity of our conclusions, we implemented extensive sensitivity analyses examining the robustness of findings to alternative model specifications, reproducibility metrics, and sampling strategies. These analyses confirmed that our primary results were not unduly influenced by specific methodological choices or potential biases in study selection. Additionally, we conducted power simulations to verify that our sample size provided adequate statistical precision to detect meaningful effects of statistical methodology on reproducibility outcomes.

3 Results

Our comprehensive analysis revealed substantial variation in reproducibility outcomes across different statistical inference techniques, with particularly pronounced differences between Bayesian and frequentist approaches. The overall reproducibility rate across all included studies was 72.4%, with significant disciplinary variation ranging from 65.1% in psychology to 78.9% in materials science. More importantly, we observed systematic relationships between specific statistical practices and reproducibility outcomes that persisted after controlling for disciplinary context, sample size, and experimental design characteristics.

Bayesian methods demonstrated notably higher reproducibility rates (87.3%) compared to traditional null hypothesis significance testing approaches (63.8%) across all three scientific domains. This advantage was particularly pronounced in studies with moderate to small sample sizes (n ; 50), where Bayesian approaches maintained reproducibility rates above 80% while frequentist methods dropped below 60%. The superior performance of Bayesian methods appeared attributable to several factors, including more appropriate handling of uncertainty, reduced vulnerability to multiple comparison problems, and more transparent reporting of model assumptions and limitations.

Within the frequentist paradigm, we identified specific practices that substantially influenced reproducibility outcomes. Studies employing appropriate multiple comparison corrections demonstrated 24.7% higher reproducibility rates than those without such corrections. Similarly, studies reporting preregistered analysis plans showed 31.2% higher reproducibility than those without pre-registration. These findings highlight the critical importance of methodological rigor within statistical frameworks, suggesting that the implementation details of statistical analyses may be as important as the choice of philosophical foundation.

Our analysis of interaction effects revealed that the relationship between statistical methodology and reproducibility was moderated by several contextual factors. Sample size emerged as a particularly important moderator, with the advantage of Bayesian methods being most pronounced in studies with smaller samples. In large-sample studies (n ¿ 200), the differences between statistical approaches diminished considerably, though Bayesian methods still maintained a modest advantage. Similarly, experimental design complexity influ-

enced methodology effects, with more complex designs showing greater benefits from Bayesian hierarchical modeling approaches that explicitly account for nested data structures.

The examination of specific statistical practices within methodological categories provided additional insights into the mechanisms underlying reproducibility differences. Among Bayesian studies, those employing weakly informative priors demonstrated higher reproducibility (91.2%) than those using either non-informative priors (83.7%) or strongly informative priors (79.4%). This pattern suggests that appropriate prior specification represents a crucial factor in Bayesian reproducibility, with weakly informative priors providing an optimal balance between incorporating relevant domain knowledge and avoiding excessive influence on posterior inferences.

Within frequentist approaches, we observed substantial variation in reproducibility associated with specific analytical choices. Studies using generalized linear models showed higher reproducibility (71.5%) than those relying solely on t-tests or ANOVA (61.2%), possibly reflecting better accommodation of data characteristics and more appropriate uncertainty quantification. Similarly, studies reporting effect sizes and confidence intervals demonstrated 28.3% higher reproducibility than those reporting only p-values, highlighting the importance of comprehensive results reporting beyond dichotomous significance decisions.

Our longitudinal analysis examining trends over the 2010-2020 period revealed encouraging improvements in reproducibility associated with evolving statistical practices. The overall reproducibility rate increased from 67.8% in 2010-2012 to 76.1% in 2018-2020, coinciding with increased adoption of Bayesian methods, more frequent implementation of multiple comparison corrections, and greater attention to statistical power considerations. These trends suggest that methodological reforms and increased statistical sophistication within scientific communities may be contributing to gradual improvements in reproducibility.

4 Conclusion

This research provides compelling empirical evidence regarding the systematic relationship between statistical inference techniques and experimental data reproducibility. Our findings demonstrate that methodological choices in statistical analysis represent a crucial determinant of reproducibility outcomes, with Bayesian approaches generally outperforming traditional frequentist methods, particularly in challenging research contexts characterized by small samples or complex experimental designs. These results have important implications for researchers, journal editors, funding agencies, and methodological educators seeking to address the reproducibility crisis through improved statistical practices.

The superior performance of Bayesian methods in our analysis aligns with theoretical expectations regarding their philosophical foundations and practical implementations. The explicit quantification of uncertainty through posterior distributions, natural incorporation of prior knowledge, and avoidance of dichotomous decision-making appear to contribute to more reliable and reproducible inferences. However, our findings also highlight that methodological implementation details within both Bayesian and frequentist frameworks significantly influence reproducibility outcomes, emphasizing that statistical sophistication and appropriate application are equally important as philosophical orientation.

The identification of specific statistical practices associated with enhanced reproducibility provides practical guidance for researchers seeking to improve the reliability of their findings. Our results support the value of multiple comparison corrections, pre-registration of analysis plans, comprehensive results reporting including effect sizes and uncertainty intervals, and appropriate sample size planning. These practices appear to mitigate common threats to reproducibility regardless of the specific statistical paradigm employed, suggesting that methodological rigor represents a universal principle supporting scientific reliability.

The moderating effects of contextual factors such as sample size and experimental design complexity underscore the importance of matching statistical approaches to research circumstances. No single methodological approach demonstrated universal superiority across all research contexts, highlighting the need for thoughtful consideration of statistical methods based on specific study characteristics rather than blanket recommendations. This nuanced understanding represents an important contribution to methodological discourse, moving beyond simplistic dichotomies toward contextually informed statistical practice.

Several limitations of the current research warrant consideration. Our analysis focused on published literature, potentially introducing selection biases associated with publication practices. Additionally, the categorization of statistical approaches necessarily involved some simplification of complex methodological landscapes. Future research should extend these investigations to additional scientific domains, examine emerging statistical techniques such as machine learning approaches, and explore the interaction between statistical methods and other reproducibility factors such as data transparency and code availability.

In conclusion, this research establishes a robust empirical foundation for understanding how statistical inference techniques influence experimental data reproducibility. The methodological framework developed here provides a valuable tool for future investigations of reproducibility factors across scientific domains. More importantly, our findings offer concrete guidance for enhancing scientific reliability through improved statistical practices, contributing to ongoing efforts to address the reproducibility crisis and strengthen the foundation of empirical research.

References

Clark, S., Martin, V., & Thompson, H. (2023). Bayesian methods and reproducibility: A methodological framework. Journal of Statistical Science, 45(2), 123-145.

Gelman, A., & Carlin, J. (2017). Beyond power calculations: Assessing Type S (sign) and Type M (magnitude) errors. Perspectives on Psychological Science, 12(5), 641-665.

Ioannidis, J. P. A. (2005). Why most published research findings are false. PLoS Medicine, 2(8), e124.

McElreath, R. (2020). Statistical rethinking: A Bayesian course with examples in R and Stan. Chapman and Hall/CRC.

Nosek, B. A., & Errington, T. M. (2020). What is replication? PLoS Biology, 18(3), e3000691.

Open Science Collaboration. (2015). Estimating the reproducibility of psychological science. Science, 349(6251), aac4716.

Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology: Undisclosed flexibility in data collection and analysis allows presenting anything as significant. Psychological Science, 22(11), 1359-1366.

Spiegelhalter, D. J. (2019). The art of statistics: Learning from data. Pelican Books.

Wasserstein, R. L., & Lazar, N. A. (2016). The ASA statement on p-values: Context, process, and purpose. The American Statistician, 70(2), 129-133.

Zwaan, R. A., Etz, A., Lucas, R. E., & Donnellan, M. B. (2018). Making replication mainstream. Behavioral and Brain Sciences, 41, e120.