

1 Scoping Review Protocol: Statistical Models for Longitudinal Data
2 in Health and Biomedical Research: Current State, Challenges,
3 and Opportunities

4 Ariel I. Mundo Ortiz

5 2022-09-09

6 **Table of contents**

7	1 Notes	2
8	2 Registration	3
9	3 Author Contributions	3
10	4 Amendements	3
11	5 Support	3
12	5.1 Sources	3
13	6 Introduction	3
14	6.1 Rationale	3
15	7 Objectives	5
16	8 Review Question	5

17	9 Methods	6
18	9.1 Types of Studies	6
19	9.2 Eligibility Criteria	6
20	9.2.1 For the Application of Modern Statistical Models on Longitudinal Biomedical/Health	
21	Data (Aims 1a and 1b)	6
22	9.2.2 For Methods on Longitudinal Data (Aim 2)	6
23	9.3 Information Sources	7
24	9.4 Search Strategy	7
25	9.4.1 For the Application of Modern Models on Longitudinal Biomedical/Health Data . .	7
26	9.4.2 For Methods on Longitudinal Data	8
27	9.5 Data Collection and Analysis	9
28	9.5.1 Selection Process and Data Management	9
29	9.5.2 Data Collection Process	9
30	9.6 Data Items	9
31	9.7 Risk of Bias in Individual Studies	10
32	9.8 Data Synthesis	10
33	9.9 Meta-Biases	10
34	10 References	10

1 Notes

As of Sept 7, 2022 this document follows the structure recommended by PRISMA-P
<https://prisma-statement.org/documents/PRISMA-P-checklist.pdf>

38 **2 Registration**

39 This section will be populated with the registration number and registry name once the protocol is submitted
40 for peer review.

41 **3 Author Contributions**

- 42 • AM: Writing, query design, data extraction and analysis . . .

43 Other authors to add later

44 **4 Amendments**

45 Protocol amendments resulting from peer review will be indicated in this section indicating the date of each
46 amendment.

47 **5 Support**

48 This section will indicate the sources of financial or other support for the review

49 **5.1 Sources**

50 **6 Introduction**

51 **6.1 Rationale**

52 Longitudinal studies are frequently used in the health sciences (biomedical research, epidemiology, public
53 health, among others) as they allow to examine how the temporal effect of a treatment or an intervention,
54 in contrast to a cross-sectional study, which only allows to examine the effect of the intervention at a single
55 time point. When compared their cross-sectional counterparts, longitudinal studies allow for increased
56 statistical power and more cost efficient strategies^{1,2}. However, the statistical analysis of longitudinal data
57 requires to take into consideration factors such as data missingness, correlation, and non-linear trends,

58 which do not occur on cross-sectional data^{3,4}. In other words, there is an “analytic cost” associated with
59 the increased complexity of longitudinal data².

60 This additional layer of complexity has led to a problem of model misspecification in the statistical analysis
61 of the data (i.e., the use of a statistical model that is not coherent with the data), which has been reported to
62 occur in many fields, including the health sciences⁵. For example, in a landmark study Liu et al. showed that
63 in a subset of papers in the biomedical sciences, the most popular model used to analyze longitudinal data
64 was the analysis of variance (ANOVA, an approach that fails to take into account the correlation between
65 measures over time), and that only 18% of the studies analyzed used models intended for longitudinal
66 analysis while checking that the assumptions of the model were satisfied by the data⁶.

67 Historically, the repeated measures ANOVA (rm-ANOVA, a statistical model for longitudinal data) has
68 been the preferred method in the health sciences to analyze longitudinal data, despite the fact that the
69 multiple assumptions required by this model are frequently not satisfied by the data collected in longitudinal
70 studies⁴. On the other hand, the last 30 years have seen incredible progress in the field of Statistics with the
71 development of statistical models for longitudinal data that relax the assumptions of rm-ANOVA. Linear
72 mixed models, generalized additive mixed models, and generalized estimating equations are among these
73 modern statistical models developed for longitudinal data^{7–11}. From these statistical methods, linear mixed
74 models and generalized estimating equations are the two classes of models that have been frequently applied
75 to analyze longitudinal data in the health sciences during the last decade^{12–14}.

76 However, modern statistical methods that are suited to analyze longitudinal data have been the exception
77 rather than the norm in the health sciences. In 2001, a study reported that only 30% of the clinical trials
78 analyzed used linear mixed models to analyze their results, and that the preferred method of analysis
79 continued to be rm-ANOVA¹⁵ (in comparison, McCullagh and Nelder’s seminal book on the generalized
80 linear model (GLM) was published in 1989¹⁶, and there was ongoing work on the extension of the GLM
81 framework to the mixed model case by 1993¹⁷). Apart from the aforementioned study, there are not recent
82 papers that examine the use of modern statistical methods for longitudinal data in the health sciences.
83 Such information is critical to understand if the use of these methods has increased or decreased in the field
84 over the last 20 years, and the reasons behind such changes.

85 Additionally, the reproducibility crisis is an ongoing issue in the health sciences^{18,19}, a major component of
86 it being the misuse and lack of reproducibility of statistical analyses^{20,21}. Despite the fact that the landscape
87 of statistical software has vastly increased in the last decade with many statistical computational tools now
88 available to researchers, reproducibility standards vary between each computational tool²². Furthermore,

there is still high variability in the amount of statistical reporting across journals²³. Understanding what statistical computational tools are used nowadays by researchers in the health sciences can provide an assessment of the advances in the field towards research reproducibility, while identifying limitations that might still be in place.

7 Objectives

This study aims to:

- Identify the different statistical models for longitudinal data that are used in the health sciences in order to measure the current extent in the adoption of modern statistical methods by the field (Aim 1a)
- Summarize the computational tools used by researchers in the health sciences to statistically analyze longitudinal data to understand the current status of the field with regards to reproducibility. (Aim 1b)
- List statistical methods for longitudinal data developed within the last decade in order to showcase newer methods that may be applicable for longitudinal data in a biomedical/health context. (Aim 2)

8 Review Question

- What are the statistical methods used in biomedical/health sciences research?
- Has the use of modern statistical methods increased in the field during the last 20 years?
- What computational tools are most commonly used by researchers to analyze longitudinal data, and how in turn this affects reproducibility?
- What are most recent statistical methods developed for longitudinal data, and how can they be applied in the health sciences?

110 **9 Methods**

111 **9.1 Types of Studies**

112 For all the study aims, studies included in the analysis correspond to peer-reviewed publications in English.

113 **9.2 Eligibility Criteria**

114 **9.2.1 For the Application of Modern Statistical Models on Longitudinal Biomedical/Health** 115 **Data (Aims 1a and 1b)**

116 **9.2.1.1 Inclusion Criteria**

- 117 • Articles that:
 - 118 – Are written in English
 - 119 – Belong to the biomedical/health sciences fields
 - 120 – Describe the collection and analysis of continuous or discrete longitudinal data
 - 121 – Indicate the statistical model used to analyze the data
 - 122 – Report the results of their statistical analyses

123 **9.2.1.2 Exclusion Criteria**

- 124 • Cross-sectional studies
- 125 • Tutorials that present the application of existing statistical methods to biomedical/health data
- 126 • Reviews, meta-analyses, or systematic reviews on existing statistical methods for longitudinal data
- 127 • Studies that use only descriptive statistics to summarize/analyze the data
- 128 • Studies that collect and analyze categorical data

129 **9.2.2 For Methods on Longitudinal Data (Aim 2)**

130 **9.2.2.1 Inclusion Criteria**

- Articles that:

- Are written in English

- Present new methodologies or significant improvements to existing methods for longitudinal data

9.2.2.2 Exclusion Criteria

- Systematic reviews, meta-analyses, or reviews of statistical methods for longitudinal data
- Tutorials that present the application of existing statistical methods to biomedical/health longitudinal data

9.3 Information Sources

Studies will be retrieved from PubMed and Web of Science.

9.4 Search Strategy

9.4.1 For the Application of Modern Models on Longitudinal Biomedical/Health Data

9.4.1.1 PubMed

9.4.1.1.1 Query:

(biomedical OR health) AND ((repeated measures) OR (longitudinal study) OR (longitudinal data))
AND ((statistical analyses) OR (statistical analysis)) NOT (Review[Publication Type] OR Meta
analy*[Publication Type]) NOT (“Statistics as Topic/methods”[Majr] OR “Statistics as Topic/statistics
and numerical data”[Majr] OR “Models, Statistical”[Mesh] OR “Research Design”[Mesh])

Hits: 10,972

9.4.1.2 Web of Science

150 **9.4.1.2.1 Query:**

151 (ALL=(biomedical) OR ALL=(health)) AND (ALL=(repeated measures) OR ALL=(longitudinal study)
152 OR ALL=(longitudinal data)) AND (ALL=(statistical analyses) OR ALL=(statistical analysis)) AND
153 (DT=(Article)) NOT (WC=(Statistics Probability) OR WC=(Mathematics)) NOT (SU=(Agriculture)
154 AND SU=(Business Economics) AND SU=(Veterinary Sciences) AND SU=(Education Educational
155 Research) AND SU=(Business Economics) AND SU=(Social Sciences Other Topics) AND SU=(Food
156 Science Technology) AND SU=(Anthropology) AND SU=(Linguistics) AND SU=(Sociology) AND
157 SU=(Criminology Penology) AND SU=(Zoology) AND SU=(Meteorology Atmospheric Sciences) AND
158 SU=(Mathematical Methods in Social Sciences) AND SU=(Geology) AND SU=(Construction Build-
159 ing Technology) AND SU=(Geology) AND SU=(Religion) AND SU=(Marine Freshwater Biology)
160 AND SU=(Operations Research Management Science) AND SU=(Fisheries) AND SU=(Metallurgy
161 Metallurgical Engineering))

162 Hits: 12,458

163 **9.4.2 For Methods on Longitudinal Data**

164 **9.4.2.0.1 Query 1:**

165 (“Models, Statistical” [Mesh] OR “Biostatistics/methods”[Mesh]) AND (“Longitudinal Studies”[Mesh])
166 NOT (Review[Publication Type] OR Meta Analys*[Publication Type] OR “editorial”[Publication Type])
167 NOT (“survival”[Title/abstract]) NOT (“tutorial”[title/abstract] OR “orientation”[title/abstract]) NOT
168 (Humans[Mesh] OR Adolescent [Mesh] OR Animals[Mesh])

169 Hits: 142

170 **9.4.2.1 Web of Science**

171 (ALL=(longitudinal studies) OR ALL=(repeated measures)) NOT (TI=(survival)) AND (WC=(Statistics
172 Probability) OR (WC=Mathematics) OR (WC=Mathematics Applied))

173 Hits: 8,135

174 9.5 Data Collection and Analysis

175 9.5.1 Selection Process and Data Management

176 Two reviewers will independently analyze the database search results and pre-screen articles based on ti-
177 tle and abstract content following the aforementioned inclusion/exclusion criteria. Manuscripts from the
178 database(s) search will be stored in the Covidence platform, where duplicated entries will be removed. For
179 articles where pre-screening inclusion (or exclusion) is unclear based on title and abstract analysis, full-text
180 review will be used to make a decision following review by a third independent reviewer. Manuscripts
181 included after title and abstract pre-screening will be further screening by two reviewers that will indepen-
182 dently examine the full text of each article.

183 9.5.2 Data Collection Process

184 Pilot forms (electronic spreadsheets) will be tested using a representative sample of the studies to be
185 reviewed (~100 studies). Information in the forms will be independently included by each reviewer. The
186 forms will be updated (if needed), after the pilot test by consensus between the reviewers.

187 Information obtained from each study (statistical method used, software, etc.) will be tabulated indepen-
188 dently by the reviewers in an electronic spreadsheet.

189 9.6 Data Items

190 Aims 1a and 1b:

- 191 • Statistical method used
- 192 • Sub-area of application (oncology, psychology, public health, etc)
- 193 • Computational tool used
- 194 • Congruence between statistical method used and the data
- 195 • Year of publication

196 Aim 2:

- 197 • Statistical method reported

- Assumptions of the model
- Computational tools available for its implementation
- Year of publication

9.7 Risk of Bias in Individual Studies

N/A

9.8 Data Synthesis

The data from the results of each included study will be extracted into electronic spreadsheets. Summary measures for Aims 1a and 1b include plots (pie, bar, etc) to show the relative use of each statistical method reported, computational tool, and congruence between statistical method and the data. Each plot will be segmented by year to show trends over time.

For Aim 2, a table will be created where statistical method, year of publication, assumptions of the model, and applicability to health data is reported.

9.9 Meta-Biases

N/A

10 References

1. Edwards LJ. Modern statistical techniques for the analysis of longitudinal data in biomedical research. Pediatric Pulmonology. 2000;30(4):330-344. doi:[https://doi.org/10.1002/1099-0496\(200010\)30:4%3C330::AID-PPUL10%3E3.0.CO;2-D](https://doi.org/10.1002/1099-0496(200010)30:4%3C330::AID-PPUL10%3E3.0.CO;2-D)
2. Zeger SL, Liang K-Y. An overview of methods for the analysis of longitudinal data. Statistics in Medicine. 1992;11(14-15):1825-1839. doi:<https://doi.org/10.1002/sim.4780111406>
3. Caruana EJ, Roman M, Hernández-Sánchez J, Solli P. Longitudinal studies. Journal of Thoracic Disease. 2015;7(11):E537-40.

4. Mundo AI, Tipton JR, Muldoon TJ. Generalized additive models to analyze nonlinear trends in biomedical longitudinal data using r: Beyond repeated measures ANOVA and linear mixed models. Statistics in Medicine. Published online July 2022.
5. Thiese MS, Arnold ZC, Walker SD. The misuse and abuse of statistics in biomedical research. Biochem Med (Zagreb). 2015;25(1):5-11.
6. Liu C, Cripe TP, Kim M-O. Statistical issues in longitudinal data analysis for treatment efficacy studies in the biomedical sciences. Molecular Therapy. 2010;18(9):1724-1730. doi:<https://doi.org/10.1038/mt.2010.127>
7. Linear mixed-effects models: Basic concepts and examples. In: Mixed-Effects Models in s and s-PLUS. Springer New York; 2000:3-56. doi:[10.1007/0-387-22747-4_1](https://doi.org/10.1007/0-387-22747-4_1)
8. Jiang J, Nguyen T. Linear and Generalized Linear Mixed Models and Their Applications. 2nd ed. Springer; 2021.
9. Hastie TJ. Statistical Models in S. (Chambers JM, Hastie TJ, eds.). Routledge; 2017.
10. Rosa GJM, Gianola D, Padovani CR. Bayesian longitudinal data analysis with mixed models and thick-tailed distributions using MCMC. Journal of Applied Statistics. 2004;31(7):855-873.
11. Ballinger GA. Using generalized estimating equations for longitudinal data analysis. Organizational Research Methods. 2004;7(2):127-150.
12. Wang M. Generalized estimating equations in longitudinal data analysis: A review and recent developments. Advances in Statistics. 2014;2014:1-11.
13. Tian Q, Qin L, Zhu W, Xiong S, Wu B. Analysis of factors contributing to postoperative body weight change in patients with gastric cancer: Based on generalized estimation equation. PeerJ. 2020;8(e9390):e9390.
14. Şevik M, Doğan M. Epidemiological and molecular studies on lumpy skin disease outbreaks in turkey during 2014-2015. Transboundary and Emerging Diseases. 2017;64(4):1268-1279.
15. Gueorguieva R, Krystal JH. Move Over ANOVA: Progress in Analyzing Repeated-Measures Data and Its Reflection in Papers Published in the Archives of General Psychiatry. Archives of General Psychiatry. 2004;61(3):310-317. doi:[10.1001/archpsyc.61.3.310](https://doi.org/10.1001/archpsyc.61.3.310)
16. McCullagh P, Nelder JA. Generalized Linear Models. Routledge; 2019.

- 245 17. Breslow NE, Clayton DG. Approximate inference in generalized linear mixed models. Journal of the
246 American Statistical Association. 1993;88(421):9-25. doi:[10.1080/01621459.1993.10594284](https://doi.org/10.1080/01621459.1993.10594284)
- 247 18. Jarvis MF, Williams M. Irreproducibility in preclinical biomedical research: Perceptions, uncer-
tainties, and knowledge gaps. Trends in Pharmacological Sciences. 2016;37(4):290-302. doi:<https://doi.org/10.1016/j.tips.2015.12.001>
248
- 249 19. Turkiewicz A, Luta G, Hughes HV, Ranstam J. Statistical mistakes and how to avoid them –
lessons learned from the reproducibility crisis. Osteoarthritis and Cartilage. 2018;26(11):1409-1411.
250 doi:[10.1016/j.joca.2018.07.017](https://doi.org/10.1016/j.joca.2018.07.017)
- 251 20. Gosselin R-D. Statistical analysis must improve to address the reproducibility crisis: The ACcess to
252 transparent statistics (ACTS) call to action. Bioessays. 2020;42(1):e1900189.
- 253 21. Lang TA, Altman DG. Basic statistical reporting for articles published in biomedical journals: The
“statistical analyses and methods in the published literature” or the SAMPL guidelines. Int J Nurs
254 Stud. 2015;52(1):5-9.
- 255 22. Gentleman R, Lang DT. Statistical analyses and reproducible research. Journal of Computational
and Graphical Statistics. 2007;16(1):1-23. Accessed August 16, 2022. [http://www.jstor.org/stable/](http://www.jstor.org/stable/27594227)
256 [27594227](http://www.jstor.org/stable/27594227)
- 257 23. Indrayan A. Reporting of basic statistical methods in biomedical journals: Improved SAMPL guide-
258 lines. Indian Pediatrics. 2020;57(1):43-48. doi:[10.1007/s13312-020-1702-4](https://doi.org/10.1007/s13312-020-1702-4)