

1 Scoping Review Protocol: Statistical Models for Longitudinal Data

2 Ariel I. Mundo Ortiz

3 2022-08-11

4 **Table of contents**

5	<b>1 Background</b>	<b>2</b>
6	<b>2 Objective</b>	<b>3</b>
7	<b>3 Review Question</b>	<b>3</b>
8	<b>4 Databases</b>	<b>3</b>
9	<b>5 Search Terms</b>	<b>4</b>
10	<b>6 Criteria</b>	<b>4</b>
11	6.1 Inclusion Criteria . . . . .	4
12	6.2 Exclusion Criteria . . . . .	4
13	<b>7 Additional Resources</b>	<b>4</b>
14	<b>8 Comparison (?)</b>	<b>4</b>
15	<b>9 Data Extraction</b>	<b>4</b>
16	<b>10 Data Synthesis Strategy</b>	<b>4</b>
17	<b>11 References</b>	<b>4</b>

# 1 Background

Longitudinal studies are frequently used in the health sciences (biomedical research, epidemiology, public health, among others) as they allow to examine how the temporal effect of a treatment or an intervention, in contrast to a cross-sectional study, which only allows to examine the effect of the intervention at a single time point. When compared to cross-sectional studies, longitudinal studies allow for increased statistical power and more cost efficient strategies<sup>1,2</sup>. However, the statistical analysis of longitudinal requires to take into consideration factors such as data missingness, correlation, and non-linear trends, which do not occur on cross-sectional data<sup>3,4</sup>.

This additional layer of complexity in the analysis of longitudinal data has led to a well documented problem of model misspecification (the use of a statistical model that is not coherent with the data) in the health sciences<sup>4</sup>, which can be partly explained by the fact that researchers have a tendency to use the same statistical analysis, methods and tests from other papers without having a clear understanding of the limitations, assumptions, and applicability of the model in each situation<sup>5,6</sup>. For example, in a landmark study Liu et al. showed that in a subset of papers in the biomedical sciences, the most popular model used to analyze longitudinal data was ANOVA (an approach that fails to take into account the correlation between measures over time), and that only 18% of studies used models intended for longitudinal analysis while checking that the assumptions of the model were satisfied by the data<sup>7</sup>.

Historically, the repeated measures analysis of variance (rm-ANOVA) has been the preferred method in the health sciences to analyze longitudinal data, despite the fact that frequently, the assumptions required for its use are not satisfied by the data<sup>4</sup>. On the other hand, over the last 30 years the field of Statistics has been able to develop models for longitudinal data that overcome the limitations of rm-ANOVA, such as linear mixed models, generalized additive models, Bayesian models, and generalized estimating equations<sup>8–12</sup>. However, the adoption of these modern statistical techniques has been slow, as showcased by Gueorguieva et al., who showed that by 2001, only 30% of clinical trials reported in the *Archives of General Psychiatry* used linear mixed models to analyze their results and that rm-ANOVA continued to be the preferred method of analysis in most cases<sup>13</sup>.

During the last decade, the increased availability of computational tools to analyze longitudinal data has lead to increased adoption of modern statistical methods to analyze longitudinal data in the health sciences<sup>14–17</sup>. Despite this, it is not known how much the adoption of these modern statistical methods has increased over the last 20 years, and what are the reasons that may continue to limit the knowledge and application of these statistical methods by researchers in the health sciences. Because research reproducibility continues

49 to be at the center of the debate on biomedical research [citation], there is a need to better understand the  
50 current status of statistical practices in the health sciences in order to implement changes that can lead to  
51 a harmonized used of statistics.

52 To answer this question, in this study we surveyed the statistical methods used in papers dealing with  
53 longitudinal data in health sciences over the last 20 years, in order to gain a better understanding of: 1)  
54 the trends in adoption of modern statistical methods, 2) identify the most frequent pitfalls in statistical  
55 analysis, and 3) provide a rationale for situations where these methods are still not widely adopted.

## 56 **2 Objective**

57 This study aims to summarize the different statistical models for longitudinal data that are used in the  
58 health sciences, identify the extent of the adoption of modern statistical methods in the field, and determine  
59 if in each case, model assumptions are checked by researchers to ensure congruency between the data and  
60 the model.

## 61 **3 Review Question**

62 Summarize the statistical methods used to analyze longitudinal data in the health sciences to identify  
63 which methods are most commonly used, the applicability of such methods in the context of each study,  
64 and gaps that might exist that prevent the adoption of modern statistical methods that can be better suited  
65 to analyze the data. Additionally, identify if studies check for model assumptions, and how this in turn  
66 impacts the reported results.

## 67 **4 Databases**

- 68 • PubMed
- 69 • Web of Science

## 70 5 Search Terms

## 71 6 Criteria

### 72 6.1 Inclusion Criteria

- 73 • methods paper see new methods developed
- 74 • application

### 75 6.2 Exclusion Criteria

## 76 7 Additional Resources

## 77 8 Comparison (?)

## 78 9 Data Extraction

## 79 10 Data Synthesis Strategy

## 80 11 References

- 81 1. Edwards LJ. Modern statistical techniques for the analysis of longitudinal data in biomedical re-  
search. *Pediatric Pulmonology*. 2000;30(4):330-344. doi:[https://doi.org/10.1002/1099-0496\(200010\)](https://doi.org/10.1002/1099-0496(200010)30:4%3C330::AID-PPUL10%3E3.0.CO;2-D)  
82 [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)
- 83 2. Zeger SL, Liang K-Y. An overview of methods for the analysis of longitudinal data. *Statistics in*  
84 *Medicine*. 1992;11(14-15):1825-1839. doi:<https://doi.org/10.1002/sim.4780111406>
- 85 3. Caruana EJ, Roman M, Hernández-Sánchez J, Solli P. Longitudinal studies. *Journal of Thoracic*  
86 *Disease*. 2015;7(11):E537-40.
- 87 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.  
88 *Statistics in Medicine*. Published online July 2022.

5. Ercan I, Yazici B, Yaning Y, et al. Misusage of statistics in medical research. *European Journal of General Medicine*. 2007;4(3):128-134.
6. Thiese MS, Arnold ZC, Walker SD. The misuse and abuse of statistics in biomedical research. *Biochem Med (Zagreb)*. 2015;25(1):5-11.
7. 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>
8. 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)
9. Jiang J, Nguyen T. *Linear and Generalized Linear Mixed Models and Their Applications*. 2nd ed. Springer; 2021.
10. Hastie TJ. *Statistical Models in S*. (Chambers JM, Hastie TJ, eds.). Routledge; 2017.
11. 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.
12. Ballinger GA. Using generalized estimating equations for longitudinal data analysis. *Organizational Research Methods*. 2004;7(2):127-150.
13. 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)
14. Mundo AI, Muhammad A, Balza K, Nelson CE, Muldoon TJ. Longitudinal examination of perfusion and angiogenesis markers in primary colorectal tumors shows distinct signatures for metronomic and maximum-tolerated dose strategies. *Neoplasia*. 2022;32:100825. doi:[10.1016/j.neo.2022.100825](https://doi.org/10.1016/j.neo.2022.100825)
15. Wang M. Generalized estimating equations in longitudinal data analysis: A review and recent developments. *Advances in Statistics*. 2014;2014:1-11.
16. 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.
17. 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.