Why OLS Is A Bad Model For Longitudinal Data

Andy Grogan-Kaylor

2025-01-02

Table of contents

1	Some Beginning Ideas	2				
2	An Empirical Example					
3 Introduction						
4	A First Longitudinal Model	2				
5 What About Change Scores?						
6 What If We Have More Than Two Time Points?						
7	Two Conceptual Diagrams7.1OLS or MLM for 2 Timepoints7.2Cross-Lagged Model	3 3 4				
8 Additionally						
9	Our Answer To the Problem 9.1 Data in Long Format 9.2 Equation 9.2.1 Simple MLM 9.2.2 Longitudinal MLM	4 5 5 5				
10	This Has The Following Advantages:10.1 First10.2 How To Address Missing Data?10.3 Further10.4 Appropriate Metric For Time10.5 Visually10.6 Lastly10.7 Let's continue to explore how this model works.	5 6 6 6 6 7				
Re	eferences	7				

1 Some Beginning Ideas

"Despite the incredible diversity existing among and within human cultures, there are many phenomena that occur regularly in all known societies. These commonalities, or universals, while deriving in part from human nature, may also have specific social, cultural, and systemic sources. We need to develop a working understanding of these universals so that we might advance legitimate, empirically based human science set on creating knowledge that is politically relevant to fostering real solutions to the problems that complicate human co-existence in the Age of the Anthropocene." (Antweiler 2016)

"The language we have in that world is not large enough for the territory that we've already entered." (Whyte and Tippett 2016)

2 An Empirical Example



Figure 1: Happiness as a Function of Time and Pizza

3 Introduction

We are all familiar with the idea of:

$$y_i = \beta_0 + \beta_1 x + e_i \text{ (OLS)}$$

get substantive example

Table 1: Data in WIDE format

id	x1	x2	x3	y1	y2	y3
1						
2						
3						

4 A First Longitudinal Model

We could imagine a longitudinal model where we regress y_i at time 2 on y_i at time 1....

$$y_{i2} = \beta_0 + \beta_1 x + \beta_2 y_{i1} + e_i$$

And we could even make this (*perhaps confusingly*) a multilevel model for individual *i* in social unit *j*: $y_{ij} = -\beta_{ij} + \beta_{ij} x_{ij} + \beta_{ij} y_{ij} + e_{ij}$

$$y_{i2j} = \rho_0 + \rho_1 x + \rho_2 y_{i1j} + u_{0j} + e_{ij}$$

... and add all of the usual random slope terms...

Any problems yet?

5 What About Change Scores?

 $y_{i2}-y_{i1}=\beta_0+\beta_1x+e_i$

? What Happens To The Regression Coefficients in a Change Score Model?

 βy_{i1}

6 What If We Have More Than Two Time Points?

 $y_{i3}=\beta_0+\beta_1x+\beta_2y_{i1}+\beta_3y_{i2}+e_i$

💡 Tip

What is the problem here? We have 2 terms that are likely to be collinear: $\beta_2~\&~\beta_3$

This issue only becomes worse the more time points we add.

As a result, we are not really modeling y_2 and y_1 .

7 Two Conceptual Diagrams

7.1 OLS or MLM for 2 Timepoints



Figure 2: An OLS Or Multilevel Model For 2 Timepoints

7.2 Cross-Lagged Model



Figure 3: A Cross Lagged Model For 3 Timepoints

8 Additionally ...

🛕 No Explicit Function of Time

Additionally, we do not have an explicit function of time. We don't know really have a clear idea of whether our outcome increases with time, or decreases with time. Or whether the effect is curvilinear e.g. t^2 or $\ln(t)$.

1 Unbalanced Data Are A Problem

Additionally, any data that is unbalanced i.e. study participants enter the study late, or leave the study early are going to be difficult for this kind of model to deal with.

🛕 Missing Data Are A Problem

Similarly, data that is missing at one time point, but present at other time points, is going to be a problem for this kind of model. (and it is going to be difficult for many of our colleagues to see how we can get around this issue.)

9 Our Answer To the Problem

We Reshape The Data and Use the SAME Notation!!!

"Mathematics is the art of giving the same name to different things." (Poincare 1908)

9.1 Data in Long Format

id	t	х	у
1	1		
1	2		
1	3		
2	1		
2	2		
2	3		
3	1		
3	2		
3	3		

9.2 Equation

So.... we take our standard multilevel notation.

9.2.1 Simple MLM

$$y_{ij} = \beta_0 + \beta_1 x + u_{0j} + e_{ij} \tag{1}$$

cross out j write in t.

9.2.2 Longitudinal MLM

$$y_{it} = \beta_0 + \beta_1 t + u_{0i} + e_{it} \tag{2}$$

Person-Observations

Every row is a *person-observation* (person i observed at time t). Every person has *multiple rows*.

10 This Has The Following Advantages:

10.1 First...

- 1. No multicollinearity issue. By inspection of Equation 2, we see that there is only a single β coefficient for each variable, \therefore no multicollinearity problem.
- 2. Unbalanced data is less of a problem, the data structure and estimation are robust to these possibilities (Singer and Willett 2003; Raudenbush and Bryk 2002).
- 3. Missing data is less of a problem (assuming MCAR). When a person observation is missing, that person simply has fewer rows of data (J. Hox 2010; Luke 2004; Raudenbush and Bryk 2002; Rabe-Hesketh and Skrondal 2012). But all rows of data are "matched" to the same person by *i*.

Addressing Missing Data is Complicated!!!

It is sometimes best to (a) do nothing; (b) do something complicated.

- Ignore it.
- Fill in the mean.
- Use previous observation.
- Use next observation.
- Linearly interpolate previous and next observation.
- Regression imputation.
- Multiple imputation.

10.3 Further...

- 3. We now have an *explicit function of time* $\beta_1 t$, and could treat time more flexibly, by creating a polynomial function of time e.g. by adding $\beta_2 t^2$, etc. (Raudenbush and Bryk 2002; Singer and Willett 2003). (We could even substitute $\beta \ln(t)$.)
- 4. Again, by inspection of Equation 2, we see that *multiple or many time-points are not a problem*. Same algebra for 2 time points as for 10,000 time points. (Helpful when we start to think about intensive longitudinal data *e.g.* George Holden's *recording study*).
- 5. We are measuring exactly the time at which events take place for each individual (Singer and Willett 2003; Luke 2004). Not simply saying Wave 1, Wave 2, Wave 3, etc...
- 6. Every individual could have a *completely different set of time points* and even a *completely different number of time points* (J. Hox 2010; J. J. Hox, Moerbeek, and van de Schoot 2018; Singer and Willett 2003; Luke 2004).

And we can even add βx back into the model.

10.4 Appropriate Metric For Time

Caution

We do need to think carefully about what is the appropriate variable for time. Is it the variable we used to reshape the data-often wave-or some other more appropriate metric, like age or time in study (Singer and Willett 2003)?

10.5 Visually

10.6 Lastly

Caution

Generating appropriate descriptive statistics can be a problem.



Figure 4: A Multilevel Model For Longitudinal Data

10.7 Let's continue to explore how this model works.

References

Antweiler, Christoph. 2016. Our Common Denominator: Human Universals Revisited. Berghahn.

Hox, Joop. 2010. Multilevel Analysis: Techniques and Applications. 2nd ed. Routledge.

- Hox, Jop J, Mirjam Moerbeek, and Rens van de Schoot. 2018. Multilevel Analysis: Techniques and Applications. Multilevel Analysis: Techniques and Applications. Third edition. Routledge, Taylor & Francis Group,.
- Luke, Douglas. 2004. Multilevel Modeling. SAGE Publications, Inc. https://doi.org/10.4135/ 9781412985147.

Poincare, Henri. 1908. Science Et Methode. Flammarion.

- Rabe-Hesketh, Sophia, and Anders Skrondal. 2012. Multilevel and Longitudinal Modeling Using Stata Volume i: Continuous Responses. Stata Press. 3rd ed. Stata Press.
- Raudenbush, Stephen W, and Anthony S Bryk. 2002. *Hierarchical Linear Models: Applications and Data Analysis Methods*. Sage Publications.
- Singer, Judith D, and John B Willett. 2003. Applied Longitudinal Data Analysis : Modeling Change and Event Occurrence. Applied Longitudinal Data Analysis : Modeling Change and Event Occurrence. Oxford University Press.
- Whyte, David, and Krista Tippett. 2016. "David Whyte: Seeking Language Large Enough." The On Being Project. https://onbeing.org/programs/david-whyte-seeking-language-large-enough/.