ACRM

AMERICAN CONGRESS OF REHABILITATION MEDICINE

# Tips, tricks, and FAQs for getting started in longitudinal data analysis.

Keith Lohse, PhD

Department of Health, Kinesiology, & Recreation

Department of Physical Therapy and Athletic Training

University of Utah

rehabinformatics@gmail.com

Allan J. Kozlowski, PhD

Director of Outcomes Research, Mary Free Bed Rehabilitation Hospital

Department of Epidemiology and Biostatistics, Michigan State University College of Human Medicine
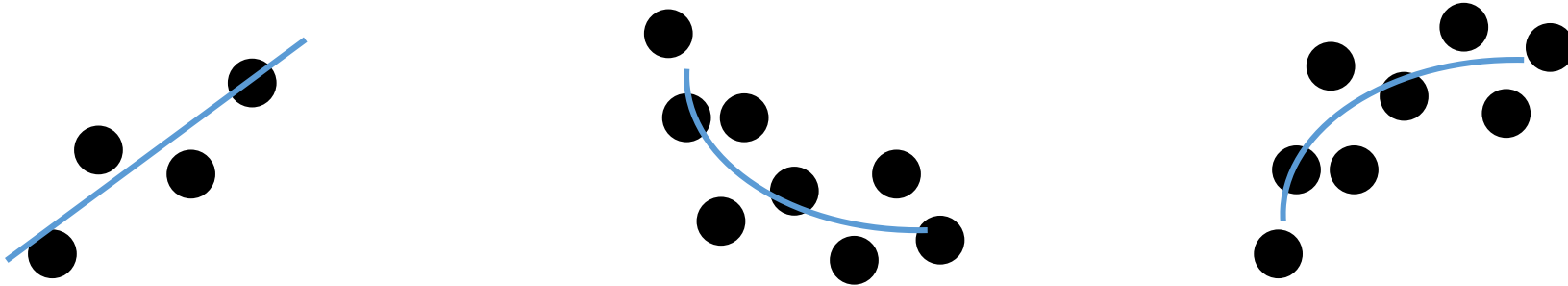
allan.kozlowski@maryfreebed.com

NRIL

NEUROREHABILITATION INFORMATICS LAB

UNIVERSITY OF UTAH COLLEGE OF HEALTH

Restoring Hope and Freedom
Mary Free Bed
Rehabilitation Hospital

MICHIGAN STATE UNIVERSITY

1

# 1. Know your types of variables.

- ***Static/ "Fixed"*** variables are variables that keep the same value over the course of the study.
  - For most longitudinal studies, these are variable that vary between people but stay constant within a person (e.g., gender and age at start of study are example static variables).
  - Can be continuous or categorical.

- ***Dynamic/ "Time Varying"*** variables are variables that change value over the course of the study.
  - Our principle dynamic independent variable is Time (but this could be seconds, months, or years, depending on the resolution over your data).
  - Most of our dependent variables are also dynamic (i.e., we might have BBS, WMFT, or 10m WT scores at each time point).
  - Can be continuous or categorical.

# 2. How will you model time?

- Do most people tend to change linearly or non-linearly?
  - Is there a between-subjects variable associated with different change curves?
  - Exploratory data visualization I really helpful here and can inform subsequent model building.
- Remember that the more complicated your hypotheses about time, the more time-points you will need to collect.

# 3. What does zero mean in your model?

**Continuous Variables**

- Do you have an interpretable zero in your independent variables:
  - Age versus Onset days (Age = 0 doesn't make sense; Onset = 0 might).

- Have you mean-centered the variables in your model?
  - If all variables are mean-centered, you can interpret the effects of one variable "on average" across the other variables.

- Is there a separate value you want to center your variables on?
  - I.e., look at group differences at the end rather than beginning by making the terminal point the intercept.
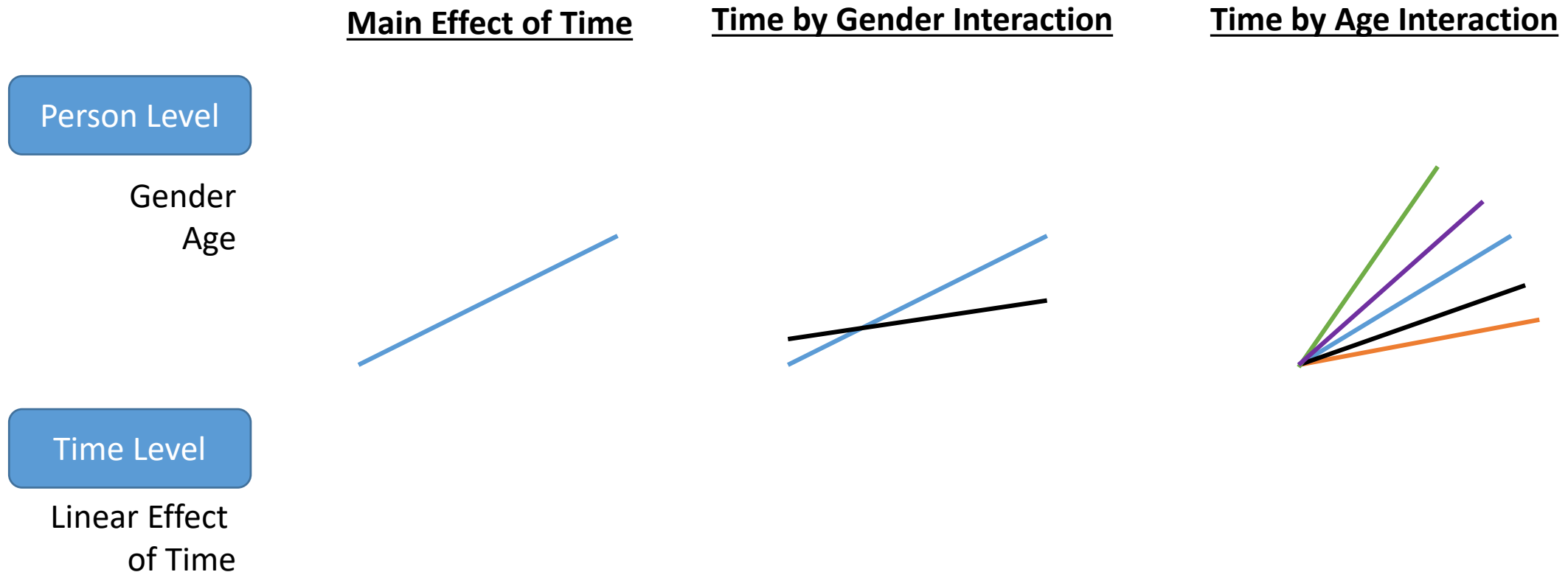
**Categorical variables**:

- Contrast coded versus dummy coded variables.
  - Contrast codes make zero the average, dummy codes make the reference group zero.

# 4. Levels of measurement

- Do you really have an interval level variable?
  - The errors that result from treating non-interval data as interval data actually get worse across time (e.g. FIM scores).
    - Especially anytime you break a scale into subscales!
  - Unequal differences in scale warp the shape of the time function.

- Look into Rasch Scaling as an approach to "intervalizing" ordinal data.
  - Pragmatically, check your residuals and the assumptions of the MLM.

# 5. What effects are you interested in?

**Main Effect of Time**     **Time by Gender Interaction**     **Time by Age Interaction**

Person Level

Gender

Age

Time Level

Linear Effect
of Time



- Often, we are interested in interactions between the person-level and the time-level, but we can also test main-effects and interactions within the person-level or within the time-level!
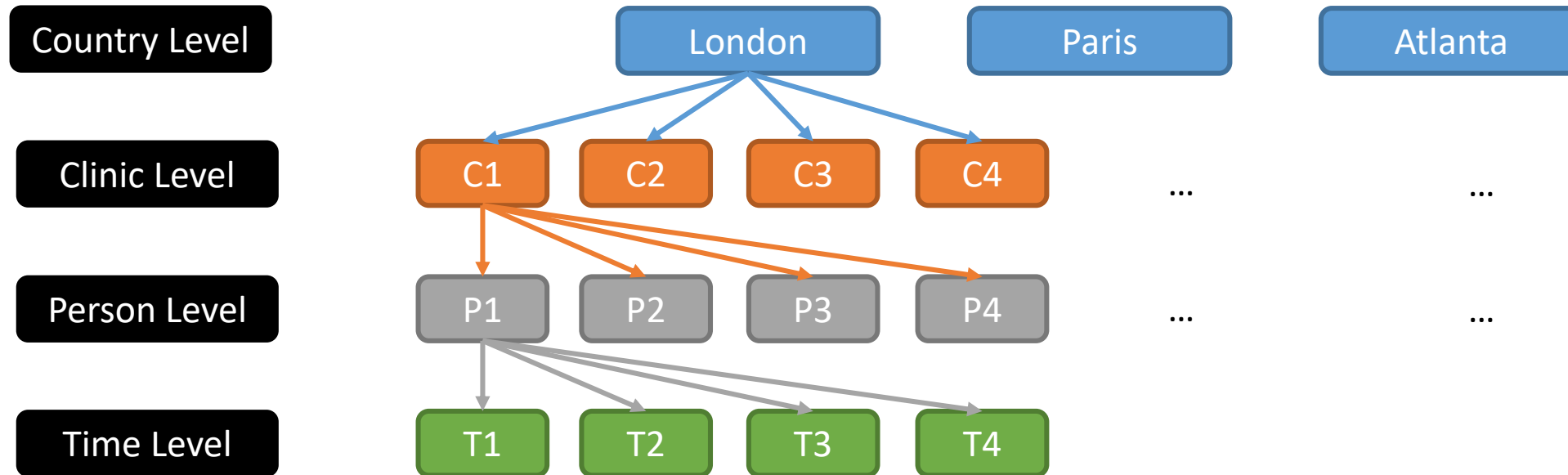
# 6. How do I compare between models?

- Models can have different methods of estimation in order to fit their parameters:
  - ML – maximum likelihood estimation.
  - REML – restricted maximum likelihood estimation.


- Often we prefer ML to REML because it allows us to compare nested models using likelihood based methods like the change in deviance or the Akaike Information Criterion (AIC).
  - Deviance is a measure of the amount of error in a model, so lower deviance means a better model.
    - This can be tested statistically with the Wald Test of the change in deviance.
  - AIC is also a measure of error in a model, so lower AIC means a better model.
    - However, the AIC also introduces a penalty for the number of parameters in a model. This makes the AIC more conservative and helps prevent "over-fitting" of the model.

# 7. How do I statistically power a longitudinal study?

- Statistical power for multi-level models gets pretty complicated, so it is highly recommended that you talk to a statistical consultant. In preparation for that meeting, you'll want to be able to phrase your main narrative hypothesis as a statistical hypothesis like the following:
  - "I am interested in the main-effect of time."
    - You will need to estimate how much you expect participants to change over time, estimate the average standard deviation at each time point, and the average correlation between time points.
  - "I am interested in the interaction of time and group."
    - You will need to estimate all of the same information as above, but you will need to estimate it for each group.

- As a rule of thumb, increasing the number of **time-points** will improve power for effects at the time-level and person by time interactions.
  - Increasing the number of **participants** will improve power for effects at the person-level and person by time interactions.

# 8. What if I have multiple levels?

- Multi-level models can do that!
  - Let's say that you are running large international study…
  - Or combining data from lot's of different studies in secondary analysis…

# 9. What are Fixed-Effects and Random-Effects?

Remember the general concept of DATA = MODEL + Error.
This can be more elaborately written as:

$$y_{ij} = B_0 + U_{0j} + (B_1 + U_{1j}) * (TIME_{ij}) + \epsilon_{ij}$$

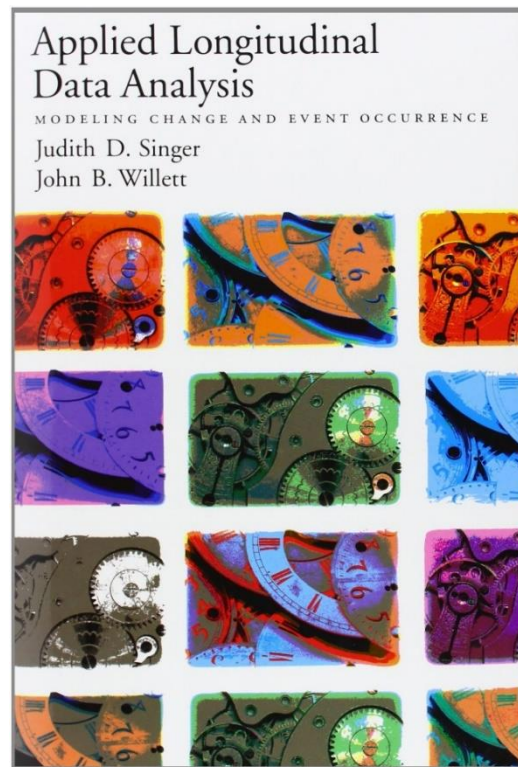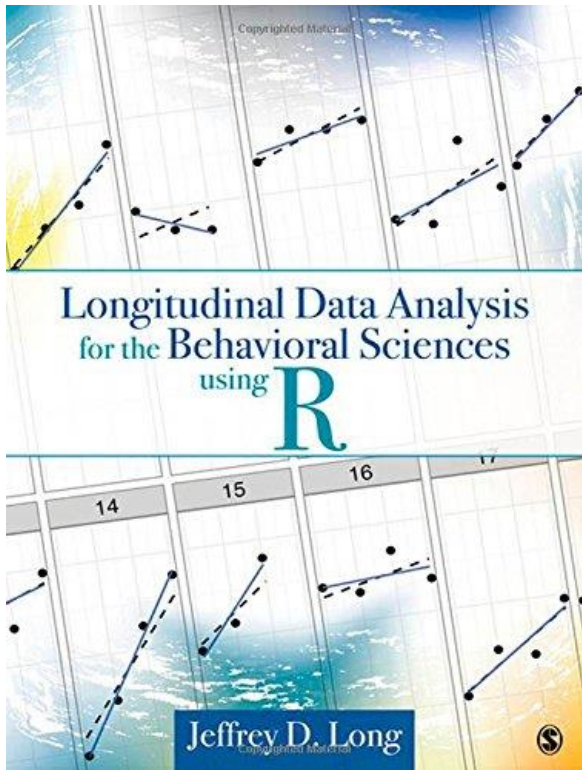Thus, we have the following terms in our *DATA* ($y_{ij}$'s):

- The *MODEL* includes fixed effects and random effects.
- ***Fixed-Effects*** are the group-level $B$'s, these effects parallel the traditional main-effects and interactions that you have probably encountered in other statistical analyses.
- ***Random-Effects*** are the participant-level $U_j$'s that remove statistical dependency from our data. (This is bit of a simplification, but you can think of not including the appropriate random-effects like running a between-subjects ANOVA when you should be running a repeated-measures ANOVA.)
- The *ERRORS*, or more specifically *Random Errors ($\epsilon_{ij}$'s)*, are the difference between our *MODEL*'s predictions and the actual *DATA*.

# 10. Assumptions

- Normality
  - Does transforming the DV change the model?
  - Make it clear to readers that you tested transformed and raw DVs.
- Homoscedasticity
- Scale Invariance
  - Is there bias in the models predictions?
  - Explore methods for looking at measurement variance over time.
- Influential data points/sources
  - There are tools for checking Cook's Distances and VIFs in MLMs.
  - Influential data don't always show up in univariate plots/analyses.
    - See philosophical discussions about outliers and removing influential data points.
    - Run the model both ways and be transparent about what you did in your write-up.
- Floor/Ceiling Effects

# 11. How can I actually run my multi-level models?

- There are numerous texts to help and software packages to do it. They are all slightly different, but users need the same basic understanding of fixed-effects and random-effects to make sure models run correctly.



Longitudinal Data Analysis for the Behavioral Sciences using R

Jeffrey D. Long



Applied Longitudinal Data Analysis
MODELING CHANGE AND EVENT OCCURRENCE
Judith D. Singer
John B. Willett

*We will be using:*



R and R Studio
- Packages:
  - lme4
  - ggplot
  - dplyr

*But you can also use:*



SPSS
AN IBM® COMPANY



sas
THE POWER TO KNOW.



HLM7
SSI SCIENTIFIC SOFTWARE INTERNATIONAL
HIERARCHICAL LINEAR & NONLINEAR MODELING

Most of what I will say has been said better in these resources!

# 12. Be up front about your limitations.

- Exploratory modelling.
  - You will test a lot of things you probably didn't plan on testing, but be transparent in the reporting of your analyses.
  - The dataset that generates a model/prediction cannot also be used to confirm that model/prediction.


- Are your results "robust" to the method of analysis?
  - A lot of issues about how you are modelling times, do you meet normality, or should exclude an influential data point can be addressed by running the model both ways.
  - Is the answer the same both time? Is a difference in the answer meaningful?