

RESEARCH

Open Access



# Detecting influential subjects in intensive longitudinal data using mixed-effects location scale models

Xingruo Zhang<sup>1\*</sup> and Donald Hedeker<sup>1</sup>

## Abstract

**Background** Collection of intensive longitudinal health outcomes allows joint modeling of their mean (location) and variability (scale). Focusing on the location of the outcome, measures to detect influential subjects in longitudinal data using standard mixed-effects regression models (MRMs) have been widely discussed. However, no existing approach enables the detection of subjects that heavily influence the scale of the outcome.

**Methods** We propose applying mixed-effects location scale (MELS) modeling combined with commonly used influence measures such as Cook's distance and DFBETAS to fill this gap. In this paper, we provide a framework for researchers to follow when trying to detect influential subjects for both the scale and location of the outcome. The framework allows detailed examination of each subject's influence on model fit as well as point estimates and precision of coefficients in different components of a MELS model.

**Results** We simulated two common scenarios in longitudinal healthcare studies and found that influence measures in our framework successfully capture influential subjects over 99% of the time. We also re-analyzed data from a health behavior study and found 4 particularly influential subjects, among which two cannot be detected by influence analyses via regular MRMs.

**Conclusion** The proposed framework can help researchers detect influential subject(s) that will be otherwise overlooked by influential analysis using regular MRMs and analyze all data in one model despite influential subjects.

**Keywords** Cook's distance, Influential data, Intensive longitudinal data, Mixed-effects location scale model, Variance modeling

## Background

In the past decades, intra-individual variability, also called within-subject (WS) heterogeneity or level-1 heterogeneity, of health behaviors and conditions has received increased attention [1–3]. Accordingly, new developments in statistical modeling, including mixed-effects

location scale (MELS) models proposed by Hedeker et al. [4] and Nordgren et al. [5], allow researchers to model the mean, or location, and variability, or (square of the) scale of responses simultaneously.

Although these models can accommodate between-subject (BS) heterogeneity, also called level-2 heterogeneity, through random subject effects, there often exist subjects who are so different from the others that they need additional analyses and/or may affect model estimation. For example, in ecological momentary assessments (EMA) studies, a typical type of careless responses is that subjects give many consecutive items the same answer

\*Correspondence:

Xingruo Zhang  
xrzhang@uchicago.edu

<sup>1</sup> Department of Public Health Sciences, The University of Chicago, 5841 South Maryland Ave MC2000, Chicago 60637, IL, USA



© The Author(s) 2023. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

[6]. In this case, even if the data from these subjects do not deviate from any pattern in mean responses that subjects correctly responding to questions may show, they exhibit exceptionally low variability. Such abnormal behavior can distort the parameter estimates, especially the estimated variance of random scale effects. On the other hand, subjects who display exceptional consistency during behavioral studies, such as a study that aims at increasing subjects' physical activities, might be of interest as such a pattern may indicate good adherence to the study protocol. Being able to identify such subjects is important in terms of informing personalized health-care. Especially, there have been discussions about using mixed-effects regression models (MRMs) to find individuals who behave differently from others regarding their health conditions and outcomes [7, 8]. Since MELS models are an extension of standard MRMs, they can also perform similar functions and offer additional information about WS variability.

While there has been extensive literature discussing influential subjects and observations in MRMs [7, 9, 10] and software developed to implement these methods [11], to the best of our knowledge, there lack similar methods for the scale model of the response. As we will show in later sections, not considering WS variability in influence analysis can leave subject(s) with abnormal WS variability undetected. Hence, in the topic of influence analysis, WS variability deserves equal attention as the mean of the outcomes.

Therefore, we propose a framework for detecting influential subjects using MELS models. We examine each subject's influence on model fit as well as point estimates of parameters and their efficiency. The proposed method enables researchers to identify subjects who exhibit dubious patterns in their WS variability and/or mean, thus facilitating research on intra-individual variability. Unlike the traditional case-deletion approach of influential data detection, we adopt the method proposed by Langford and Lewis [9]. If a subject is being examined for its influence, it is removed from the estimation of level-2 effects and given subject-specific fixed effects. We will demonstrate how the estimation of the leave-one-out models can be carried out in SAS and R. Also, a health behavior study will be used as an example to illustrate the proposed methodology. Finally, we will assess the performance of our method and demonstrate its benefits over influence analysis using MRMs via simulated examples.

## Methods

### Leave-one-out MELS model

In this section, we first briefly review MELS models developed by Hedeker et al. [4]. For simplicity, only random intercepts with scalar variances are included in both

the location and the scale models described below. Nevertheless, more complicated models that include random slopes of time-varying covariates and/or have covariates influence the variances of the random effects are possible.

Suppose that subject  $i$  ( $i = 1, 2, \dots, N$ ) is measured at visit  $j$  ( $j = 1, 2, \dots, n_i$ ). The model for the response  $y_{ij}$  can be expressed as:

$$y_{ij} = \beta_0 + v_i + x_{ij}^T \beta' + \epsilon_{ij}, \tag{1}$$

where  $x_{ij}$  is a  $p \times 1$  vector of time-varying covariates influencing the mean of  $y_{ij}$ , and  $v_i$  is subject  $i$ 's random intercept, indicating subject  $i$ 's deviation from the fixed part of the model.  $\beta_0$  is the fixed intercept, i.e. average response when all covariates equal 0, and  $\beta'$  is a  $p \times 1$  vector of coefficients corresponding to  $x_{ij}$ .  $\epsilon_{ij}$  is the level-1 residual and follows a normal distribution with a mean of 0 and a variance of  $\sigma_{\epsilon_{ij}}^2$ . Later on, the notation  $\beta$  in influence measures represents a vector of which the first entry is  $\beta_0$  and the remaining entries come from  $\beta'$ .

In the MELS model, a log-linear sub-model is applied to the WS variance of the response to ensure a positive value:

$$\sigma_{\epsilon_{ij}}^2 = \exp\left(\tau_0 + \omega_i + w_{ij}^T \tau'\right), \tag{2}$$

where  $w_{ij}$  is an  $r \times 1$  vector of time-varying covariates influencing the scale of  $y_{ij}$ , and  $\omega_i$  is the  $i^{th}$  subject's deviation from the average log WS variance of  $y_{ij}$ .  $\tau_0$  is the average log WS variance when every covariate in  $w_{ij}$  is 0, and  $\tau'$  is an  $r \times 1$  vector of coefficients of  $w_{ij}$ . Likewise,  $\tau$  in influence measures refers to a vector of which the first entry is  $\tau_0$ , and the remaining entries come from  $\tau'$ .

The random location and scale effects can correlate with each other, and they are assumed to follow a bivariate normal distribution:

$$\begin{pmatrix} v_i \\ \omega_i \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_v^2 & \sigma_{v\omega} \\ \sigma_{v\omega} & \sigma_\omega^2 \end{pmatrix}\right). \tag{3}$$

Given our focus on influence analysis at the subject level, also known as level 2 in longitudinal models, it's important to pause here and clarify that level-2 effects in MELS models comprise fixed effects of time-invariant covariates (including the fixed intercepts) and all random effects. This distinction lays the foundation for our forthcoming exploration of influence analysis using MELS models. The following paragraphs present how the leave-one-out MELS models are formed, i.e., how each subject is separated from the random effects and subject-level fixed effects and given subject-specific fixed effects, for subsequent influence analyses.

As suggested by Langford and Lewis [9], in order to separate the subject under evaluation, denoted as  $i^*$

( $i^* \in 1, 2, \dots, N$ ), from level 2 of the location model, we exclude subject  $i^*$  from the estimation of the random effects and level-2 fixed effects. In this case, the only level-2 fixed effect is  $\beta_0$ . Then, subject  $i^*$  is given a subject-specific fixed effect denoted as  $c_{i^*}$ . Namely, the location model becomes

$$y_{ij} = \mathbb{1}(i \neq i^*) \times (\beta_{0(-i^*)} + v_i) + x_{ij}^T \beta'_{(-i^*)} + \mathbb{1}(i = i^*) \times c_{i^*} + \epsilon_{ij}, \tag{4}$$

where  $\mathbb{1}(i = i^*)$  equals 1 for subject  $i^*$  and 0 for all other subjects, and  $\mathbb{1}(i \neq i^*)$  equals 0 for subject  $i^*$  and 1 for all other subjects.

Note that for simplicity in illustrating leave-one-out MELS models, we assume that every covariate in  $x_{ij}$  and  $w_{ij}$  is time-varying. When some time-invariant covariates are also present, subject  $i^*$  should be separated from the estimation of their associated coefficients as well.

To also separate subject  $i^*$  from level 2 of the scale model described in Eq. 2, the leave-one-out scale model is as follows:

$$\sigma_{\epsilon_{ij}}^2 = \exp(\mathbb{1}(i \neq i^*) \times (\tau_{0(-i^*)} + \omega_i) + w_{ij}^T \tau'_{(-i^*)} + \mathbb{1}(i = i^*) \times d_{i^*}). \tag{5}$$

Here,  $d_{i^*}$  is the  $i^{*th}$  subject's subject-specific fixed scale effect, and the variances of random location effects and random scale effects become  $\sigma_{v(-i^*)}^2$  and  $\sigma_{\omega(-i^*)}^2$  while their covariance is  $\sigma_{v\omega(-i^*)}$  without subject  $i^*$ . Note that the subscript ( $-i^*$ ) here as well as in Eqs. 4 and 5 does not mean that subject  $i^*$  is entirely removed from modeling but instead only removed from the estimation of level-2 effects.

Overall, the algorithm to estimate a leave-one-out MELS model includes the following steps: (1) create an indicator variable that equals 1 for subject  $i^*$  and 0 otherwise; (2) use the opposite of the indicator variable in (1) as the covariate for the random intercepts; (3) following step (2), subject  $i^*$  will not be included in the random effect estimation, and the fixed location effect for the indicator variable would be  $c_{i^*}$  while the fixed scale effect is  $d_{i^*}$ . Our programming specifics and pertinent material will be explained in Section "Results".

A separate leave-one-out model is estimated for every subject and compared with the model in which all subjects are treated the same. Besides allowing the detection of highly influential subjects, this leave-one-out structure can be viewed as a way to keep all subjects in the model despite influential subject(s). Using all available data to estimate the level-1 fixed effects has the benefit of increasing statistical power.

### Influence analysis

#### Influence on model fit

Given that the model described in Eqs. 1 and 2 are nested in the leave-one-out model described in Eqs. 4 and 5, likelihood-ratio test was used to examine subject  $i^*$ 's influence on the model fit. The test statistic, difference in deviance, denoted as  $LR_{i^*}$  for subject  $i^*$ , can be calculated as follows:

$$LR_{i^*} = 2 \ln \left( \frac{\mathcal{L}_{(-i^*)}}{\mathcal{L}} \right), \tag{6}$$

where  $\mathcal{L}_{(-i^*)}$  represents the likelihood of the MELS model in which subject  $i^*$  is separated, and  $\mathcal{L}$  represents the likelihood of the naive MELS model. Since N tests are conducted simultaneously, the false discovery rate (FDR) procedure [12] is applied to adjust for multiple comparisons. This specific multiple testing correction method is chosen because of its ability to control for both Type I and Type II errors [13].

#### Influence on parameter estimates

Cook's distance [14] is a well-known measure for the influence of data on the point estimates of a group of parameter estimates. Based on the structure of MELS models, we calculate the Cook's distances for three groups of parameters, namely, fixed location effects ( $\beta$ ), fixed scale effects ( $\tau$ ), and variances and covariance of random effects, denoted as  $\eta$  ( $\eta = [\sigma_v^2, \sigma_\omega^2, \sigma_{v\omega}]$ ). Following the formula of Cook's distance for multilevel models described by Snijder and Bosker [15], Cook's distance can be calculated as

$$C_{i^*}^\gamma = \frac{1}{r_\gamma} (\hat{\gamma} - \hat{\gamma}_{(-i^*)})^T \hat{\Sigma}_{\hat{\gamma}_{(-i^*)}}^{-1} (\hat{\gamma} - \hat{\gamma}_{(-i^*)}), \tag{7}$$

where  $\gamma$  can be  $\beta$ ,  $\tau$ , or  $\eta$ . Here,  $r_\gamma$  is the number of parameters being examined, and  $\Sigma_{\hat{\gamma}_{(-i^*)}}$  is the variance-covariance matrix of  $\hat{\gamma}$  after subject  $i^*$  is separated from the random effect estimation. A large Cook's distance indicates a heavy influence on a specific group of parameter estimates.

Once a subject is determined to be influential on a particular group of parameter estimates, it is often of interest to investigate this subject's influence on each specific parameter estimate within this group. In this case, DFBETAS can be used. DFBETAS is the difference between the estimate of a parameter obtained when all subjects are kept in the random effect and level-2 effect estimation and the parameter estimate when a subject is separated, divided by the standard error of the parameter estimate [16]. Its mathematical representation is as follows:

$$DFBETAS_{i^*}^\theta = \frac{\hat{\theta} - \hat{\theta}_{(-i^*)}}{SE(\hat{\theta}_{(-i^*)})}, \tag{8}$$

where  $\theta$  can be any single parameter in the model. DFBETAS can have both positive and negative values. A highly positive value indicates that the inclusion of the influential subject creates a positive bias for  $\hat{\theta}$ , and vice versa.

**Influence on the precision of parameter estimates**

Influential subjects can also affect the variances of parameter estimates besides their point estimates. COVTRACE proposed by Christensen et al. [10] and COVRATIO described by Belsley et al. [17] are often used to measure such changes.

COVTRACE $_{i^*}^\gamma$  is calculated as the absolute value of the difference between the trace of the inverse ratio of variance-covariance matrix estimates with and without subject  $i^*$  in level-2 model and the number of parameters under investigation:

$$COVTRACE_{i^*}^\gamma = \left| \text{Tr}(\hat{\Sigma}_\gamma^{-1} \hat{\Sigma}_{\gamma(-i^*)}) - r_\gamma \right|. \tag{9}$$

The larger the COVTRACE value, the greater the influence of subject  $i^*$  on the precision of  $\hat{\gamma}$ .

Meanwhile, COVRATIO $_{i^*}^\gamma$  is the inverse ratio of the determinant of the estimated variance-covariance matrix of  $\gamma$  with and without subject  $i^*$  in level-2 model:

$$COVRATIO_{i^*}^\gamma = \frac{\det(\hat{\Sigma}_{\hat{\gamma}(-i^*)})}{\det(\hat{\Sigma}_{\hat{\gamma}})}. \tag{10}$$

Again,  $\gamma$  is one of  $\beta$ ,  $\tau$ , and  $\eta$ . The determinant of a variance-covariance matrix is also known as the generalized variance, a measure of multidimensional scatter [18]. The examples in the next two sections will focus on COVRATIO because it provides information on both the magnitude and the direction of the influence. A value of COVRATIO above 1 indicates a loss in precision when separating subject  $i^*$ , and on the contrary, a value below 1 indicates a gain in precision.

All the procedures of influence analysis described in Section “Influence analysis” are summarized in Table 1.

**Results**

**Application to health behavior study example**

Data collected by Flueckiger et al. [19] for a study on the association between health behaviors and learning goal

**Table 1** Influence analysis framework for MELS models

Influence measure	Influence sub-category	Formula
Influence on model fit		
Difference in deviance		$LR_{i^*} = 2 \ln \left( \frac{\mathcal{L}_{(-i^*)}}{\mathcal{L}} \right)$
Influence on point estimates of a group of parameters		
Cook’s distance	Fixed location effect estimates	$C_{i^*}^\beta = \frac{1}{r_\beta} (\hat{\beta} - \hat{\beta}_{(-i^*)})^T \hat{\Sigma}_{\hat{\beta}(-i^*)}^{-1} (\hat{\beta} - \hat{\beta}_{(-i^*)})$
	Fixed scale effect estimates	$C_{i^*}^\tau = \frac{1}{r_\tau} (\hat{\tau} - \hat{\tau}_{(-i^*)})^T \hat{\Sigma}_{\hat{\tau}(-i^*)}^{-1} (\hat{\tau} - \hat{\tau}_{(-i^*)})$
	Variances and covariances of random effects	$C_{i^*}^\eta = \frac{1}{r_\eta} (\hat{\eta} - \hat{\eta}_{(-i^*)})^T \hat{\Sigma}_{\hat{\eta}(-i^*)}^{-1} (\hat{\eta} - \hat{\eta}_{(-i^*)})$
Influence on point estimate of a single parameter		
DFBETAS		$DFBETAS_{i^*}^\theta = \frac{\hat{\theta} - \hat{\theta}_{(-i^*)}}{SE(\hat{\theta}_{(-i^*)})}$
Influence on variances and covariances of a group of parameters		
COVTRACE	Fixed location effect estimates	$COVTRACE_{i^*}^\beta = \left  \text{Tr}(\hat{\Sigma}_{\hat{\beta}}^{-1} \hat{\Sigma}_{\hat{\beta}(-i^*)}) - r_\beta \right $
	Fixed scale effect estimates	$COVTRACE_{i^*}^\tau = \left  \text{Tr}(\hat{\Sigma}_{\hat{\tau}}^{-1} \hat{\Sigma}_{\hat{\tau}(-i^*)}) - r_\tau \right $
	Variances and covariances of random effects	$COVTRACE_{i^*}^\eta = \left  \text{Tr}(\hat{\Sigma}_{\hat{\eta}}^{-1} \hat{\Sigma}_{\hat{\eta}(-i^*)}) - r_\eta \right $
COVRATIO	Fixed location effect estimates	$COVRATIO_{i^*}^\beta = \frac{\det(\hat{\Sigma}_{\hat{\beta}(-i^*)})}{\det(\hat{\Sigma}_{\hat{\beta}})}$
	Fixed scale effect estimates	$COVRATIO_{i^*}^\tau = \frac{\det(\hat{\Sigma}_{\hat{\tau}(-i^*)})}{\det(\hat{\Sigma}_{\hat{\tau}})}$
	Variances and covariances of random effects	$COVRATIO_{i^*}^\eta = \frac{\det(\hat{\Sigma}_{\hat{\eta}(-i^*)})}{\det(\hat{\Sigma}_{\hat{\eta}})}$

achievements are used as an example to illustrate the proposed method.

During the 32-day span of the study, 72 students answered questions about their sleep quality (SQ), physical activity, positive and negative affect, learning goal achievement (LGA), and examination grades. The following analysis focuses on complete observations with values for all variables, so our analysis includes 1788 observations from 62 subjects. One of the major findings from the original study is that better SQ is positively associated with LGA. Hence, we chose SQ as the covariate for both the mean and the scale models. Namely, the location model for the example can be expressed as

$$\text{LGA}_{ij} = \beta_0 + v_i + \beta_{\text{SQ}}\text{SQ}_{ij} + \epsilon_{ij}, \quad (11)$$

and the scale model, in which the variance of  $\epsilon_{ij}$  is being modeled, is

$$\sigma_{\epsilon_{ij}}^2 = \exp(\tau_0 + \omega_i + \tau_{\text{SQ}}\text{SQ}_{ij}), \quad (12)$$

where SQ was measured on a 4-point Likert scale in which 1 means very bad, and 4 means very good. LGA was measured on a 5-point Likert scale in which 0 represents having not achieved the goals at all and 4 represents having achieved the goals completely. We treated LGA and SQ as continuous, consistent with the approach taken in the original study. The approximate normal distribution of LGA has been validated through our exploratory analysis. Nevertheless, it is important for the readers to exercise caution when generalizing the results beyond the original LGA range.

The Cook's distances and COVRATIOs of all subjects are visualized in Fig. 1. After applying the FDR procedure, subjects 7, 12, 49, and 69 are determined to have a statistically significant influence on the model fit based on likelihood-ratio tests. Influence analysis results of these four subjects are summarized in Table 2. According to the Cook's distances, we can see that subjects 7 and 12 have high influence on the fixed scale effects, subject 69 has high influence on the fixed location effects, while subject 49 has high influence on both. In particular, subject 7 has a large DFBETAS for  $\tau_{\text{SQ}}$ , subject 49 has the largest DFBETAS for  $\beta_0$  and the smallest DFBETAS for  $\tau_0$ , and subject 69 has the smallest DFBETAS for  $\beta_0$  and the 2<sup>nd</sup> largest DFBETAS for  $\beta_{\text{SQ}}$ . All four subjects have high influence on the point estimates of the variances and covariance of the random effects. Specifically, separation of subjects 7, 12, and 49 shrinks  $\hat{\sigma}_{\omega}^2$ ; separation of subjects 49 and 69 shrinks  $\hat{\sigma}_v^2$ ; separation of subject 7 and 69 decreases  $\hat{\sigma}_{v\omega}$  while separation of subject 12 and 49 increases  $\hat{\sigma}_{v\omega}$ . Moreover, excluding each of these subjects

from the random effect estimation also shrinks the corresponding generalized variances of the groups of parameters that they have influence on the point estimates, suggesting that these four subjects have caused a loss in model precision in the original model.

We visualize the SQ and LGA of these four subjects in Fig. 2 to see how the data correspond with the influence analysis findings. While all except two of subject 7's self-reported SQ ratings have a value of 4, the LGA of this subject was highly variable when SQ equaled 4. This is consistent with the influence analysis result that giving subject 7 subject-specific fixed effects shrinks  $\hat{\tau}_{\text{SQ}}$  and  $\hat{\sigma}_{\omega}^2$ . While subject 49 had consistently high LGA, subject 69 had all low LGA values at 0 and 1. Such patterns in the data correspond to subject 49's large DFBETAS $_{\beta_0}$  and subject 69's small DFBETAS $_{\beta_0}$ . Given subject 12, 49, and 69 have such great influence on different components of the model and the fact that they have the lowest variability in LGA across all subjects, it could be worthwhile to investigate whether the similar answers for LGA is a result of careless responses or simply outstanding consistency in actual LGA.

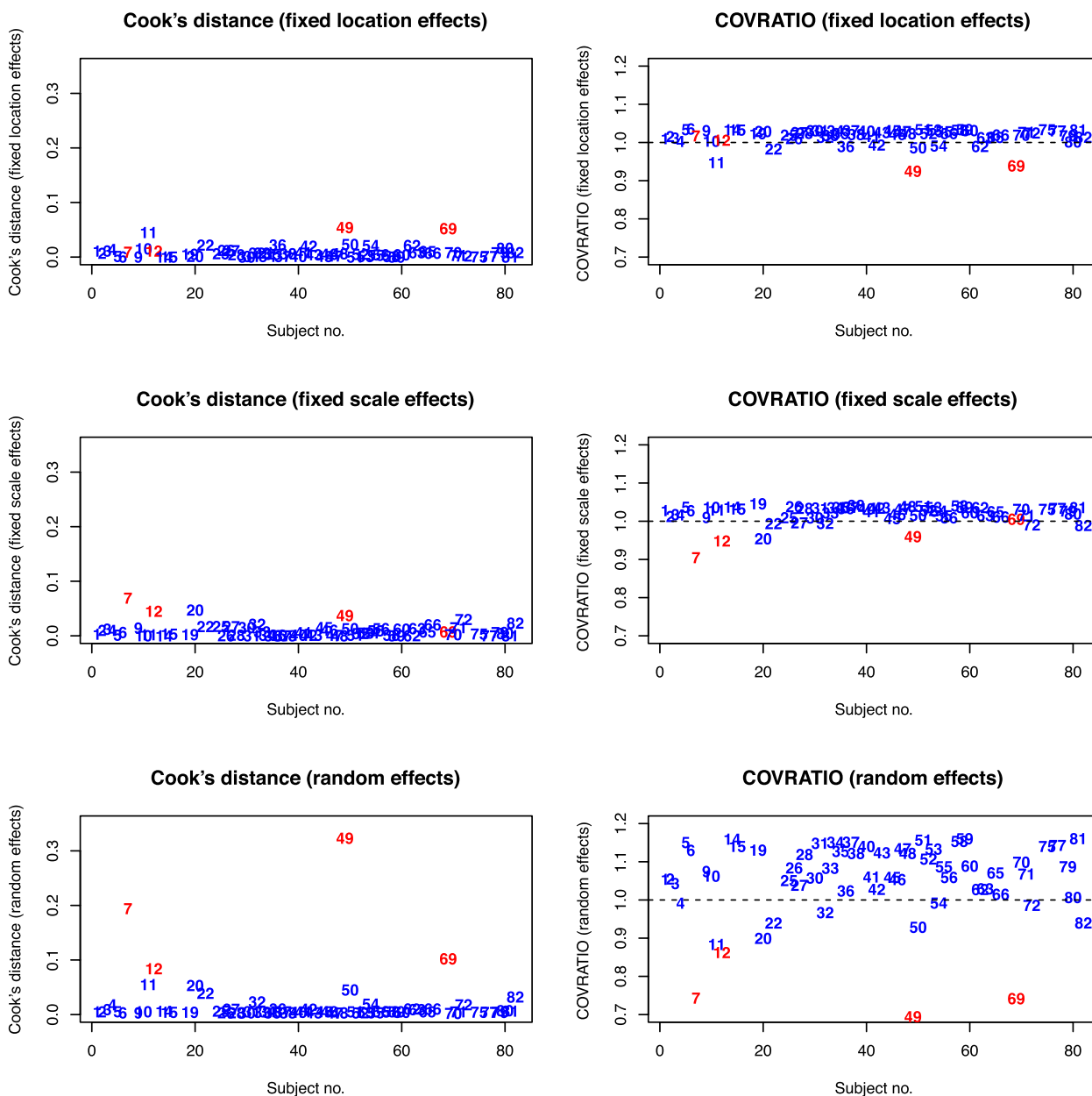
Sensitivity analyses were conducted by excluding the scale model, i.e., using standard MRMs. The sensitivity analysis results still acknowledge the strong influence of subjects 49 and 69 on the location model but do not reveal significant influence from subjects 7 and 12. These findings are consistent with results in Section "Simulation study", that is, subjects with influential data in terms of scale can often be neglected if influence analyses are conducted via MRMs only.

PROC NL MIXED in SAS OnDemand for Academics (SAS Institute Inc.) was used to estimate both the original and the leave-one-out MELS models. All the influence analyses using the results from SAS were carried out in R version 4.2.0 (R Core Team). Both the SAS codes and the R codes we used are included in the [supporting information](#).

### Simulation study

To further illustrate the advantages of conducting influence analyses using MELS models compared to analyses using MRMs with no scale components, we generated simulated examples in two different scenarios, with 500 datasets created for each scenario. Every dataset contains 50 subjects that follow the structure described in Eqs. 1 and 2 before designating some subject(s) to be influential. Both scenarios' responses resemble daily physical activity time in minutes, so negative values were removed. Examples of one regular subject and one artificially influential subject simulated in each scenario are illustrated in Fig. 3.





**Fig. 1** Cook's distances and COVRATIOs for the health behavior data. Points in red are the four subjects influential on the fit of the MELS model of health behavior data

The proposed method requires fitting a MELS model, for which estimation can be time-consuming, for every subject in the dataset. Hence, all the MELS models involved in the simulation study were estimated using a fast estimation algorithm for MELS models, FastRegLS, developed by Gill and Hedeker [20]. The [supporting information](#) contains codes that

utilize FastRegLS to carry out the influence analysis. One or more of the models failed to converge properly for 1 to 18 datasets (mean = 7.67) in the six simulation studies, and these datasets are not included in the result summary. One drawback of the FastRegLS algorithm is that it doesn't offer log-likelihoods. Hence, the simulation studies primarily concentrate on point

**Table 2** Influence analysis results for health behavior data

Subject	Influence measure	Results
7	Cook's distance	- Largest $C^\tau$ (0.069)
		- 2 <sup>nd</sup> largest $C^\eta$ (0.194)
	DFBETAS	- 3 <sup>rd</sup> largest DFBETAS <sup><math>\beta_{SQ}</math></sup> (0.055)
		- Largest DFBETAS <sup><math>\tau_{SQ}</math></sup> (0.163)
		- Largest DFBETAS <sup><math>\sigma_\omega^2</math></sup> (0.604)
	COVRATIO	- Largest DFBETAS <sup><math>\sigma_{vw}</math></sup> (0.325)
- Smallest COVRATIO <sup><math>\tau</math></sup> (0.905)		
12	Cook's distance	- 3 <sup>rd</sup> smallest COVRATIO <sup><math>\eta</math></sup> (0.743)
		- 3 <sup>rd</sup> largest $C^\tau$ (0.045)
	DFBETAS	- 4 <sup>th</sup> largest $C^\eta$ (0.084)
		- 4 <sup>th</sup> largest DFBETAS <sup><math>\beta_0</math></sup> (0.096)
		- 2 <sup>nd</sup> largest DFBETAS <sup><math>\sigma_\omega^2</math></sup> (0.392)
	COVRATIO	- 2 <sup>nd</sup> smallest DFBETAS <sup><math>\sigma_{vw}</math></sup> (-0.343)
- 2 <sup>nd</sup> smallest COVRATIO <sup><math>\tau</math></sup> (0.948)		
49	Cook's distance	- 4 <sup>th</sup> smallest COVRATIO <sup><math>\eta</math></sup> (0.862)
		- Largest $C^\beta$ (0.054)
	DFBETAS	- 4 <sup>th</sup> largest $C^\tau$ (0.037)
		- Largest $C^\eta$ (0.323)
		- Largest DFBETAS <sup><math>\beta_0</math></sup> (0.194)
	COVRATIO	- Smallest DFBETAS <sup><math>\tau_0</math></sup> (-0.140)
- Largest DFBETAS <sup><math>\sigma_v^2</math></sup> (0.553)		
69	Cook's distance	- 4 <sup>th</sup> largest DFBETAS <sup><math>\sigma_\omega^2</math></sup> (0.325)
		- Smallest DFBETAS <sup><math>\sigma_{vw}</math></sup> (-0.737)
	DFBETAS	- Smallest COVRATIO <sup><math>\beta</math></sup> (0.925)
		- 4 <sup>th</sup> smallest COVRATIO <sup><math>\tau</math></sup> (0.959)
		- Smallest COVRATIO <sup><math>\eta</math></sup> (0.694)
	COVRATIO	- 2 <sup>nd</sup> largest $C^\beta$ (0.052)
- 3 <sup>rd</sup> largest $C^\eta$ (0.102)		
69	Cook's distance	- Smallest DFBETAS <sup><math>\beta_0</math></sup> (-0.281)
		- 2 <sup>nd</sup> largest DFBETAS <sup><math>\beta_{SQ}</math></sup> (0.121)
	DFBETAS	- 2 <sup>nd</sup> largest DFBETAS <sup><math>\sigma_\omega^2</math></sup> (0.422)
		- 3 <sup>rd</sup> largest DFBETAS <sup><math>\sigma_{vw}</math></sup> (0.269)
		- 2 <sup>nd</sup> smallest COVRATIO <sup><math>\beta</math></sup> (0.939)
	COVRATIO	- 2 <sup>nd</sup> smallest COVRATIO <sup><math>\eta</math></sup> (0.741)

estimates, along with the variances and covariances of parameter estimates.

**Simulation scenario 1**

In the first scenario, the time-varying covariate, portions of fiber intake, is continuous and was simulated based on a  $\mathcal{N}(3, 1)$  distribution, after which absolute values were taken to ensure non-negativity. The values of the parameters used to generate the simulated datasets are as

follows:  $\beta = [100, 10]$ ;  $\tau = [7, 0.5]$ ;  $\eta = [20.09, 1, 0.45]$ . In other words, the model used to generate data except for the influential one(s) is as follows:

$$\begin{aligned}
 y_{ij} &= 100 + v_i + 10x_{1ij} + \epsilon_{ij}, \\
 \sigma_{\epsilon_{ij}}^2 &= \exp(7 + \omega_i + 0.5x_{1ij}), \\
 \begin{pmatrix} v_i \\ \omega_i \end{pmatrix} &\sim \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 20.09 & 0.45 \\ 0.45 & 1 \end{pmatrix}\right).
 \end{aligned}
 \tag{13}$$

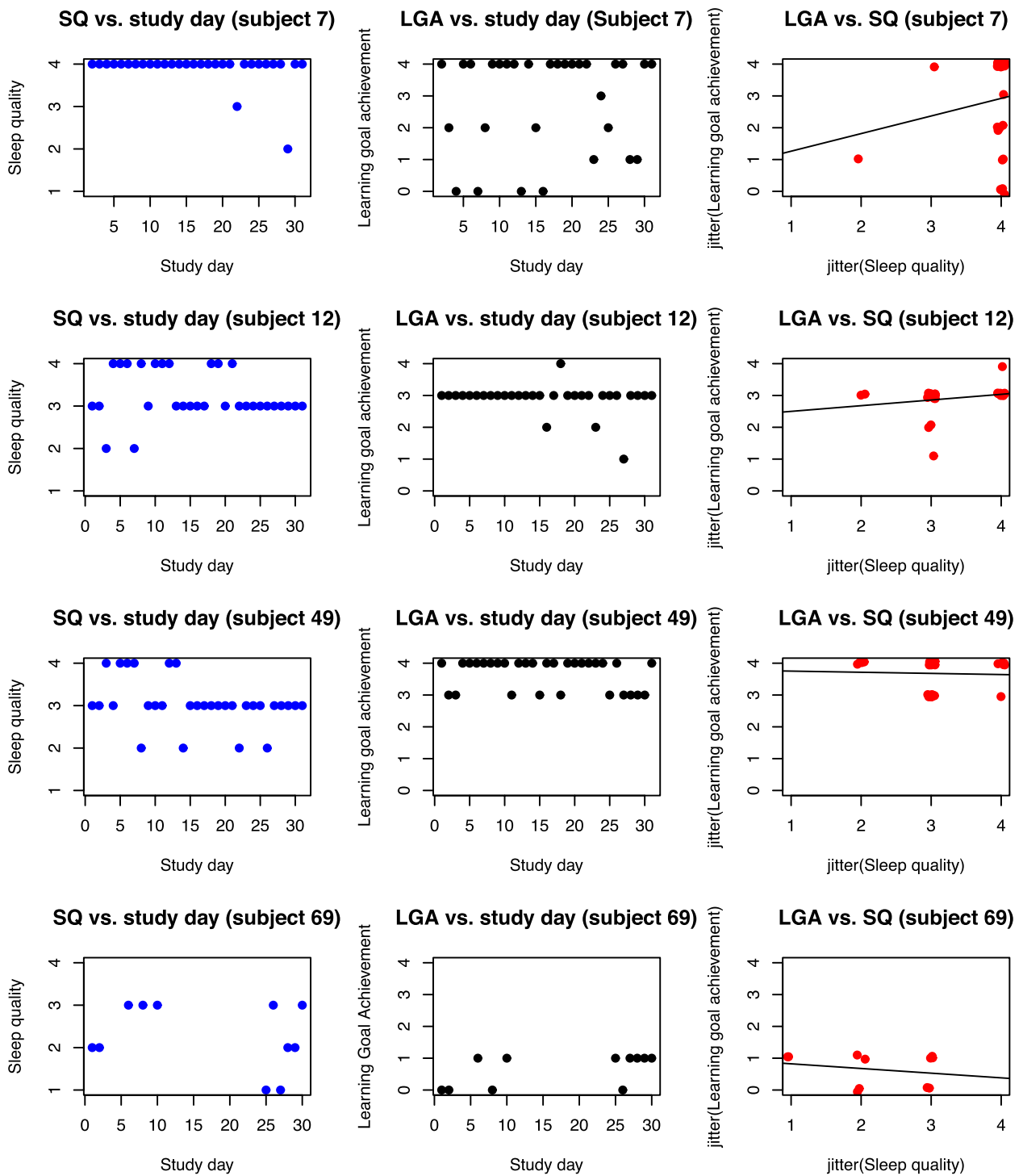
There are 500 observations from each subject, and the covariate is standardized for influence analysis to facilitate model estimation. This scenario focuses on the case of careless responses, so we first designated subject 1 to have all responses at 150 and 155 and all covariates at 3 and 4 before standardization. The percentage of simulations in which each influence measure detects subject 1 to be the most influential are summarized in the first column of Table 3. In summary, subject 1 always has the smallest COVRATIO <sup>$\tau$</sup>  and the largest DFBETAS <sup>$\sigma_\omega^2$</sup> . It also almost always has the largest DFBETAS <sup>$\tau_0$</sup> . The second column in Table 3 shows the percentage of detection when we specify the first three subjects to have such careless responses, and the three subjects still almost always have the smallest COVRATIO <sup>$\tau$</sup> 's and the largest DFBETAS <sup>$\sigma_\omega^2$</sup> 's. Note that only all three influential subjects among the top three are counted as one detection.

**Simulation scenario 2**

For scenario 2, the covariate values were simulated to represent discrete time periods with values that range from 0 to 7, and each subject had 25 observations from each time period before removing negative values. The simulated datasets were generated according to the following structure:

$$\begin{aligned}
 y_{ij} &= 100 + v_i + 15x_{1ij} + \epsilon_{ij}, \\
 \sigma_{\epsilon_{ij}}^2 &= \exp(3 + \omega_i + 0.6x_{1ij}), \\
 \begin{pmatrix} v_i \\ \omega_i \end{pmatrix} &\sim \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 20.09 & 0.22 \\ 0.22 & 0.25 \end{pmatrix}\right).
 \end{aligned}
 \tag{14}$$

The influential subject(s) are subject(s) with exceptional consistency throughout the behavioral study and are simulated with an intercept of 0 in the scale model. In the simulation study with the first subject like this, the subject almost always has the largest  $C^\tau$ , the smallest COVRATIO <sup>$\tau$</sup> , the largest DFBETAS <sup>$\tau_0$</sup> , and the largest DFBETAS <sup>$\sigma_\omega^2$</sup> . The same results remain in terms of COVRATIO <sup>$\tau$</sup>  and DFBETAS <sup>$\sigma_\omega^2$</sup>  when the number of artificial influential subjects increases to three.

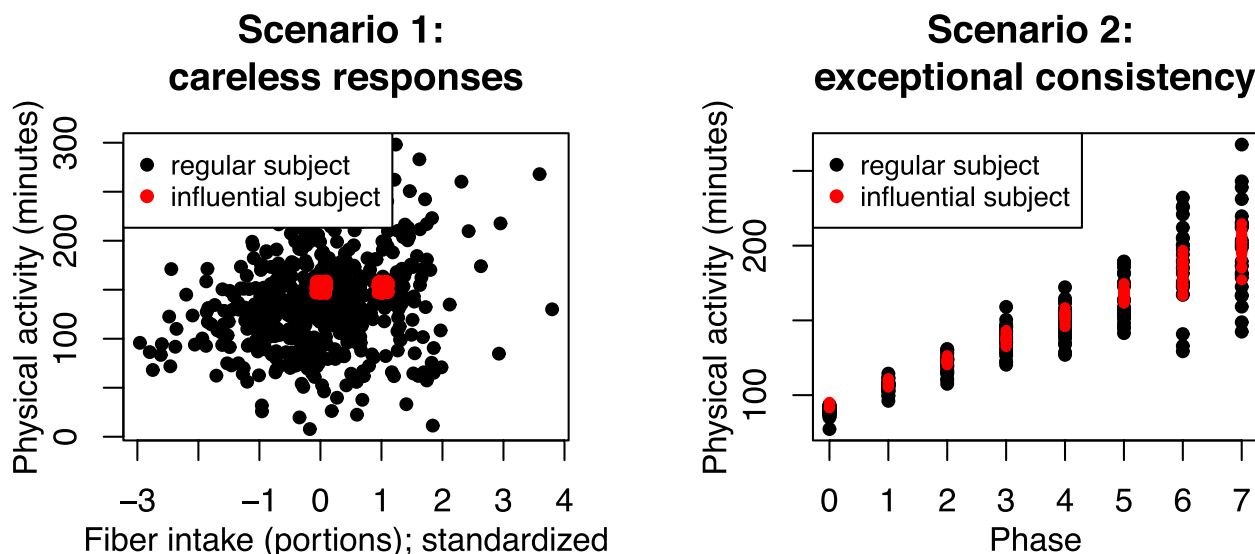


**Fig. 2** Data from subjects influential on the fit of the MELS model of health behavior data

As shown in the third column of Table 3, using standard MRMs, none of the influence measures recognize the influences of the designated influential subjects by more than

half of the time in either scenario. Therefore, in order to successfully identify an influential subject in these scenarios, it is essential to take the scale model into consideration.





**Fig. 3** Simulated data examples. Points in black represent data from one regular subject, and points in red represent data from one influential subject

**Discussion**

In this paper, we have discussed procedures to detect subject(s) influential in model fit and parameter estimates in MELS models. This approach allows researchers to identify subjects influential on the scale structure of the outcome being modeled in addition to the location structure. We hope that our method is able to help researchers, especially researchers interested in studying intra-individual variability, better identify interesting or troublesome subjects that they want to study further in an EMA study or a large-scale longitudinal clinical trial. After determining which subjects are considered influential, researchers can keep these subjects in the analysis using our leave-one-out MELS model described in Section “Influence analysis”, and a summary of other common ways of dealing with influential data can be found in a recent work by Aguinis et al. [21]. The proposed method can also benefit analyses not carried out in MELS models by providing researchers with a better understanding of both the location and the scale structures of their data during exploratory analyses.

Our study has focused on subject-level influence analyses, so one possible extension is the detection of influential observations in MELS model. Given that there could be influential observation(s) within a non-influential subject, estimating a separate model for each observation might be required. Because of the enormous number of observations in intensive longitudinal data, such methods can be extremely computationally intensive.

This article has also focused on maximum likelihood estimates of 2-level MELS models for normally distributed continuous outcomes. Further development can extend the proposed framework to accommodate models with 3 or more levels [22, 23], models on outcomes with more complicated structures [24], and Bayesian estimation approaches [25]. Future work will also extend to ordinal MELS models [26].

To plan for further data analysis based on the influence analysis results, we recommend readers go through all the results from our framework and use their domain knowledge to decide whether specific subject(s) are considered influential and need further analysis. However, we understand that one might want cutoffs to guide their judgment. Rule-of-thumb cut-off values are  $\frac{4}{N}$  for Cook’s distances,  $1 \pm 3(\frac{r_y}{N})$  for COVRATIO, and  $\frac{2}{\sqrt{N}}$  for DFBE-TAS [17]. Also, a possible future direction of this work is to improve existing cut-off values to be more suitable for MELS models.

The simulation studies have showcased that at least some of the influence measures are able to capture all influential subjects in the case that multiple of them coexist. Nevertheless, not all measures perform the same, which might be attributed to the masking effect [27]. Therefore, we again recommend readers carefully examine all influence measures mentioned in the framework, and a future step of this study will be improving individual influence measures to overcome any possible masking effect.

**Table 3** Percentage of detection using MELS models and MRMs and different influence measures. For scenarios with the first subject as the influential subject, "large(%)" represents the percentage of simulations in which this subject has the largest value of the respective influence measure among all subjects, and "small(%)" represents the percentage of simulations in which this subject has the smallest value of the respective influence measure among all subjects. For scenarios with the first three subjects as the influential subjects, "large(%)" represents the percentage of simulations in which all three of these subjects have the largest values of the respective influence measure among all subjects, and "small(%)" represents the percentage of simulations in which all these three subjects have the smallest values of the respective influence measure among all subjects. Results in the column named "MRMs" were obtained through analyses of using MRMs with no scale components. The simulated datasets used in the MELS model and MRM analyses are the same in each scenario

	MELS models		MRMs
	Single	Multiple (3)	Single
<b>Scenario 1</b>			
$C^\tau$ [large(%)]	75.9	16.7	-
$COVRATIO^\tau$ [large(%)/small(%)]	0/100	0/99.6	-
$DFBETAS^{\tau_0}$ [large(%)/small(%)]	99.4/0	25.4/0.4	-
$DFBETAS^{\tau_1}$ [large(%)/small(%)]	5.4/34.6	1.4/15.7	-
$DFBETAS^{\sigma_\omega^2}$ [large(%)/small(%)]	100/0	99.6/0	-
$C^\beta$ [large(%)]	3.5	5.8	0
$COVRATIO^\beta$ [large(%)/small(%)]	15.1/0.8	2.0/0.8	41.6/0
$DFBETAS^{\beta_0}$ [large(%)/small(%)]	0.8/1.7	0.6/0.6	0/0
$DFBETAS^{\beta_1}$ [large(%)/small(%)]	21.0/44.4	3.6/22.7	6.8/13.9
$DFBETAS^{\sigma_\nu^2}$ [large(%)/small(%)]	0.6/0	1.0/0	0/20.1
<b>Scenario 2</b>			
$C^\tau$ [large(%)]	99.8	62.3	-
$COVRATIO^\tau$ [large(%)/small(%)]	0/100	0/99.8	-
$DFBETAS^{\tau_0}$ [large(%)/small(%)]	99.8/0	69.9/0	-
$DFBETAS^{\tau_1}$ [large(%)/small(%)]	11.3/16.0	0.4/0.8	-
$DFBETAS^{\sigma_\omega^2}$ [large(%)/small(%)]	100/0	100/0	-
$C^\beta$ [large(%)]	21.9	0	1.2
$COVRATIO^\beta$ [large(%)/small(%)]	0/2.6	0/0	2.0/1.2
$DFBETAS^{\beta_0}$ [large(%)/small(%)]	1.4/2.0	0/0	1.4/1.6
$DFBETAS^{\beta_1}$ [large(%)/small(%)]	38.9/44.9	0/0.8	1.4/0.6
$DFBETAS^{\sigma_\nu^2}$ [large(%)/small(%)]	1.4/6.3	0/0	1.2/2.0

**Conclusion**

The proposed influence analysis framework using MELS models enables detection of influential subjects on the scale structure and/or location structure of intensive longitudinal data. Thus, it can facilitate modeling that accounts for the abnormality of certain subject(s). Such benefits of the proposed methods are revealed in both the real-life example and the simulated examples.

**Abbreviations**

MELS	Mixed-effects location scale
MRMs	Mixed-effects regression models
BS	Between-subject
WS	Within-subject
EMA	Ecological momentary assessments
FDR	False discovery rate
SQ	Sleep quality
LGA	Learning goal achievement

**Supplementary Information**

The online version contains supplementary material available at <https://doi.org/10.1186/s12874-023-02046-9>.

Additional file 1.

**Acknowledgements**

An earlier version of the abstract of this manuscript has been previously presented and published at the 2023 International Conference on Health Policy Statistics and the 2023 Midwest Biopharmaceutical Statistics Workshop. This work utilized the computational resources of the Center for Research Informatics' Gardner HPC cluster at the University of Chicago (<http://cri.uchicago.edu>).

**Authors' contributions**

DH and XZ developed the influential analysis framework and contributed to writing the manuscript. XZ analyzed and interpreted the health behavior data and simulated data. All authors read and approved the final manuscript.

**Funding**

This project was supported by National Institute of Diabetes and Digestive and Kidney Diseases (grant number R01 DK125414). The content of this study is solely the responsibility of the authors and does not represent the official views of National Institute of Diabetes and Digestive and Kidney Diseases.

**Availability of data and materials**

The dataset used in the health behavior study example are publicly available in Harvard Dataverse, <https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/27388>. Code used to generate data in the simulated study is included in the [Supporting Information](#).

**Declarations**

**Ethics approval and consent to participate**

Not applicable.

**Consent for publication**

Not applicable.

**Competing interests**

The authors declare no competing interests.

Received: 7 July 2023 Accepted: 26 September 2023

Published online: 18 October 2023

**References**

- Tudor-Locke C, Ham SA, Macera CA, Ainsworth BE, Kirtland KA, Reis JP, et al. Descriptive epidemiology of pedometer-determined physical activity. *Med Sci Sports Exerc.* 2004;36(9):1567–73.
- MacDonald SW, Nyberg L, Bäckman L. Intra-individual variability in behavior: links to brain structure, neurotransmission and neuronal activity. *Trends Neurosci.* 2006;29(8):474–80.
- Yu L, Wang T, Wilson RS, Leurgans S, Schneider JA, Bennett DA, et al. Common age-related neuropathologies and yearly variability in cognition. *Ann Clin Transl Neurol.* 2019;6(11):2140–9.

4. Hedeker D, Mermelstein RJ, Demirtas H. An application of a mixed-effects location scale model for analysis of ecological momentary assessment (EMA) data. *Biometrics*. 2008;64(2):627–34.
5. Nordgren R, Hedeker D, Dunton G, Yang CH. Extending the mixed-effects model to consider within-subject variance for Ecological Momentary Assessment data. *Stat Med*. 2020;39(5):577–90.
6. Meade AW, Craig SB. Identifying careless responses in survey data. *Psychol Methods*. 2012;17(3):437.
7. Shi L, Chen G. Detection of outliers in multilevel models. *J Stat Plan Inference*. 2008;138(10):3189–99.
8. Sharifi M, Marshall G, Goldman R, Rifas-Shiman SL, Horan CM, Koziol R, et al. Exploring innovative approaches and patient-centered outcomes from positive outliers in childhood obesity. *Acad Pediatr*. 2014;14(6):646–55.
9. Langford IH, Lewis T. Outliers in multilevel data. *J R Stat Soc Ser A Stat Soc*. 1998;161(2):121–60.
10. Christensen R, Pearson LM, Johnson W. Case-deletion diagnostics for mixed models. *Technometrics*. 1992;34(1):38–45.
11. Nieuwenhuis R, Te Grotenhuis H, Pelzer B. Influence. ME: tools for detecting influential data in mixed effects models. *R J*. 2012;4(2):38–47.
12. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J R Stat Soc Ser B Methodol*. 1995;57(1):289–300.
13. Verhoeven KJ, Simonsen KL, McIntyre LM. Implementing false discovery rate control: increasing your power. *Oikos*. 2005;108(3):643–7.
14. Cook RD. Detection of influential observation in linear regression. *Technometrics*. 1977;19(1):15–8.
15. Snijders TA, Bosker RJ. *Multilevel analysis: An introduction to basic and advanced multilevel modeling*. Thousand Oaks: Sage; 2011.
16. Fox J, Monette G. *An R and S-Plus companion to applied regression*. Thousand Oaks: Sage; 2002.
17. Belsley DA, Kuh E, Welsch RE. *Regression diagnostics: Identifying influential data and sources of collinearity*. New York: John Wiley & Sons; 2005.
18. Wilks SS. Certain generalizations in the analysis of variance. *Biometrika*. 1932;24(3/4):471–94.
19. Flueckiger L, Lieb R, Meyer AH, Mata J. How health behaviors relate to academic performance via affect: An intensive longitudinal study. *PLoS ONE*. 2014;9(10):e111080.
20. Gill N, Hedeker D. Fast estimation of mixed-effects location-scale regression models. *Stat Med*. 2023;42(9):1430–44.
21. Aguinis H, Gottfredson RK, Joo H. Best-practice recommendations for defining, identifying, and handling outliers. *Organ Res Methods*. 2013;16(2):270–301.
22. Li X, Hedeker D. A three-level mixed-effects location scale model with an application to ecological momentary assessment data. *Stat Med*. 2012;31(26):3192–210.
23. Lin X, Mermelstein RJ, Hedeker D. A 3-level Bayesian mixed effects location scale model with an application to ecological momentary assessment data. *Stat Med*. 2018;37(13):2108–19.
24. Blozis SA, McTernan M, Harring JR, Zheng Q. Two-part mixed-effects location scale models. *Behav Res Methods*. 2020;52(5):1836–47.
25. Rast P, Hofer SM, Sparks C. Modeling individual differences in within-person variation of negative and positive affect in a mixed effects location scale model using BUGS/JAGS. *Multivar Behav Res*. 2012;47(2):177–200.
26. Hedeker D, Demirtas H, Mermelstein RJ. A mixed ordinal location scale model for analysis of Ecological Momentary Assessment (EMA) data. *Stat Interface*. 2009;2(4):391.
27. Atkinson A. Masking unmasked. *Biometrika*. 1986;73(3):533–41.

## Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Ready to submit your research? Choose BMC and benefit from:**

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

**At BMC, research is always in progress.**

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

