

# Advances in statistical inference for longitudinal generalized linear models

## flipscores Two-Stage Summary Statistics

Angela Andreella  
[angela.andreella@unitn.it](mailto:angela.andreella@unitn.it)

BEYOND Summer School

May 29, 2025

# So ...

We introduce some concepts about GLMMs and GEEs. We have seen that both have advantages and disadvantages. Let's focus on the disadvantages again 😊:

## GLMM:

- Uncertainty about how to specify the random effects structure,
- Invalid inference under model misspecification (both random effects and fixed effects structure),
- Convergence issues.


## GEE:

- Problems with small sample sizes,
- Difficulties with unbalanced designs,
- Presence of endogenous covariates — i.e., variables correlated with the error term — when assuming independence,
- Non-randomly missing data.



## APPLICATION REVIEWS AND CASE STUDIES

# Robust Inference for Generalized Linear Mixed Models: A “Two-Stage Summary Statistics” Approach Based on Score Sign Flipping

Angela Andreella<sup>1</sup> , Jelle Goeman<sup>2</sup>, Jesse Hemerik<sup>3</sup> and Livio Finos<sup>4</sup>

<sup>1</sup>University of Trento, Italy; <sup>2</sup>Leiden University Medical Center, The Netherlands; <sup>3</sup>Erasmus University Rotterdam, The Netherlands; <sup>4</sup>University of Padova, Italy

**Corresponding author:** Angela Andreella; Email: [angela.andreella@unitn.it](mailto:angela.andreella@unitn.it)

(Received 24 October 2024; accepted 11 December 2024)

This manuscript is part of the special section *Model Identification and Estimation for Longitudinal Data in Practice*. We thank Drs. Carolyn J. Anderson and Donald Hedeker for serving as co-Guest Editors.

Ingredients for the **flip2sss** (flipscores two-stage summary statistics):

- Permutation theory (flip)
- Score test (scores)
- Two-stage summary statistics approach (2sss)

The proposed method is efficiently implemented in the R package `jointest`, which is compatible with large datasets and complex statistical models.

## flipscores (brief recap)

Consider  $n$  independent observations  $y_1, \dots, y_n$  following a GLM from the exponential dispersion family:

$$f(y_i; \theta_i; \phi_i) = \exp \left\{ \frac{y_i \theta_i - b(\theta_i)}{a(\phi_i)} + c(y_i, \phi_i) \right\}$$

where  $\theta_i$  is the canonical parameter,  $\phi_i$  is the dispersion parameter, and  $\mu_i = \mathbb{E}(y_i) = b'(\theta_i)$ . The GLM is then defined as:

$$\boxed{g(\boldsymbol{\mu}) = \boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta}}$$

where  $g(\cdot)$  is the link function,  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)^\top$ ,  $\mathbf{X}$  is the design matrix, and  $\boldsymbol{\beta}$  is the vector of  $q$  parameters.

## flipscores (brief recap)

The main focus is to test the null hypothesis for a given element  $d$  of  $\beta$ , while still accounting for all nuisance parameters, i.e.,

$$H_0 : \beta_d = \beta_0 \mid \beta_1, \dots, \beta_{d-1}, \beta_{d+1}, \dots, \beta_q, \phi_1, \dots, \phi_n.$$

The *effective score* is given by:

$$S = n^{-1/2} \mathbf{X}_d^\top \mathbf{W}^{1/2} (\mathbf{I} - \mathbf{H}) \mathbf{V}^{1/2} (\mathbf{y} - \hat{\boldsymbol{\mu}}) \quad (1)$$

where:

- $\mathbf{V} = \text{diag}\{\text{Var}(y_i)\}$ ,
- $\mathbf{W} = \text{diag}\left\{\frac{\partial \mu_i}{\partial \eta_i}\right\} \mathbf{V}^{-1} \text{diag}\left\{\frac{\partial \mu_i}{\partial \eta_i}\right\}$ ,
- $\mathbf{H} = \mathbf{W}^{1/2} \mathbf{X}_{-d} (\mathbf{X}_{-d}^\top \mathbf{W} \mathbf{X}_{-d})^{-1} \mathbf{X}_{-d}^\top \mathbf{W}^{1/2}$ ,
- $\hat{\boldsymbol{\mu}} = (\hat{\mu}_1, \dots, \hat{\mu}_n)$  are the fitted values of the model under  $H_0$ .

## flipscores (brief recap)

To improve small-sample reliability, we consider the standardized version:

$$S^{\star} = \frac{S}{\text{var}\{S\}^{1/2}},$$

where

$$\text{Var}\{S\} = n^{-1} \mathbf{X}_d^{\top} \mathbf{W}^{1/2} (\mathbf{I} - \mathbf{H}) \mathbf{W}^{1/2} \mathbf{X}_d + o_p(1).$$

$S^{\star}$  is asymptotically **valid under any variance misspecification**.

The statistic  $S$  in Equation (1) can also be rephrased as a sum of  $n$  components: the **effective score contributions**.

## flipscores (brief recap)

The  $p$ -value is computed by randomly flipping the signs of these score contributions.  $\Rightarrow$  matrix  $\mathbf{F}$  with diagonal elements  $-1$  or  $1$  sampled with equal probability.

Sign-flipping the score contribution: multiplying the effective score by  $\mathbf{F}$ :

$$S(\mathbf{F}) = n^{-1/2} \mathbf{X}_d^\top \mathbf{W}^{1/2} (\mathbf{I} - \mathbf{H}) \mathbf{V}^{-1/2} \mathbf{F} (\mathbf{y} - \hat{\boldsymbol{\mu}}).$$

The corresponding standardized version is

$$S^*(\mathbf{F}) = \frac{S(\mathbf{F})}{\text{var}\{S(\mathbf{F})\}^{1/2}},$$

where

$$\text{Var}\{S(\mathbf{F})\} = n^{-1} \mathbf{X}_d^\top \mathbf{W}^{1/2} (\mathbf{I} - \mathbf{H}) \mathbf{F} (\mathbf{I} - \mathbf{H}) \mathbf{F} (\mathbf{I} - \mathbf{H}) \mathbf{W}^{1/2} \mathbf{X}_d + o_p(1).$$



## flipscores (brief recap)

Considering  $B$  independent **sign flip transformations**, where  $S_1^* = S^*(\mathbf{I})$  is the observed test statistic, we reject the null hypothesis  $H_0 : \beta_d = \beta_0$  versus the alternative  $H_1 : \beta_d > \beta_0$  at significance level  $\alpha$  if:

$$S_1^* > S_{(\lceil (1-\alpha)B \rceil)}^*,$$

where  $S_{(1)}^* \leq S_{(2)}^* \leq \dots \leq S_{(B)}^*$  are the sorted statistics and  $\lceil \cdot \rceil$  is the ceiling function.

In the same way, we reject  $H_0$  versus  $H_1 : \beta_d < \beta_0$  if:

$$S_1^* < S_{(\lceil \alpha B \rceil)}^*,$$

and versus  $H_1 : \beta_d \neq \beta_0$  if:

$$S_1^* < S_{(\lceil (\alpha/2)B \rceil)}^* \cup S_1^* > S_{(\lceil (1-\alpha/2)B \rceil)}^*.$$

# adolong data

Adopted children measured longitudinally over Time:

- Health issues: response variable `Unhealth` (binary variable)
- Sex (binary variable)
- Age (age of the child when they arrived in the family, in years)
- Country (country of origin)

A possible model of interest is:

$$\text{Unhealth} \sim 1 + \text{Country} + \text{Age} + \text{Sex} + \text{Time} + \text{Sex:Time}.$$

Consider again  $n$  observations  $y_1, \dots, y_n$  and  $n_j$  observations in cluster (i.e., child)  $j$ , where  $n = \sum_j^N n_j$  and  $N$  is the total number of clusters (i.e., children).

The **exchangeability assumption** used to compute the null distribution of the standardized score test statistic  $S^*$  does not hold  
⇒ **"two-stage summary statistics" approach.**

In short, we reduce the hierarchical complexity of the data by computing summary measures at the cluster level in the first stage. These summary measures are then analyzed as a response random variable in the second stage.

Therefore, the main assumption is that we have (at least asymptotically) unbiased estimators of these summary measures.

## First stage

A GLM is fitted separately for each subject  $j$  including only the  $h$  covariates  $\mathbf{K}_{ij} \in \mathbb{R}^{1 \times h}$  that vary within-cluster  $j$ , i.e.,

$$g(\mu_{ij}) = \mathbf{K}_{ij} \boldsymbol{\tau}_j \quad (2)$$

where  $\boldsymbol{\tau}_j \in \mathbb{R}^h$  is the vector of  $h$  parameters for cluster  $j$ .

Fitting the model defined in Equation (2) leads to a vector of estimated parameters  $\hat{\boldsymbol{\tau}}_j \in \mathbb{R}^h$  for each cluster  $j$  corresponding to the design matrix  $\mathbf{K}_j$  that varies within-subject  $j$ .

## First stage

We fit a logistic regression for each child where Time is the only within-cluster covariate. So:

$$\text{logit}[\text{Pr}(\text{Unhealth}_{ij} = 1)] = \mathbf{K}_{ij}\boldsymbol{\tau}_j$$

where  $\mathbf{K}_{ij} \in \mathbb{R}^{1 \times 2}$ :

$$\mathbf{K}_{ij} = \begin{bmatrix} 1 & \text{Time}_{ij} \end{bmatrix}$$

→  $\hat{\boldsymbol{\tau}}_j$  is our “**summary statistics**”: estimated intercept and slope for the variable Time for child  $j$ .

## Second stage

The  $N$  vectors  $\hat{\tau}_j$  are collected in a  $N \times h$  matrix  $\mathbf{T}$  and modeled by a linear model with the between-cluster variables as predictors:

$$\begin{bmatrix} \hat{\tau}_1^\top \\ \vdots \\ \hat{\tau}_N^\top \end{bmatrix} = \mathbf{T} = \mathbf{X}\boldsymbol{\beta} + \mathbf{E} \quad (3)$$

where

- $\mathbf{X} \in \mathbb{R}^{N \times l}$  is the matrix of between-cluster variables,
- $\boldsymbol{\beta} \in \mathbb{R}^{l \times h}$  are the corresponding parameters,
- $\mathbf{E} \in \mathbb{R}^{N \times h}$  is the matrix of errors.

## Second stage

Recalling our initial model:  $\text{Unhealth} \sim 1 + \text{Country} + \text{Age} + \text{Sex} + \text{Time} + \text{Sex}:\text{Time}$

- The intercept of within-cluster effects (i.e., the first column of  $\mathbf{T}$ ) is modeled by the between-cluster effects Intercept, Country, Age, and Sex.
- The estimated slope related to Time (i.e., the second column of  $\mathbf{T}$ ) is modeled only by Intercept and Sex covariates.

Therefore, in the second column of  $\mathbf{T}$ , only those coefficients of Intercept and Sex covariates will be estimated, while the ones associated with Country and Age are set to 0.

# adolong data

## Second stage

$$\beta \in \mathbb{R}^{12 \times 2}$$

$$\beta = \begin{bmatrix} \beta_{\text{Intercept}} & 0 \\ \beta_{\text{Age}} & 0 \\ \beta_{\text{Sex Male}} & 0 \\ \beta_{\text{Country Cambodia}} & 0 \\ \beta_{\text{Country China}} & 0 \\ \beta_{\text{Country Colombia}} & 0 \\ \beta_{\text{Country Ethiopia}} & 0 \\ \beta_{\text{Country India}} & 0 \\ \beta_{\text{Country Thailand}} & 0 \\ \beta_{\text{Country Vietnam}} & 0 \\ 0 & \beta_{\text{Intercept}} \\ 0 & \beta_{\text{Sex}} \end{bmatrix}$$



# adolog data

## Second stage

$\mathbf{X} \in \mathbb{R}^{N \times 12}$ , here we represent the first three child:

Unhealth	Age	Sex	Time	Country	Subj
1	7	Male	3	Ethiopia	1
0	2	Female	6	Vietnam	3
0	8	Male	1	Colombia	5
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$

$$\mathbf{X} = \begin{bmatrix} 1 & 7 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 1 \\ 1 & 2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 0 \\ 1 & 8 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

## Second stage

In summary, the problem is now rewritten as a multiple multivariate linear model where each observation is the **summary statistics**:

$$\hat{\tau}_j \sim (x_j\beta, \psi_j)$$

with

- $x_j$  is the  $j$ -th row of  $\mathbf{X}$
- $\psi_j = \text{Var}(\epsilon_j)$  is the variance of the independent rows of the error matrix  $\mathbf{E} = [\epsilon_1^\top, \dots, \epsilon_N^\top]^\top$

$\psi_j$  can be decomposed as a sum of two sources of variability:

- $\Sigma$ : variance of the random effects  $\Rightarrow$  captures the between-subject variability due to the presence of **random effects**
- $\Sigma_j$ : conditional variance of  $\hat{\tau}_j \Rightarrow$  the second represents the **within-cluster variability**.

Outside the fully balanced designs with homoscedastic errors, the  $\psi_j$  varies among clusters, **making the standard LM tools unreliable**.

The literature on 2sss focuses on the estimate of  $\psi_j$

⇒ making assumptions on the data and being constrained to a precise specification of the random quantities in the model.

⇒ choose whether intercept and slope are random coefficients and also among independent or correlated random effects.

⇒ possible **misspecification!**

➔ ***standardized flipscores** approach in the second stage!*

$H_0 : \beta = \beta_0$  is then tested with the standardized score test, applying the same sign-flipping transformation across all  $h$  models, still avoiding the complexities associated with formulating the random component of the model.

Any **summary measures** at the cluster level can be used in the second stage; the only required property is the (asymptotic) **unbiasedness** of these summary measures.

As an example, in the application and in the simulations, we will adopt and evaluate the use of the Firth correction in the fitting of the cluster-level binomial models; this helps in reducing the finite-sample bias of the maximum likelihood estimates.

Although the unbiasedness property holds, the selection of summary measures is guided by the **researcher's specific objectives**. For instance, in clinical trials, common choices include post-treatment means or mean changes relative to baseline.

In our example, we focus on the effect of time on health conditions.

## Simulated model

$y$  is simulated as a Bernoulli rv with the following mean:

$$\mu_{ij} = \text{logit}^{-1}(X_{ij}\beta + Z_j\gamma + U_j + D_jX_{ij})$$

where

- $Z$  design matrix of the nuisance parameters  $\gamma$  with  $Z_j \sim \mathcal{N}(0, 1)$
- $X$  design matrix of the tested parameters  $\beta$  with  $X_{ij} = O_{ij} + Z_j/2$  where  $O_{ij} \sim \mathcal{N}(0, 1)$
- $U$  defines the subject-specific random effect (i.e., random intercept) where  $U_j \sim \mathcal{N}(0, 0.5)$
- $D$  the subject-specific random slope of  $X$  where  $D_j \sim \mathcal{N}(0, 0.5)$

# Simulations

## Setting

The aim is to test  $H_0 : \beta = 0$ . We compare GLMM, GEE and flip2sss in terms of:

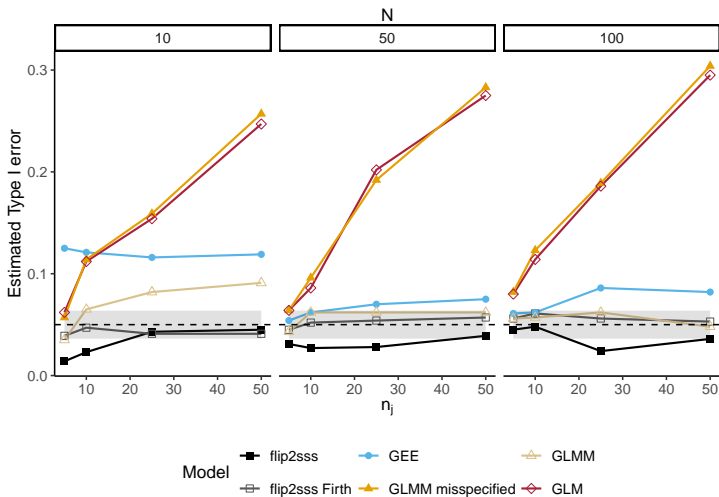
- type I error control  $\Rightarrow \beta = 0, \gamma = 2$
- power  $\Rightarrow \beta = 2, \gamma = 2$

In both scenarios, 1000 simulations and 100 permutations are performed.

- **flip2sss**: using MLE of GLM and Firth correction in the first stage
- **GLMM**: correctly specified (random intercept and slope) and misspecified (only random slope)
- **GEE**: independent working correlation matrix (to assure consistency) which is typically a safe choice

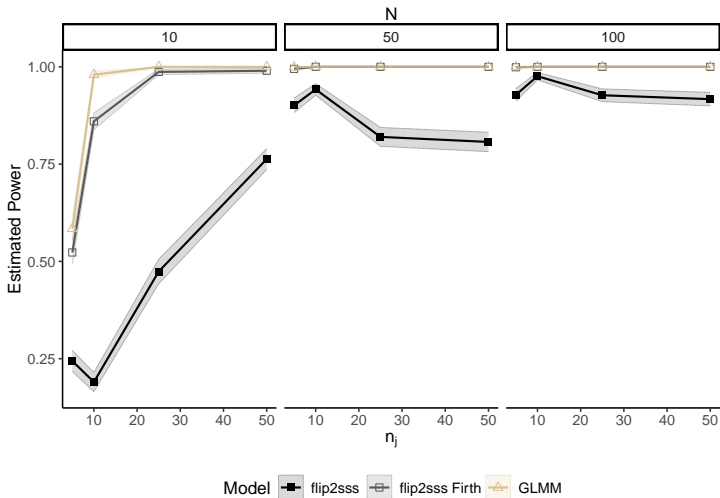
# Simulations

## Balanced data



# Simulations

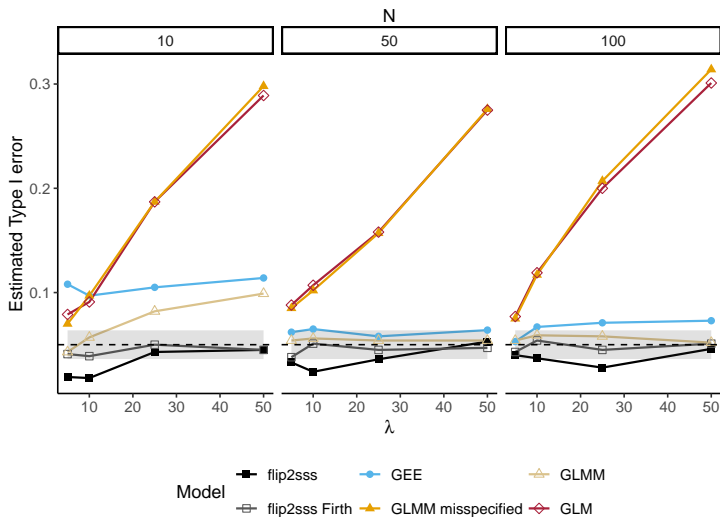
## Balanced data





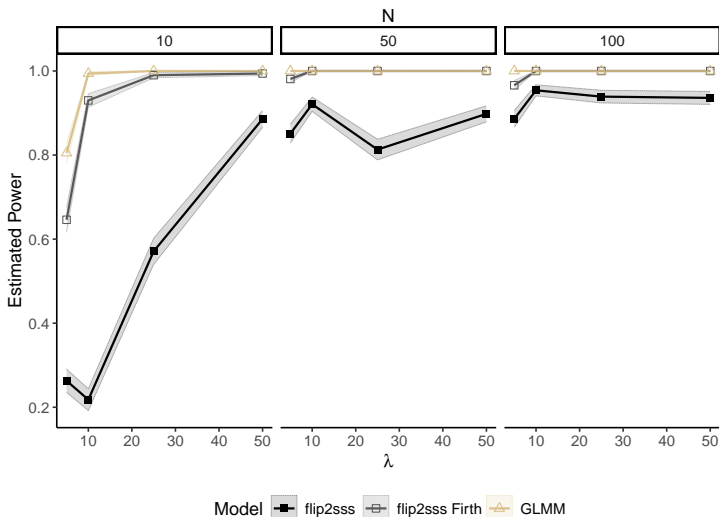
# Simulations

**Unbalanced data:**  $n_j \sim \text{Pois}(\lambda)$  with  $\lambda \in \{5, 10, 25, 50\}$ .



# Simulations

**Unbalanced data:**  $n_j \sim \text{Pois}(\lambda)$  with  $\lambda \in \{5, 10, 25, 50\}$ .



# To sum up

flip2sss, addresses heteroscedasticity, variance misspecification, and within-subject dependence, key aspects in statistical modeling.

★ Key **advantages** include:

- No need to specify random effect structures but only the clusters.
- Flexibility in the first stage with different estimators (e.g., maximum likelihood and Firth correction).
- Easy extension to multivariate cases (multiple dependent variables).

💡 **Limitations** and **future research**:

- Bias in the first stage summary measure estimator impacts statistical power in the second stage.
- Handling more complex correlation structures, such as crossed-random effects.



Let's see it in R!

# Bibliography

- Diggle, P. (2002). Analysis of longitudinal data. Oxford university press.
- Andreella, A., Goeman, J., Hemerik, J., & Finos, L. (2024). Robust inference for Generalized Linear Mixed Models: a “Two-Stage Summary Statistics” approach based on score sign flipping. Psychometrika, 1-42.
- De Santis, R., Goeman, J. J., Davenport, S. J., Hemerik, J., & Finos, L. (2025). Inference in generalized linear models with robustness to misspecified variances. Journal of the American Statistical Association, 1-16.