

# Longitudinal data analysis: Assignment 3

Team B: Kendall Brown *r0773111*  
Stefan Velez *r0924289*

Raïsa Carmen *s0204278*  
Adhithya Unni Narayanan *r0776057*

---

## 1 Introduction

This report assesses the evolution of hearing thresholds over time for a sample of 546 healthy male volunteers. The data originates from the famous Baltimore Longitudinal Study of Aging (BLSA). Previous research showed a change in hearing threshold for all age groups but especially the older population (Brant and Fozard 1990). In this report, we will especially take care of missingness in the data.

First, the *TIME* variable is rounded to the nearest integer value. As such, we aim to balance the dataset with equally-spaced time instances when hearing thresholds are measured.

Hearing thresholds will be explored, both as a continuous variable and as a trichotomized (ordinal) variable with the following three levels:

- $\leq 6$  dB: Excellent hearing
- over 6 and  $\leq 25$  dB: Normal hearing
- $\geq 26$  dB: Hearing loss

### 1.1 Missingness exploration

After discretizing the *TIME* variable, we consider a subject to be missing at a certain time instance if there is no measurement for that subject at that time. It should be noted that, if the subject is not missing (18.16% of TIME-subject instances), we usually (16.88% of TIME-subject instances) have two measurements (one for each ear) at each time instance. In fact, the average number of measurements per subject at each time instance is 0.35, and maximum 4.

Figure 1 was created using the *visdat* package. It shows all subjects, ordered from youngest (in the top) to oldest (in the bottom) and whether or not their data is missing at a certain time instance (on the x-axis). The percentages on top shows the percentage of missingness at each time instance. It is clear from Figure 1 that the missingness is not monotone; subjects may be missing at one time instance and come back later. Since there are too many possible missingness patterns with 23 time instances ( $2^{23}$ ), we do not give an overview of the number of subjects that follow each possible pattern. Instead, figure 2 shows, for each time instance, the number of subjects that:

- are *present*: when the subject's hearing is measured at time  $t$  and  $t - 1$
- are *missing*: when the subject is missing at time  $t$  and  $t - 1$
- *drop out*: when the subject's hearing is measured at time  $t - 1$  but not at time  $t$
- *return*: when the subject's hearing is measured at time  $t$  but not at time  $t - 1$

Table 1: A mixed model to predict missingness.

Variable	Estimate
Intercept	2.22 ***
TIME	0.24 ***
sideright	-0.19 ***
age	0.10
$R_{t-1}$	-2.67 ***
sigma	2.31

Figure 2 clearly shows that subject rarely are measured two years in a row, most are not measured at  $t = 1$ , and the number of subjects that stay missing gradually increases as time passes.

Lastly we explore whether the missingness can be explained by the data by fitting a mixed model to a dataset where  $R_{it}$  is equal to one if the hearing threshold is missing and zero otherwise:

$$\begin{cases} \text{logit}(R_{it}) = \beta_0 + \beta_1 \text{TIME}_{it} + \beta_3 \text{side}_{it} + \beta_4 \text{age}_{it} + \beta_5 R_{it-1} + b_i \\ b_i \sim N(0, \sigma^2) \end{cases} \quad (1)$$

The variable *age* was standardized to get convergence in the model. Table 1 shows that TIME is significant; as time increases, subjects are more likely to be missing. We can therefore assume missingness at random (MAR). Left ear measurements are also more likely to be missing. A subject is also less likely to be missing at time  $t$  if he was missing at time  $t - 1$ . This can be seen especially in the first couple of years in figure 1: all but 3 subjects are measured at  $t = 0$ , almost no-one is measured at  $t = 1$  and many are measured again at  $t = 2$ .

## 2 Methodology

First, a direct likelihood analysis is compared with multiple imputation in the ?? continuous/discrete??? case. Next, weighted generalized estimating equations are compared with ‘multiple-imputation generalized estimating equations’. Lastly, a sensitivity analysis is performed.

For imputation, the *mice* library is used (Buuren and Groothuis-Oudshoorn 2011) and different imputation techniques were tested: Predictive mean matching, Bayesian linear regression, Unconditional mean imputation, and imputation by random forests.

All analysis was done in R. All scripts are freely available at this git repository.

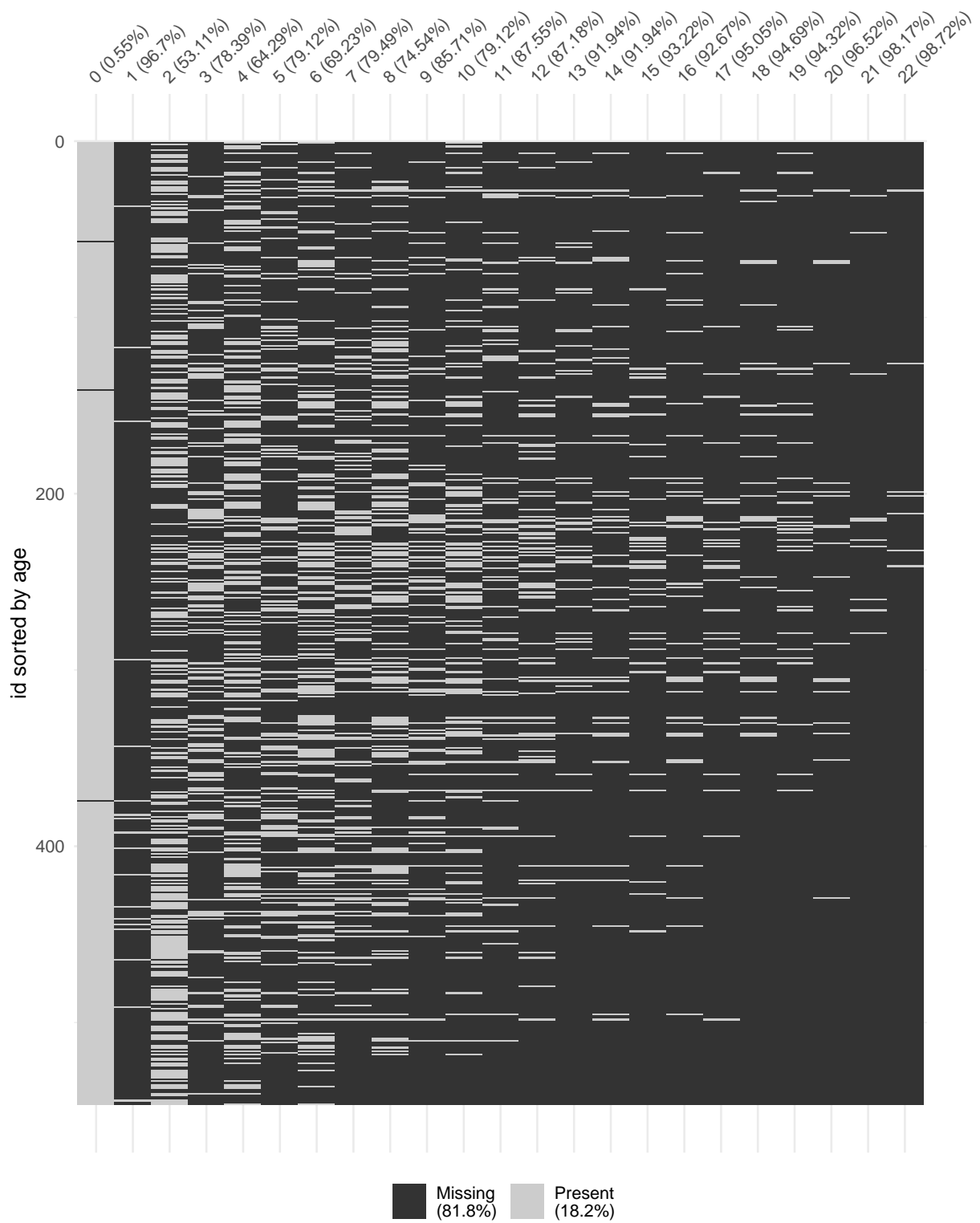


Figure 1: Visual inspection of missingness for different ages at different time instances.

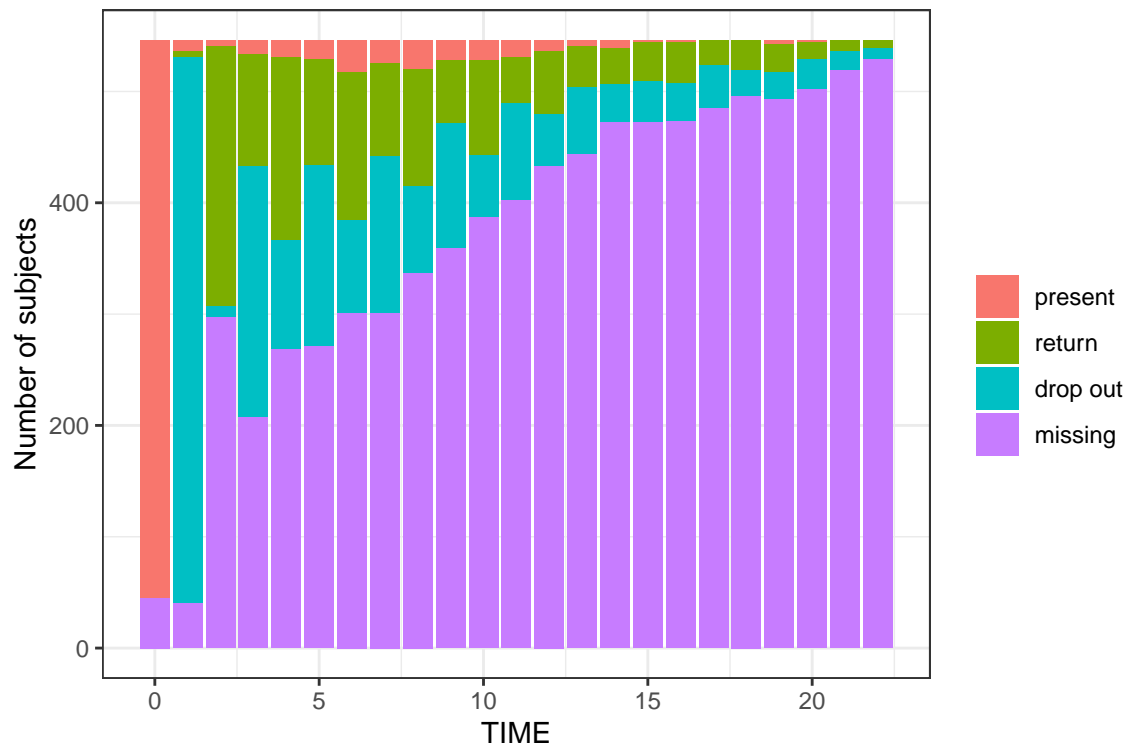


Figure 2: Number of subjects the are present, return, drop out or are missing at each time instance.

### 3 Results

#### 3.1 Direct likelihood analysis versus multiple imputation

Q4

##### 3.1.1 Direct likelihood

##### 3.1.2 Multiple imputation

#### 3.2 Weighted generalized estimating equations versus ‘multiple-imputation generalized estimating equations’

Using frequentist methods, Q5

#### 3.3 Sensitivity analysis

Q6

### Bibliography

Brant, Larry J., and James L. Fozard. 1990. “Age Changes in Pure-tone Hearing Thresholds in a Longitudinal Study of Normal Human Aging.” *The Journal of the Acoustical Society of America* 88 (2): 813–20. <https://doi.org/10.1121/1.399731>.

Buuren, Stef van, and Karin Groothuis-Oudshoorn. 2011. "Mice: Multivariate Imputation by Chained Equations in r." *Journal of Statistical Software* 45 (3): 1–67. <https://doi.org/10.18637/jss.v045.i03>.