

Speech intelligibility measurement

A latent variable approach

Jose Rivera¹, Sven de Maeyer², and Steven Gillis³

¹ Department of Training and Education Sciences,
University of Antwerp, Antwerp, Belgium
E-mail: JoseManuel.RiveraEspejo@uantwerpen.be

(corresponding author)

² Department of Training and Education Sciences,
University of Antwerp, Antwerp, Belgium
E-mail: sven.demaeyer@uantwerpen.be

³ Computational Linguistics, & Psycholinguistics Research Centre
University of Antwerp, Antwerp, Belgium
E-mail: steven.gillis@uantwerpen.be

May 14, 2022

Contents

1. Introduction	5
2. Materials and Methods	7
2.1. Children	7
2.2. Stimuli	8
2.3. Experimental setup	8
3. Results	10
3.1. Data collection	10
4. Discussion	11
5. Author contributions	11
6. Financial support	11
7. Conflicts of interest	11
8. Research transparency and reproducibility	12
A. Supplementary	13
A.1. About Speech Intelligibility	13
A.2. Sampling bias	13
A.3. Children characteristics	13
A.4. DAG: factors influencing Intelligibility	13
A.5. Model details	15
A.6. Simulation	19
A.7. Model selection	19
Bibliography	23

List of Figures

1.	Slider for the HJ task.	9
2.	22
3.	22
4.	DAG for the hierarchical beta regression model for entropy.	22

List of Tables

1. Characteristics of selected children. 21

Abstract

1. Introduction

Intelligible speech can be defined as the extent to which the elements in an speaker’s acoustic signal, e.g. phonemes or words, can be correctly recovered by a listener [13, 25, 22, 9]. Because intelligible spoken language requires all core components of speech perception, cognitive processing, linguistic knowledge, and articulation to be mastered [9], its attainment carries an important societal value, as it is a milestone in children’s language development, the ultimate checkpoint for the success of speech therapy, and has been qualified as the “gold standard for assessing the benefit of cochlear implantation” [4].

For instance, in a transcription task, intelligibility refers to the extent to which a transcriber can identify the words contained in an utterance

The literature suggest there are two perspectives from which speech intelligibility (SI) can be assessed: the message and listener’s perspective [1, 2]. The first, also known as acoustic studies, center its focus on assessing separately particular characteristics of speech samples, e.g. their pitch, duration or stress (supra segmental characteristics), or the articulation of vowels and consonants (segmental characteristics) [20]. Whereas the second, also known as perceptual studies, center its focus on making holistic assessments of the speech stimuli, e.g. measure their perceived quality [1, 2]. On both instances, the stimuli or representations (children’s utterances) can be generated from reading at loud, contextualized utterances, or spontaneous speech tasks¹.

Moreover, perceptual studies can use multiple approaches to measure SI. However, they can be largely grouped into two: objective and subjective ratings (OR and SR, respectively) [11]. In OR, listeners transcribe children’s utterances orthographically or phonetically, and use such information to construct an index of SI. In contrast, under SR, listeners directly produce the SI index using one or a combination of the following procedures: absolute holistic (HJ), analytic (AJ), or comparative (CJ) judgments, the last, also known as the relative holistic method.

It is easy to infer from the previous description, that under perceptual studies, OR methods are more valid² and reliable³ than SR methods, and therefore as their

¹ordered on increasing level of ecological validity [8, 6]

²defined as the extent to which scores are appropriate for their intended interpretation and use [14, 21].

³the extend to which a measure would give us the same result over and over again [21], i.e. measure something, free from error, in a consistent way.

name imply, are usually used as an objective measure of SI [2, 7]. However, given the demanding process in terms of number of listeners required, time, and ultimately cost entailed by OR methods, SR methods can be regarded as an efficient alternative, as long as we can ensure they provide equally valid and reliable SI measures.

Furthermore, within the SR methods, the literature evidence indicate that while HJ procedures are less time consuming than any other alternative [2], they suffer from a lack of intra- and inter-rater reliability⁴ [16, 12, 11, 2]. Additionally, the literature inform us the procedure does not allow to assess subtle differences in the representations [2], while the scales derived from them are usually coarse, where children reach the higher levels fairly quickly [19], e.g. the Speech Intelligibility Rating (SIR) [5, 15].

In this context, CJ has received a growing attention, because it directly tackles the issues with the HJ procedures: it fosters reliable [24] and valid scores [3, 14], while the judges can focus only on the relevant aspects of the compared representations, i.e. the "just noticeable difference" [14]. Moreover, depending on the task, it provides a set of additional benefits, e.g. judges feel more comfortable using comparisons, which foster more accurate judgments [10], it does not require a high level of expertise [14, 1], it encourages to tackle hard to operationlize or open-ended tasks [17, 18, 14], and the measurement of competencies [23], among others.

In the current study, the speech intelligibility of children with a cochlear implant (CI) is investigated in comparison with that of peers with normal hearing. A CI partially restores a severe-to-profound sensorineural hearing loss. Even though the signal provided by a CI is still degraded compared to the signal in normal hearing (Drennan, Rubinstein, 2008), the device enables children with severe-to-profound hearing impair- ment to perceive speech and other environmental sounds. After cochlear implantation,

As expected, model (1) allow us to construct the *most objective* measure of a child's SI ranking. Moreover, using such model we will be able to test some research hypothesis of our interest.

⁴the lack of *intra-rater reliability* happens when the listener rates the same speech recording (representation) after a time lapse, and does not arrive at exactly the same score. On the other hand, the lack of *inter-rater reliability* happens if two listeners, who independently rate the same representation, does not arrive to the same score [21].

2. Materials and Methods

2.1. Children

Thirty two (32) children were selected using a large corpus of *spontaneously spoken speech*, collected by the Computational Linguistics, Psycholinguistics and Sociolinguistics research center (CLiPS). The selection followed a two step procedure, similar to one outlined in Faes et al. [7]. First, a [convenient sample](#) of hearing impaired children was selected based on the quality of their registered stimuli (utterances). And second, a [matched sample](#) of normal hearing children was also selected.

For the first step, the [convenient sample](#) of 16 hearing impaired children with cochlear implant (HI/CI) were all native speakers of Belgian Dutch, living in Flanders, the Dutch speaking area of Belgium. They were all raised in monolingual Dutch with a limited support of signs, and all were screened as hearing impaired by the Universal Neonatal Hearing Screening (UNHS) using automated Auditory Brainstem Response hearing tests for newborns. After the identification of their hearing loss, the children were referred to a specialized audiological center for further audiological workup.

For the second step, 12 normal hearing children (NH) are matched on gender, age, and regional background, to the groups selected in the previous step. The matching is done through a [Manual or Propensity Score Matching \(PSM\)](#) procedure, [explain the appropriate procedure](#).

Finally, the researcher considers important to highlight two relevant points from the children's selection process. First, while the matching procedure for the NH group uses the child's *age* (at recording), the method cannot use the same variable for the other two groups. This is due to the fact that *age* is merely used as a proxy, for the amount of time a child has been developing his(her) language. In that sense, more appropriate variables to use under the HI/CI and HI/HA groups would be e.g. the *device length of use*, which approximates the "hearing age" of such children, or their *vocabulary size*, which resembles their "lexical age" [7]. For this research, we consider the *device length of use* as the simplest one to implement. Second, due to the nature of the sample selection procedure, we cannot ensure the HI/CI and HI/HA, nor the NH group, are representative of their respective populations. Therefore, inferences beyond this particular set of children must be taken with care.

2.2. Stimuli

The stimuli consisted of the children’s utterances (sentences of similar length) recovered from previously mentioned CLiPS corpus. More specifically, we use a portion of the corpus that consisted of 10 utterances recordings, for each of the 32 selected children. The stimuli were documented when the child was telling a story cued by the picture book "Frog, where are you" [?] to a caregiver "who does not know the story". The quality of the stimuli was ensured by selecting utterances with no syntactically ill-formed or incomplete sentences, any background noise, cross-talk, long hesitations, revisions or non-words [2].

As a result, the data set consisted in a total of 320 utterances⁵ presented to the listeners in a random order, based on the adaptive pairing algorithm [18] implemented in Comproved⁶.

Similar designs were used by Boonen et al. [1] and Faes et al. [7]. However, in the former case the number of samples were low, while in the latter, the design was unbalanced and not conducive to appropriate inferences from the pairwise comparisons.

2.3. Experimental setup

On both procedures, the experimental settings for the **judgment task** followed the next steps [1, 2]:

1. the listener take a seat in front of a computer screen, located at his(her) home, workplace, or the experimental laboratory.
2. the listener open Comproved⁷ and select the rating task.
3. the listener read two set of instructions presented on the computer screen about:
 - a) how to perform the task, and

⁵under the Design of Experiments (DoE) literature, we would say we have 32 experimental units with 10 replicate runs each, making a total of 320 experimental runs. As it is defined in [?], an experimental unit is the item under study upon which something is changed, while a replicate run is the experiment conducted with the same factor settings, but using different experimental units.

⁶evidence suggest that the number of comparisons and pairing algorithm impacts the reliability, validity and efficiency of the procedure [? ? 14, 24]. However, this is not investigated in the current research.

⁷software developed by the University of Antwerp designed to perform comparative judgments: <https://comproved.com/en/a>.

- b) the aspects not to consider for the task.
- 4. the listener hear the stimuli through high quality headphones, set at a comfortable volume.
- 5. the listener rate which stimulus sounded more intelligible by selecting the appropriate button, for CJ-D and CJ-O tasks, or a score from a slider on the computer screen, for the HJ task.
- 6. the listener provide a decision statement on why the selected stimulus sounded more intelligible.

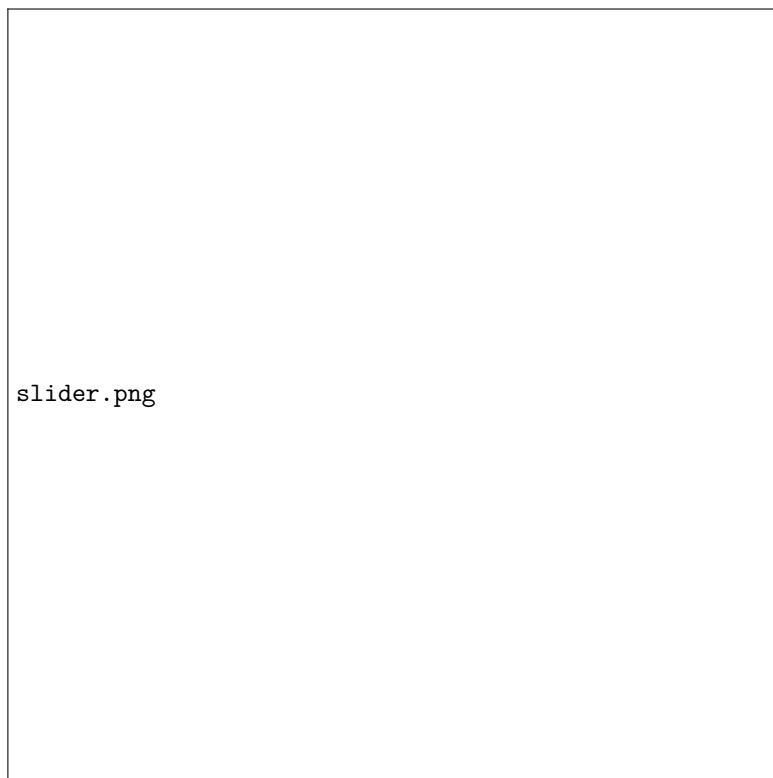


Figure 1: Slider for the HJ task. Extracted from Boonen et al. [2].

Observational evidence indicate the HJ procedure might suffer from anchoring effects⁸ or issues with the default option of the slider. About the former, the anchoring seem to happen when listeners consider the previous assessment as a reference point for the next, effectively turning the task into a comparative rating, similar to CJ. About the latter, as the default setting for the slider is located on the far left for each new assessment (as seen in Figure 1), it is likely that such setting might impact the rating procedure⁹. Considering the previous, in order to minimize both issues, care is taken to randomize the display of stimuli within each listener. However, the researcher recognizes that a better approach to face the second problem would be to randomize the default setting of the slider, but this will not be applied nor investigated on the current research.

talk about decision statements or thinking-at-loud tasks.

Finally, for the **transcription task**, the followed steps were similar to the previous task. Although, at the fourth step, the listeners did not rated the stimulus but rather wrote their orthographic transcriptions, in a free text field in the Comproved environment.

On the other hand, for the transcription task, 100 transcribers participated in the experiment. The participants and stimuli were divided into five groups, where each group of 20 students (100/5) transcribed 64 stimuli on their series (320/5), resulting in 20 transcriptions per utterance ($64 \times 100/320$). In total we registered 6400 transcriptionsfoot:doe.

3. Results

3.1. Data collection

Research has shown that children with CI can attain spoken language skills similar to those of their normal hearing peers after three to four years of device use (i.a., Bruijnzeel, Ziylan, Stegeman, Topsakal, & Grolman, 2016; Dettman, Dowell, Choo, Arnott, Abrahams, Davis, Dornan, Leigh, Constantinescu, Cowan, & Briggs, 2016; Geers, & Nicholas, 2013; Wie et al., 2020). However, the population of children with CI is characterized by remarkable variation. On the one hand, variation relates to differences between individual children: while a considerable number of children with

⁸a bias in decision that occurs when people anchor their decisions around a reference point, and adjust their choices relative to it [? ?].

⁹compelling evidence on how default settings impact several decision process can be found in ?] and ?].

CI appear to catch up with their NH peers, some do not catch up at all (Duchesne, & Marschark, 2019; Geers, Nicholas, Tobey, & Davidson, 2016; Nicholas, & Geers, 2007). On the other hand, variation also relates to differences between domains: some areas of speech and language appear to be more difficult to master than others (Duchesne, & Marschark, 2019). For instance, Faes, Gillis, & Gillis (2015) showed that in a group of children with CI acquiring Dutch, inflectional morphology and sentence length (as a proxy of syntagmatic development) were age-appropriate when the children were 7;0, but the former (and not the latter) was already age-appropriate at age 5;0. Moreover, the phonetics of the same childrens production of vowels was still significantly different from the vowels of their NH peers at the age of 7;0 (Verhoeven, Hide, De Maeyer, Gillis, & Gillis, 2016). Thus, although children with CI start with an initial delay in spoken language, a quite significant group eventually reaches age appropriate levels of linguistic functioning. But the individual variation is also quite large: while some children do catch up with their normally hearing peers, others do not achieve much language comprehension and production even after five years of device use (Barnard, Fisher, Johnson, Eisenberg, Wang, Quittner, Carson, & Niparko, 2015).

As to intelligibility, most studies found that CI childrens speech intelligibility is less well developed than that of their NH peers (i.a., Castellanos, Kronenberger, Beer Henning, Colson, & Pisoni, 2014; Chin, & Kuhns, 2014; Freeman et al., 2017; Grandon et al., 2020).

4. Discussion

5. Author contributions

Jose Rivera performed the statistical analysis, Sven de Maeyer supervised the production of the documents and statistical results, and Steven Gillis collected the data.

6. Financial support

What is the financial support of the project

7. Conflicts of interest

The authors declare they have no conflict of interest.

8. Research transparency and reproducibility

The model's simulation procedures and testing that support the findings of this study are openly available at https://github.com/jriveraespejo/PhD_UA_paper1.

Due to the privacy and confidentiality of subjects, the data set in which the model was implemented cannot be put online.

A. Supplementary

A.1. About Speech Intelligibility

As it was specified in the document Speech Intelligibility is defined as the extent to which the elements in an speaker’s acoustic signal can be correctly recovered by a listener. In the context of the ,

is interpreted here as a latent trait of individuals which underlies the probability of answering the questions in the sample correctly. Henceforth, statements such knowledge is influenced by can be read as the probability of correctly answering the questions in the sample is influenced by. Despite this practical approach, we did our best to improve the construct validity of our study, by both developing the questionnaire and managing the resulting data in collaboration with informants sharing language and culture with the interviewees. We then expect knowledge, as measured by our model, to reflect the general ecological knowledge possessed by individuals, but do not deal with general epistemological considerations on the connection between the two.

A.1.1. Definition

A.2. Sampling bias

A.3. Children characteristics

A.4. DAG: factors influencing Intelligibility

The characteristics of the selected children is detailed in Table 1 from Appendix A.3. The table includes all the variables used for the matching procedure in Section ??, and additionally, shows the child’s etiology, i.e. the cause of their hearing impairment, and their post-implant pure tone average (PTA), i.e. the child’s subjective hearing sensitivity, aided and unaided by their hearing apparatus. No other variables are included, as no known additional comorbidities, beside their hearing impairment, is suspected.

From the table, [describe summaries from the table.](#)

Many factors have been shown to contribute to the success of spoken language development of children with CI, including: (1) audiology related factors, such as the age at implantation, the duration of device use, bilateral (or contralateral) cochlear implantation and the childrens preoperative and postoperative hearing levels; (2)

child related factors, such as the cause of the hearing impairment (genetic, infections), gender, additional disabilities (mental retardation, speech motor problems); and (3) environmental factors, such as communication modality. An overview is provided in Boons, Brokx, Dhooge, Frijns, Peeraer, Vermeulen, Wouters, and van Wieringen, 2012, Fagan, Eisenberg, and Johnson, 2020, Gillis, 2018 and Niparko, Tobey, Thal, Eisenberg, Wang, Quittner, and Fink, 2010. A factor of particular importance here is age. Studies have shown that chronological age is an important factor for intelligibility: as they grow older, childrens intelligibility increases irrespective of their hearing status (Grandon et al., 2020). But in the case of children with CI, age is a complicated factor, since it can not only refer to childrens chronological age (as is the case for children with NH), but also to the childrens so-called hearing age, which is the amount of time between the activation of their device and their chronological age. For instance, a child implanted at the age of 1;0 has a hearing age of two years at the age of 3;0. In addition, the age at implantation has been shown to play a critical role in childrens spoken language achievements. In general, earlier implantation appears to lead to better results than later implantation in several domains (Boons et al., 2012; Niparko et al., 2010). But the research findings with respect to the effect of the variable age on children with CIs intelligibility are not unequivocal. In some studies, a significant effect of chronological age on childrens intelligibility was found (i.a., Flipsen, & Colvard, 2006; Grandon et al., 2020; Habib, Waltzman, Tajudeen, & Svirsky, 2010) but not in others (e.g., Khwaileh, & Flipsen, 2010). Hearing age was found to be a significant predictor of intelligibility by i.a., Flipsen and Colvard (2006), but hearing age was not always considered as a predictor. Age at implantation predicted childrens intelligibility in a considerable number of studies (i.a., Grandon et al., 2020; Habib et al., 2010; Montag, AuBuchon, Pisoni, & Kronenberger, 2014; Svirsky, Chin, & Jester, 2007) but this was not the case in other studies (i.a., Flipsen, & Colvard, 2006; Khwaileh, & Flipsen, 2010). Nevertheless, a general finding appears to be that earlier implantation leads to better results in speech and language development and in intelligibility. At present there is consistent evidence that implantation in the first two years of life leads to consistently better results in spoken language development in comparison to later implantation, and even (inconclusive) evidence for even better outcomes of implantation in the first year of life (Bruijnzeel et al., 2016; Dettman et al., 2016).

A.5. Model details

A.5.1. Definition

Previous research already used hierarchical models with the replicated entropy measures as outcomes [2, 7]. Hierarchical models are powerful to control for heterogeneity in the data, and also to avoid pre-aggregating procedures that could be pernicious for a proper statistical inference [?].

These claims are easier to understand using a thought experiment within our research. Consider we have two children with the same mean entropy, but the second child shows more variability across the 10 utterances than the first. It is clear that the average entropy measure informs about the child’s average SI, indicating that both children have similar level. However, the entropy’s heterogeneity across the 10 utterances also informs about the child’s SI, as a higher variability imply transcribers agreed less about the second child’s intelligibility.

The intuition derived from the previous thought experiment is similar to the one presented in Boonen et al. [2], and it is what justify our use of a hierarchical model. More specifically, we will use a Hierarchical (Mixed) Beta Regression model [?], for which we argue, its implementation is rather trivial under the bayesian framework, and we present it in the following lines.

First, figure 4 depicts the DAG representation of the model. For the measurement error part, section ?? reveals the (observed) entropy replicates H_{ik}^O can represent multiple realizations of a child’s *true* entropy H_i^T , measured with error e_i . As a result, we can say the k ’th entropy measure is nested within the i ’th child, where $k = 1, \dots, K$, $i = 1, \dots, I$, $K = 10$ utterances, and $I = 32$ children.

Second, for the hypothesis part, we can say the child’s *true* entropy H_i^T is inversely explained by the child’s speech intelligibility index SI_i , and in turn, the latter by a set of covariates. Notice from Figure 4, we propose two sets of models. The model in panel (a) use hearing status (HS_i) and hearing age (A_i) as covariates. The use of hearing status is justified as we are interested in comparing SI among groups, defined by the children’s hearing characteristics (NH, HI/CI, and HI/HA). On the other hand, we expect hearing age¹⁰ and its interaction with hearing status, to also have an effect on the SI index, as previous evidence have shown the speech of HI children gradually approximate that of NH children [?].

Notice the model depicted in panel (a) is interested on (what we can call) *total effects*, i.e. the effects of the hearing characteristics, not independent from the effects of the hearing apparatus (cochlear implant or hearing aid). This is important

¹⁰see section ?? to know how the variable is defined.

to understand for two reasons. Since a hearing apparatus is fitted onto a child depending on aspects such as the locus and severity of his(her) hearing impairment [?]: (1) such specific children’s characteristics could confound the (beneficial) effects of using specific hearing apparatuses, while (2) because children are selected from a convenient sample, not representative of their respective populations (see section ??), the need to control for such characteristics is paramount, if we seek to obtain effects that can generalize better and beyond our sample¹¹.

Considering the previous, we propose the model depicted in panel (b), where we control for the possible confounding variables etiology (E_i), [as a proxy of locus](#), and unaided PTA (PTA_i), as a proxy for hearing impairment severity. In that sense, the model would estimate (what we can call) the *direct effects* of the hearing apparatus, independent of the children’s characteristics.

Lastly, we proceed to use probabilistic programming to declare the algebraic structure of our models. Given the panel (a) model is nested within the panel (b) model,

¹¹follow the *notes* folder, to see a graphical though experiment.

we declare only the model structure for the latter:

Likelihood:

$$H_{ik}^O \sim \text{BetaProp}(H_i^T, M_i) \quad (1)$$

Transformed parameters:

$$H_i^T = \text{logit}^{-1}(-SI_i) \quad (2)$$

Linear predictor:

$$SI_i = a_i + \alpha + \alpha_{HS[i]} + \beta_{A,HS[i]}(A_i - \bar{A}) + \alpha_{E[i]} + \beta_P PTA_i \quad (3)$$

Priors:

$$M_i \sim \text{LN}(\mu_M, \sigma_M) \quad (4)$$

$$a_i \sim \text{N}(\mu_a, \sigma_a) \quad (5)$$

$$\alpha \sim \text{N}(0, 0.5) \quad (6)$$

$$\alpha_{HS[i]} \sim \text{N}(0, 0.5) \quad (7)$$

$$\beta_{A,HS[i]} \sim \text{N}(0, 0.3) \quad (8)$$

$$\alpha_{E[i]} \sim \text{N}(0, 0.5) \quad (9)$$

$$\beta_P \sim \text{N}(0, 0.3) \quad (10)$$

Hyper-priors:

$$\mu_M \sim \text{N}(0, 5) \quad (11)$$

$$\sigma_M \sim \text{Exp}(1) \quad (12)$$

$$\mu_a \sim \text{N}(0, 0.5) \quad (13)$$

$$\sigma_a \sim \text{Exp}(1) \quad (14)$$

$$(15)$$

where $\text{logit}(x) = \log[x/(1-x)]$, and $\text{logit}^{-1}(x) = \exp(x)/(1+\exp(x))$. Additionally, a $\text{BetaProp}(\mu, \theta)$ distribution is equal to a $\text{Beta}(\alpha, \beta)$ distribution, with $\alpha = \mu\theta$, $\beta = (1-\mu)\theta$. For our purposes, $\mu = H_i^T$ and $\theta = M_i$, the latter denoting the "sample size" of the distribution. Moreover, a_i denote the children's random effects,

α the fixed effects' intercept, $\alpha_{HS[i]}$ and $\beta_{A,HS[i]}$ the intercept and slope of "hearing age" per hearing status group, $\alpha_{E[i]}$ the intercept per etiology group, and β_P the slope for the standardized PTA levels.

Four important things need to be noticed from the previous algebraic structure. First, all the parameters are estimated in the logit scale and centered at $PTA_i = 0$ and \bar{A} , which denotes the minimum hearing age in the sample. Second, instead of a latent measurement error e_i , we use the latent "sample size" parameter M_i to model the heterogeneity/variability of the duplicate entropies. This effectively works as a measurement error model for the duplicates, as the parameter defines the shape of the distribution. Third, we use mildly informative priors to state our uncertainty regarding the direction and magnitude of the effects¹². Fourth, if we do not consider etiology and PTA values in equation (4), we obtain the panel (a) model.

A.5.2. Priors

A.5.3. Estimation

The models proposed in sections ?? and ?? will be estimated under the Bayesian framework¹³. More specifically, we will use the No-U-Turn Hamiltonian Monte Carlo algorithm (No-U-Turn HMC) [? ? ? ?]. `Stan` [?] will be the software package that will provide us with the No-U-Turn HMC machinery, while `R` [?] and its integration packages [?], the software that will allow us to analyze its outputs.

A.5.4. Pre-processing

Besides the exclusion of corrupted observations, e.g. no available rating, no other experimental run nor duplicate was eliminated before the modeling process. This decision departs from what it is observed in previous research, e.g. Boonen et al. [1] decided to eliminate "outlying" observations based on misfit analysis [14], while [?] and Boonen et al. [2] did the same based on univariate outlier analysis.

For the case of misfit analysis, we argue that such procedures cannot be used without caution. The literature points out that in the context of CJ, these statistics are always relative, i.e. they depend on other stimulus and judges included in the assessment [17, 18]. Moreover, they have been proven to be less sensitive, as they are calculated with a low number of judgments per representation [17].

On the other hand, for the case of univariate outlier analysis, we argue that outlying observations are interesting cases to analyze [?], and usually they cannot

¹²see [?] (p. 18-19) for an intuition on prior elicitation.

¹³see [?] (p. 11-13, 15-27) for a detailed description of its benefits and shortcomings.

be identified properly outside the context of a full model [?], i.e. what can behave as an outlier based on a univariate analysis, can behave as expected under the appropriate model.

Considering the previous, if we still manage to identify outlying observations within the context of the proposed models (see Section ??), the researcher would rather make the model robust against their influence, playing on the strengths of the bayesian framework, than to eliminate the observations.

A.6. Simulation

Preliminary to the data collection, we simulated data in silico to test the models and inform data collection procedure. The simulation code is available in the GitHub repository. Several functional correlation between age and knowledge have been simulated, and the model used in the analysis - which includes age as a ordinal categorical predictor of knowledge with monotonically increasing effect - has been able to recover the different shapes. Causal effect of activities, family composition and schooling have been simulated and tested.

The simulated data have been used -albeit in a previous version- to estimate the minimum number of interviewees necessary to recover the parameter values. If individuals were to name a maximum of 300 items in the freelist, 50 interviewees would have been sufficient to obtain reliable estimates of the parameters. Given that data collection in vivo is much less regular and less controllable than in silico, we roughly doubled the number of interviewees and that of questions.

A.7. Model selection

Following the successful and comprehensive analysis in [?] and Lesterhuis [14], the current research will also use the Information-Theoretic Approach (ITA) [?] for the selection of competing models. The approach considers three steps: (1) state our hypothesis into statistical models, (2) select among competing models, and (3) make inferences based on one or multiple models.

First, for the translation of our working hypotheses into statistical models, we will use Directed Acyclic Graphs (DAG) and probabilistic programming [?]. A DAG is the simplest representation of a Graphical Causal Model (GCM), a heuristic model that contains information not purely statistical, but unlike a detailed statistical model, it allow us to deduce which variable relationships can provide valid causal inferences [?]. In summary, a DAG is a reasonable way to state our hypothesis, and make our assumption more transparent. However, abide by the *no-free lunch*

rule, the causal inferences produced under the DAG will only be valid if the assumed DAG is correct. In contrast, the probabilistic programming will serve as the algebraic formalist to define our statistical models.

Second, to select among competing models, we will use the Widely Applicable Information Criterion (WAIC) [?], and the Pareto-smoothed importance sampling cross-validation (PSIS) [?]¹⁴. Two reasons justify our decision. First, both criteria allow us to embrace the full flexibility and information of our bayesian implementation (outlined in Section ??). Last, and more important, both criteria provide us with the best approximations for the out-of-sample (cross-validated) deviance [?]. The deviance is the best approximation for the Kullback-Liebler (KL) divergence [?], i.e. a measure of how far a model is from describing the *true* distribution of our data. ?] points out that is a rather benign characteristic of the model’s selection procedure that we do not need the KL divergence’s absolute value, as the *true* distribution of our data is not available (otherwise, we would not need a statistical model). But rather, using the difference in deviance between competing models, we can measure which model is the farthest from *perfect (predictive) accuracy* for our data¹⁵.

Finally, considering the evidence provided by the previous step, we proceed to make inferences based on one or multiple models.

¹⁴?] used the Akaikes Information Criterion (AIC) [?] with similar purposes.

¹⁵see ?] (p. 202-211) for the intuition and detailed derivation of the argument.

Child	Hearing	Gender	Regional background	Age (y;m)	Device use (y;m)	Etiology	PTA (dB.)	
	Status						unaided	aided
1	NH	male				genetic		
2	HI/CI	female				CMV infection		
3	HI/HA					unknown		
4								
5								
6								
7								
8								
9								
10								
11								
12								
13								
14								
15								
16								
17								
18								
19								
20								
21								
22								
23								
24								
25								
26								
27								
28								
29								
30								
31								
32								
33								

(y;m) = (years;months)

NH = normal hearing,

HI/CI = hearing impaired / cochlear implant,

HI/HA = hearing impaired / hearing aid

Table 1: Characteristics of selected children.

0.35

entropy_DAG2.png

Figure 2

0.425

22

entropy_DAG1.png

Bibliography

- [1] Boonen, N., Kloots, H. and Gillis, S. [2020]. Rating the overall speech quality of hearing-impaired children by means of comparative judgements, *Journal of Communication Disorders* **83**: 1675–1687.
doi: <https://doi.org/10.1016/j.jcomdis.2019.105969>.
- [2] Boonen, N., Kloots, H., Nurzia, P. and Gillis, S. [2021]. Spontaneous speech intelligibility: early cochlear implanted children versus their normally hearing peers at seven years of age, *Journal of Child Language* pp. 1–26.
doi: <https://doi.org/10.1017/S0305000921000714>.
- [3] Bramley, T. [2008]. Paired comparison methods, in P. Newton, J. Baird, H. Goldsteing, H. Patrick and P. Tymms (eds), *Techniques for monitoring the comparability of examination standards*, GOV.UK., pp. 246–300.
url: <https://www.gov.uk/government/publications/techniques-for-monitoring-the-comparability-of-examination-standards>.
- [4] Chin, S., Bergeson, T. and Phan, J. [2012]. Speech intelligibility and prosody production in children with cochlear implants, *Journal of Communication Disorders* **45**: 355–366.
doi: <https://doi.org/10.1016/j.jcomdis.2012.05.003>.
- [5] Cox, R., McDaniel, D., Kent, J. and Rosenbek, J. [1989]. Development of the speech intelligibility rating (sir) test for hearing aid comparisons, *Journal of Speech, Language, and Hearing Research* **32**(2): 347–352.
doi: <https://doi.org/10.1044/jshr.3202.347>.
- [6] Ertmer, D. [2011]. Assessing speech intelligibility in children with hearing loss: Toward revitalizing a valuable clinical tool, *Language, Speech, and Hearing Services in Schools* **42**(1): 52–58.
doi: [https://doi.org/10.1044/0161-1461\(2010/09-0081\)](https://doi.org/10.1044/0161-1461(2010/09-0081)).
- [7] Faes, J., De Maeyer, S. and Gillis, S. [2021]. Speech intelligibility of children with an auditory brainstem implant: a triple-case study, pp. 1–50. (submitted).
- [8] Flipsen, P. [2006]. Measuring the intelligibility of conversational speech in children, *Clinical Linguistics and Phonetics* **20**(4): 303–312.
doi: <https://doi.org/10.1080/02699200400024863>.

- [9] Freeman, V., Pisoni, D., Kronenberger, W. and Castellanos, I. [2017]. Speech intelligibility and psychosocial functioning in deaf children and teens with cochlear implants, *Journal of Deaf Studies and Deaf Education* **22**(3): 278–289.
doi: <https://doi.org/10.1093/deafed/enx001>.
- [10] Gill, T. and Bramley, T. [2013]. How accurate are examiners’ holistic judgements of script quality?, *Assessment in Education: Principles, Policy and Practice* **20**: 308–324.
doi: <https://doi.org/10.1080/0969594X.2013.779229>.
- [11] Hustad, K., Mahr, T., Natzke, P. and Rathouz, P. [2020]. Development of speech intelligibility between 30 and 47 months in typically developing children: A cross-sectional study of growth, *Journal of Speech, Language, and Hearing Research* **63**: 1675–1687.
doi: https://doi.org/10.1044/2020_JSLHR-20-00008.
- [12] Johannisson, T., Lohmander, A. and Persson, C. [2014]. Assessing intelligibility by single words, sentences and spontaneous speech: A methodological study of speech production of 10-year-olds, *Logopedics Phoniatrics Vocology* **39**: 159–168.
doi: <https://doi.org/10.3109/14015439.2013.820487>.
- [13] Kent, R., Weismer, G., Kent, J. and Rosenbek, J. [1989]. Toward phonetic intelligibility testing in dysarthria, *Journal of Speech and Hearing Disorders* **54**(4): 482–499.
doi: <https://doi.org/10.1044/jshd.5404.482>.
- [14] Lesterhuis, M. [2018]. *The validity of comparative judgement for assessing text quality: An assessors perspective*, PhD thesis, University of Antwerp.
- [15] McDaniel, D. and Cox, R. [1992]. Evaluation of the speech intelligibility rating (sir) test for hearing aid comparisons, *Journal of Speech and Hearing Research* **35**(3): 686–693.
doi: <https://doi.org/10.1044/jshr.3503.686>.
- [16] McLeod, S., Harrison, L. and McCormack, J. [2012]. The intelligibility in context scale: Validity and reliability of a subjective rating measure, *Journal of Speech, Language and Hearing Research* **55**: 648–656.
doi: [https://doi.org/10.1044/1092-4388\(2011/10-0130\)](https://doi.org/10.1044/1092-4388(2011/10-0130)).
- [17] Pollitt, A. [2012a]. Comparative judgement for assessment, *International Journal of Technology and Design Education* **22**: 157–170.
doi: <https://doi.org/10.1007/s10798-011-9189-x>.

- [18] Pollitt, A. [2012b]. The method of adaptive comparative judgement, *Assessment in Education: Principles, Policy and Practice* **19**: 281–300.
doi: <https://doi.org/10.1080/0969594X.2012.665354>.
- [19] Raeve, L. [2010]. A longitudinal study on auditory perception and speech intelligibility in deaf children implanted younger than 18 months in comparison to those implanted at later ages, *Otology and Neurotology* **31**(8): 1261–1267.
doi: <https://doi.org/10.1097/MAO.0b013e3181f1cde3>.
- [20] Rowe, B. and Levine, D. [2018]. *A Concise Introduction to Linguistics*, Routledge.
- [21] Trochim, W. [2022]. The research methods knowledge base.
url: <https://conjointly.com/kb/>.
- [22] van Heuven, V. [2008]. Making sense of strange sounds: (mutual) intelligibility of related language varieties. a review, *International Journal of Humanities and Arts Computing* **2**(1-2): 39–62.
doi: <https://doi.org/10.3366/E1753854809000305>.
- [23] Verhavert, S. [2018]. *Beyond a Mere Rank Order: The Method, the Reliability and the Efficiency of Comparative Judgment*, PhD thesis, University of Antwerp.
- [24] Verhavert, S., Bouwer, R., Donche, V. and De Maeyer, S. [2019]. A meta-analysis on the reliability of comparative judgement, *Assessment in Education: Principles, Policy and Practice* **26**(5): 541–562.
doi: <https://doi.org/10.1080/0969594X.2019.1602027>.
- [25] Whitehill, T. and Chau, C. [2004]. Single-word intelligibility in speakers with repaired cleft palate, *Clinical Linguistics and Phonetics* **18**: 341–355.
doi: <https://doi.org/10.1080/02699200410001663344>.