# Analysis of Hearing Loss Data using Correlated Data Analysis Techniques

Ruth Penman and Gillian Heller, Department of Statistics, Macquarie University, Sydney, Australia.

Correspondence: Ruth Penman, Department of Statistics, Macquarie University, New South Wales, 2109, Australia. E-mail: rpenman@efs.mq.edu.au

## Abstract

Hearing loss is usually measured at 8 frequencies in both ears, giving 16 responses for each subject. Traditionally, analysis of risk factors for hearing loss has used a single response variable based on an aggregate of only three or four of these responses. In this study, based on data from the Blue Mountains Hearing Study (Sindhusake et al, 2001), correlated data techniques have been used to model the responses at each frequency and in both ears. This method provides considerable additional information regarding the effect of risk factors on hearing at the different frequencies and on both ears, over the use of a traditional summary measure.

## Introduction

Hearing Threshold is defined by Bess and Humes (1995) as 'the lowest (softest) sound level needed for a person to detect the presence of a signal approximately 50% of the time.' Therefore a higher hearing threshold indicates a greater hearing loss. Hearing Threshold is measured in decibels hearing level (dBHL) and is measured relative to the sound pressure required for the average young adult to hear at each frequency. That is, a hearing threshold of 0 dBHL means that the subject has the same hearing level as a healthy young adult.

Although Hearing Threshold is usually measured at eight frequencies in both ears, a single measure of hearing loss has historically been used for statistical analysis of risk factors. The 'pure tone average' (PTA) is based on the hearing threshold at either three or four of the middle frequencies in either the better or worse ear, depending on the study. Hearing loss is then classified based on this measure: for example, the criteria used by Mitchell (2002), is: 'mild (>25 and <=40 dBHL); moderate (>40 and <=60 dBHL); marked (>60 and <=90 dBHL); and profound (>90 dBHL).'

Typically, logistic regression is used for the analysis of risk factors based on a binary outcome, for example, the presence or absence of any hearing loss (PTA > 25 dBHL), or the presence or absence of a moderate hearing loss (PTA > 40 dBHL). Mitchell (2002) and Cruickshanks (1998) used this method to examine risk factors associated with age related hearing loss.

This method, apart from lacking a consistent measure, fails to take into account hearing loss at the higher and lower frequencies and does not address the differing effect of risk factors at the various frequencies.

Longford (1993) provided an alternative method of analysing hearing loss data. He considered the responses as having three levels of clustering - subject, ear and frequency, with hearing loss in an individual being highly correlated for frequencies within an ear and between ears. Sex and age were the only risk factors assessed. A normal regression model for serially correlated observations, with both a raw and log-transformed response, was used.

In this study, the data from the Blue Mountains Hearing Study (Sindhusake et al, 2001) was reassessed using correlated data techniques. The gamma distribution was considered, as well as the normal error distribution for a square root transformed response. In addition, in keeping with the traditional methods, a binary variable was created using the presence or absence of a hearing loss (hearing threshold > 25 dBHL) at each frequency in each ear. Logistic regression for correlated data was then used for the analysis of risk factors. SAS version 8.2 was used for the analysis.

**The Study**

The Blue Mountains Hearing Study (BMHS) (Sindhusake et al, 2001) was a population-based survey conducted between 1997 and 1999 in the Blue Mountains area of Australia. The aim of the BMHS was to examine the prevalence and risk factors of age related hearing loss.

There were 2003 subjects aged 54 years or older. All subjects were required to complete a detailed questionnaire, providing details on many aspects of their lives, including exposure to potential risk factors and their medical histories. A comprehensive hearing test was carried out to measure hearing thresholds. Pure tone audiometry was conducted by an audiologist in sound treated facilities. The frequencies used were 250, 500, 1000, 2000, 4000, 6000 and 8000 Hz. If more than 20dB difference existed between 2000 and 4000 Hz then a measurement was also taken at 3000 Hz; otherwise this value was calculated as the average of the hearing thresholds at 2000 and 4000 Hz.

The hearing threshold is measured at 5dB intervals. The minimum value is 0 dBHL, that is the same hearing level as an average young adult, and the maximum measured value is 120 dBHL. Those who had a hearing threshold greater than 120 dBHL were assigned a value of '888'. These are effectively censored observations and should be considered as such; however this is beyond the scope of this study and therefore, for the purposes of this study, a value of 125 dBHL has been assigned to any threshold above the maximum 120 dBHL.

Risk factors used in the final analysis of the data were:- age, sex, industrial noise, family history of hearing loss, current smoking status, alcohol, stroke, diabetes, childhood ear infections, diptheria, measles, mumps and chicken pox. Age was categorised into decades for this analysis; alcohol was

categorised on quantity consumed per week (none; < 8 drinks; 8-20 drinks; 20-40 drinks; >40 drinks). All other variables were binary, based on the presence or absence of the risk factor.

Figure 1 shows the average Hearing Threshold (HT) for four of these risk factors - age, sex, industrial noise and diabetes - across the eight frequencies, for the better ear. The worse ear shows a similar pattern but with higher average HTs.

*Figure 1 here*

Hearing thresholds are higher at the higher frequencies for all covariates, which is typical of age related hearing loss. Hearing thresholds increase with age, but at a greater rate at the higher frequencies. Males and females have similar hearing thresholds at the lower frequencies but males have greater hearing loss at the higher frequencies. Exposure to industrial noise has more impact on hearing at the middle and higher frequencies while diabetes has a similar impact at all frequencies.

## **Correlated Data Models**

The interest in this study is the relationship between the hearing loss outcome and the covariates; the within-subject correlation is not of direct interest, but must be taken into account when estimating the regression parameters. A marginal model as described in Diggle, Heagerty, Liang and Zeger (2002) was used.

Figure 2 shows histograms of the hearing threshold at each of the 8 frequencies for the better ear. The worse ear has similar distributions. The response is clearly not normally distributed and has a different distribution at each frequency.

*Figure 2 here*

Two approaches were tested and compared. The first used a square root transformation of the response and the assumption of a normal error distribution. The second used Generalized

Estimating Equations (GEEs) based on a Gamma response distribution and log link. The Gamma distribution was selected for this method as the response variable was non-negative, continuous and, for most of the frequencies, right skewed.

In addition to these two methods a binary variable was created, defining a hearing loss, for each frequency and in each ear, if the hearing threshold was greater than 25 dBHL. GEEs were also used for this model, based on the binomial distribution with the logit link.

**Normal Regression**

The marginal model is

$$Y_i = X_i\beta + W_i + z_i$$

where $Y_i$ is the response for the $i^{th}$ subject (a 16 element vector in this case), $X_i\beta$ are the fixed effects, $W_i$ is the error term for the serial correlation and $z_i$ is the random error. The term $W_i + z_i$ replaces $\varepsilon_i$, the error term in the uncorrelated model, to take into account the within-subject correlation.

The distributional assumptions are that $W_i \sim N(0, \Sigma_i)$, $z_i \sim N(0, \tau^2 I_{ni})$ and that $W_i$ and $z_i$ are mutually independent. ($n_i$ is the number of observations for the $i^{th}$ subject.) Therefore

$$Var(Y_i) = \Sigma_i + \tau^2 I_{ni} \qquad \text{(Diggle, Heagerty, Liang and Zeger, 2002)}$$

The variance-covariance matrix, $\Sigma_i$ with the following structure was used:

$$\Sigma_i = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix} \otimes \begin{bmatrix} 1 & \rho & \rho^2 & \rho^3 & \rho^4 & \rho^5 & \rho^6 & \rho^7 \\ \rho & 1 & \rho & \rho^2 & \rho^3 & \rho^4 & \rho^5 & \rho^6 \\ \rho^2 & \rho & 1 & \rho & \rho^2 & \rho^3 & \rho^4 & \rho^5 \\ \rho^3 & \rho^2 & \rho & 1 & \rho & \rho^2 & \rho^3 & \rho^4 \\ \rho^4 & \rho^3 & \rho^2 & \rho & 1 & \rho & \rho^2 & \rho^3 \\ \rho^5 & \rho^4 & \rho^3 & \rho^2 & \rho & 1 & \rho & \rho^2 \\ \rho^6 & \rho^5 & \rho^4 & \rho^3 & \rho^2 & \rho & 1 & \rho \\ \rho^7 & \rho^6 & \rho^5 & \rho^4 & \rho^3 & \rho^2 & \rho & 1 \end{bmatrix}$$

where the first matrix models the covariance between ears, and the second matrix the correlation between frequencies, within ears. Thus the hierarchy of frequencies within ears and ears within subjects is modelled using a single parameter ($\sigma_{12}$) for the covariance between ears, and an autoregressive correlation structure for the frequencies (defined as the option UN@AR(1) in SAS, Proc Mixed repeated statement), implying a decaying correlation between the responses as the distance between frequencies increases. Sandwich, or asymptotically consistent, estimators, were used for the estimation of the variance-covariance matrix of the fixed-effects parameter estimates (EMPIRICAL option).

**Gamma Regression**

The marginal model in this case is based on the following, for subject $i$ and frequencies $j$ and $k$:-

$$E(Y_{ij}) = \mu_{ij}$$

$$\log(\mu_{ij}) = x'_{ij}\beta$$

$$Var(Y_{ij}) = \mu_{ij}^{2}\phi$$

$$Corr(Y_{ij}, Y_{ik}) = \rho \quad \text{(exchangeable structure)}$$

$$= \rho^{|j-k|} \quad \text{(autoregressive structure)}$$

(Diggle, Heagerty, Liang and Zeger, 2002). Proc Genmod was used for the analysis; this does not allow for a hierarchical correlation structure as Proc Mixed does. Therefore all 16 readings for the

response variable are considered to be repeated measures within a single subject, without allowing for the two levels. Using an unstructured correlation matrix resulted in some inconsistency in the inference of the estimates; specifying either an exchangeable or an autoregressive correlation structure overcame this problem. In both the exchangeable and the autoregressive correlation structures only one correlation parameter is estimated: the exchangeable correlation structure assumes that the correlation between the responses remains constant regardless of the distance between them, in this case the difference in ear/frequency. The results reported in this paper use the autoregressive correlation structure, however the results using the compound symmetric correlation structure were very similar. The regression parameters $\boldsymbol{\beta}$ were estimated using Generalized Estimating Equations (GEE).

**Binomial Distribution  (Logistic GEE)**

Proc Genmod was also used for this model: in this case the response was the binary hearing loss variable, as defined previously, the distribution was binomial and the link was logit. The model is

$\log\left(\dfrac{p_{ij}}{1-p_{ij}}\right) = x'_{ij}\beta$   where $p_{ij}$ is the probability of a hearing loss for subject $i$ at frequency $j$, and the

within-subject correlation is defined as for the Gamma model.

## Results

Overall there was reasonable consistency in results between the models. The parameter estimates and standard errors for all models are presented in Table 1.

*Table 1*

Frequency was treated as a categorical variable in all models because of the nonlinearity of the relationship between frequency and hearing threshold. The normal regression of the PTA was also computed, for the purpose of comparison with the correlated-data models. A significance level of

α=0.05 was used throughout, and only significant covariates were retained in the final models. Note that the second lowest frequency, 500Hz, is the referent category for frequency, having the lowest mean hearing threshold of all frequencies.

The results of these models are best depicted graphically. Fitted values for the four models are shown in Figure 3 by frequency and age group, in Figure 4 by frequency and gender, in Figure 5 by frequency and industrial noise and in Figure 6 by frequency and diabetes. In each of these graphs the fitted values are for a subject with all other covariates at their referent values.

*Figures 3, 4, 5, 6*

The PTA model is not able to capture the interaction of frequency with the covariates; the fitted values do not reflect the patterns seen in the data in Figure 1. Both the gamma and the normal models, however, do reflect the observed patterns, for each of these covariates. The logistic model reflects the pattern to some extent, but as this model is fitting the probability of a hearing loss, rather than the hearing threshold, some differences are expected.

The plots of the fitted values for the gamma and normal models show that hearing threshold increases with age and increases at a greater rate at the higher frequencies. They also show that hearing thresholds are similar for males and females at lower frequencies but that males have a much higher hearing threshold at higher frequencies. Industrial noise also has a greater impact on hearing at the higher frequencies, whereas diabetes has a more constant effect over all frequencies. Other factors that were shown to be significant risk factors for hearing loss are having a family history of hearing loss and having had ear infections as a child. In some models having had a stroke or chicken pox, being a current smoker and consumption of alcohol were shown to have some impact on hearing levels.

*Table 2*

The estimates for the parameters for the correlation structures are shown in Table 2. The procedure for the normal model also includes p-values, which clearly show that there is significant correlation between ears ($\sigma_{12}$) and between frequencies ($\rho$). The estimates for the correlation parameter for the gamma and logistic models are quite high at 0.82 and 0.53 respectively, suggesting a significant correlation between the ears and frequencies. Therefore, it is important that this correlation is taken into account in any joint modelling of these responses, as has been done in this study.

**Conclusions**

The aggregation of several similar measures into one global measure is a strategy that has been driven by traditional univariate methods of analysis. A global measure has appeal because of simplicity; however it must be recognised that in aggregating variables, information is of necessity lost. Correlated data models provide a method for the joint modelling of several correlated measures. These models better reflect the observed data as they model each measure separately.

The models for hearing threshold demonstrated here, provide information that was not available using traditional methods; for example, a number of the risk factors had increased effects at higher frequencies or different effects in the better or worse ear. It is also possible that there are risk factors that only affect the higher frequencies that have not been detected using the PTA.

Gender provides a good example of the benefits of using this methodology. It has been noted by Mitchell (2002) that gender does not appear to be a significant risk factor when using the three frequency PTA but is significant when using the four frequency PTA. Using correlated data analysis not only shows gender to be a significant risk factor but also shows at which frequencies there is increased risk.

Correlated data models, well supported by current software, have wide applicability to data with serial measures, both normally and non-normally distributed.
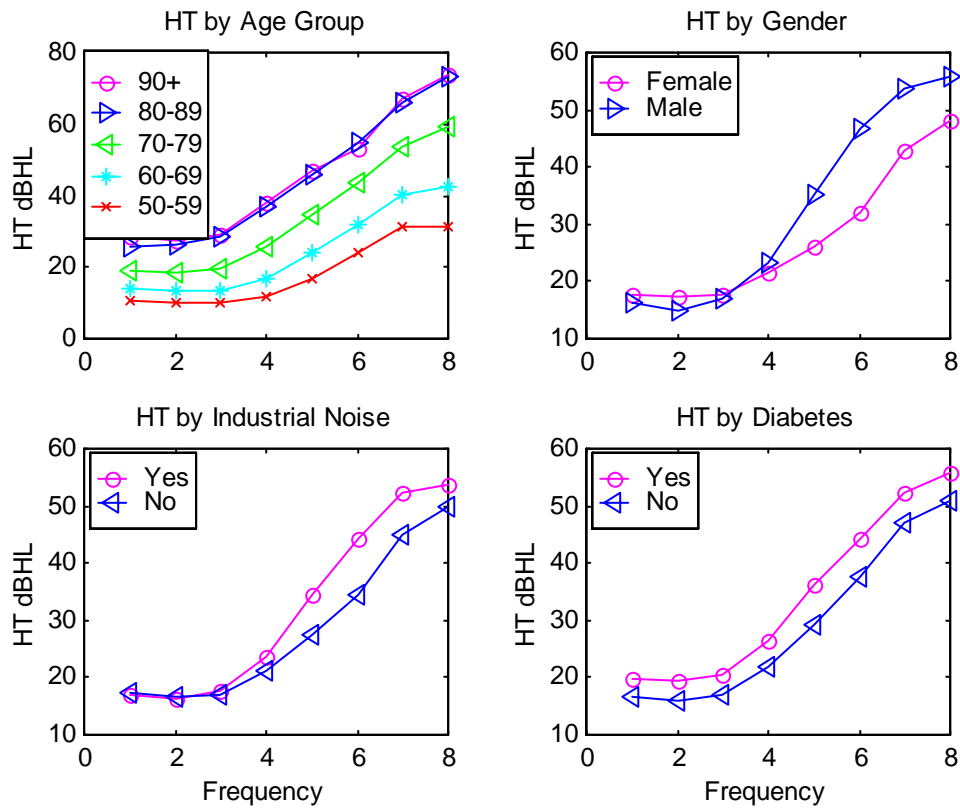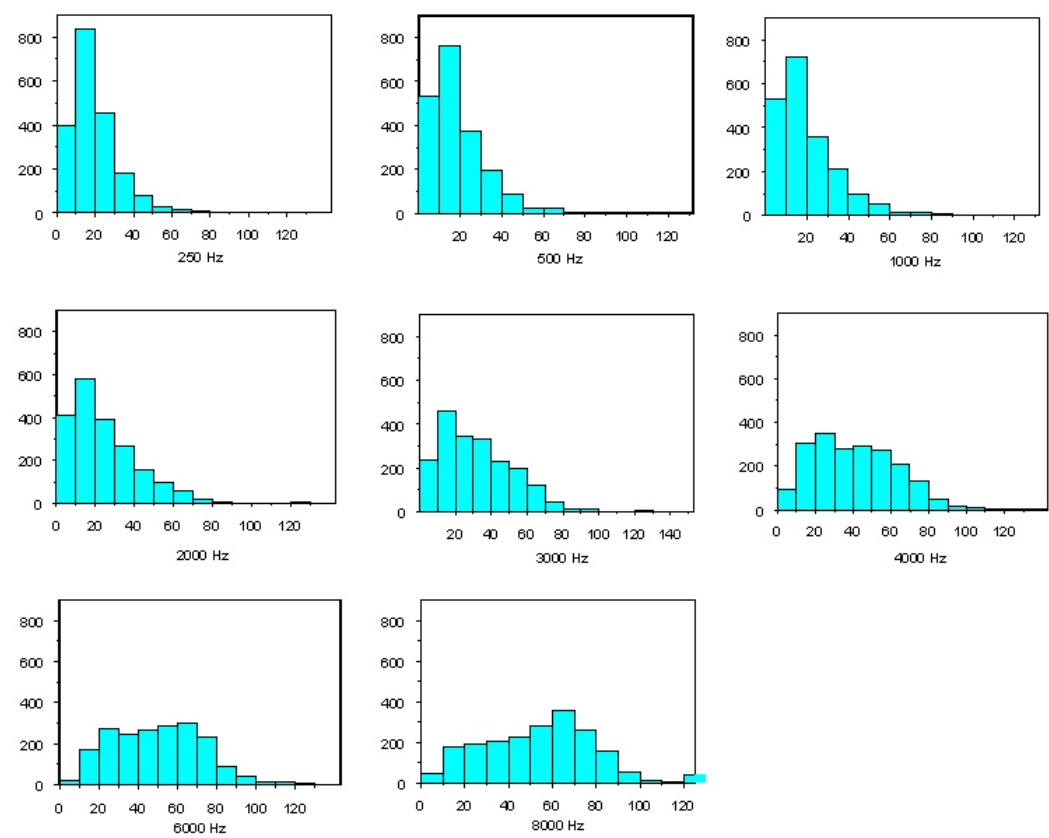
# References

Bess FH, Humes LE (1995) **Audiology: the fundamentals** 2<sup>nd</sup> Edition. Williams and Wilkins, Baltimore.

Cruickshanks KJ, Wiley TL, Tweed TS, Klein BEK, Klein R, Maresperlman JA, Nondahl DM (1998). Prevalence of Hearing Loss in Older Adults in Beaver Dam, Wisconsin: The Epidemiology of Hearing Loss Study. **American Journal of Epidemiology** 148(9): 879-886

Diggle PJ, Heagerty P, Liang KY, Zeger SL, (2002) **Analysis of Longitudinal Data**, Second Edition, Oxford : Clarendon Press ; New York : Oxford University Press

Liang KY & Zeger SL (1986). Longitudinal data analysis using generalized linear models. **Biometrika**. 73, 13-22.

Longford NT (1993). **Random Coefficient Models**. Oxford University Press

Mitchell P (2002) The Prevalence, Risk Factors and Impacts of Hearing Impairment in and Older Australian Community: the Blue Mountains Hearing Study. **The 2002 Libby Harricks Memorial Oration**.

Sindhusake D, Mitchell P, Smith W, Golding M, Newall P, Hartley D & Rubin G (2001). Validation of self-reported hearing loss. The Blue Mountains Hearing Study. **International Journal of Epidemiology** 30, 1371-1378

Zeger SL & Liang KY (1986). Longitudinal Data Analysis for Discrete and Continuous Outcomes. **Biometrics** 42, 121-130.

Zeger SL, Liang KY & Albert PS (1988). Models for Longitudinal Data: A Generalized Estimating Equation Approach. **Biometrics** 44, 1049-1060.

**Figure 1 – Average Hearing Thresholds by Risk Factors for the Better Ear**



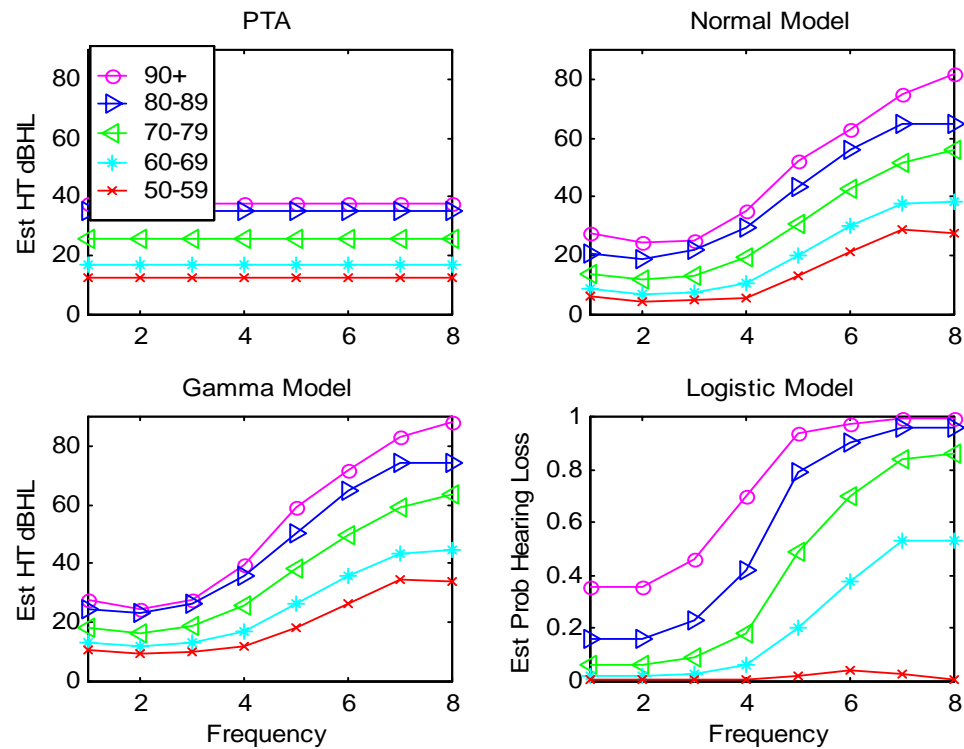The frequencies 250, 500, 1000, 2000, 3000, 4000, 6000 and 8000Hz are labeled as 1 to 8 on the horizontal axes.
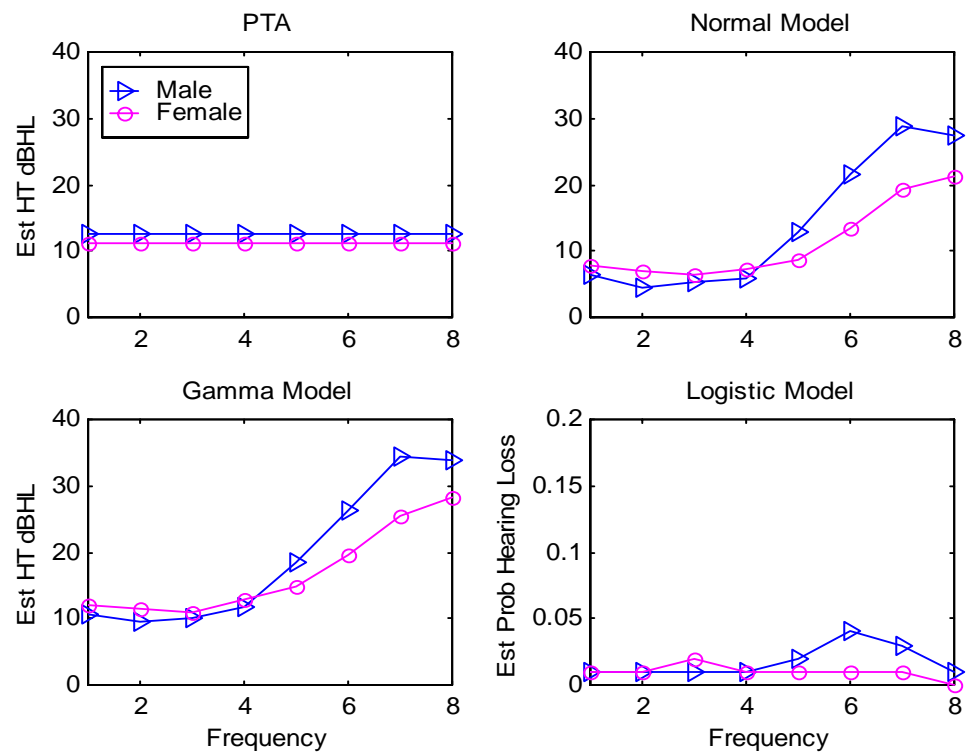
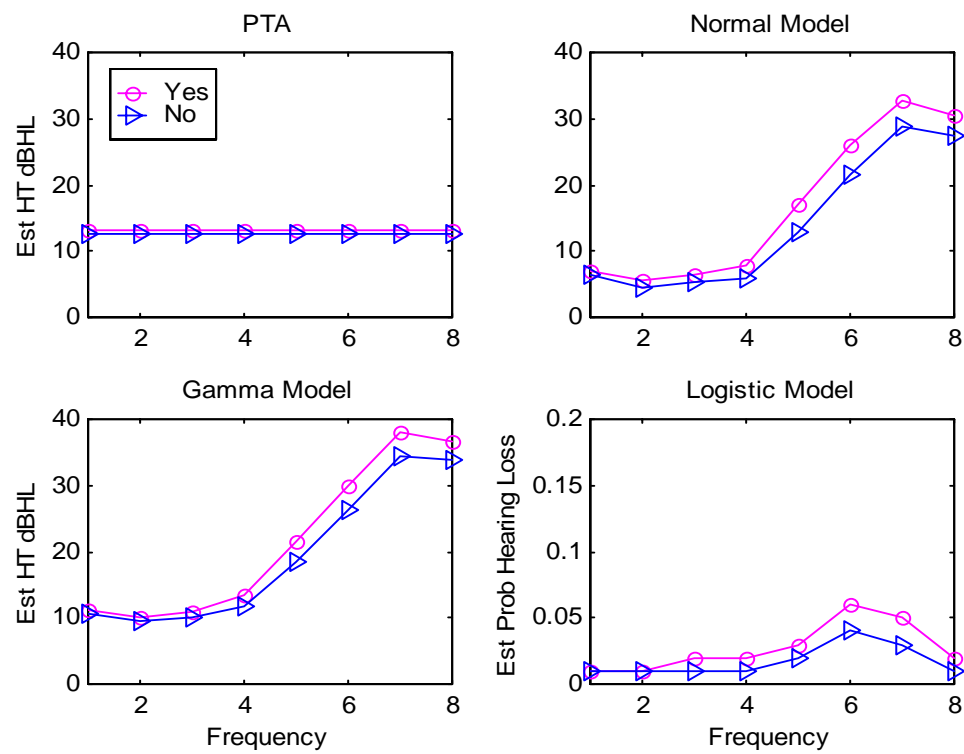**Figure 2 – Histograms of Hearing Thresholds for the Better Ear**

**Figure 3 – Fitted Values by Frequency and Age Group**

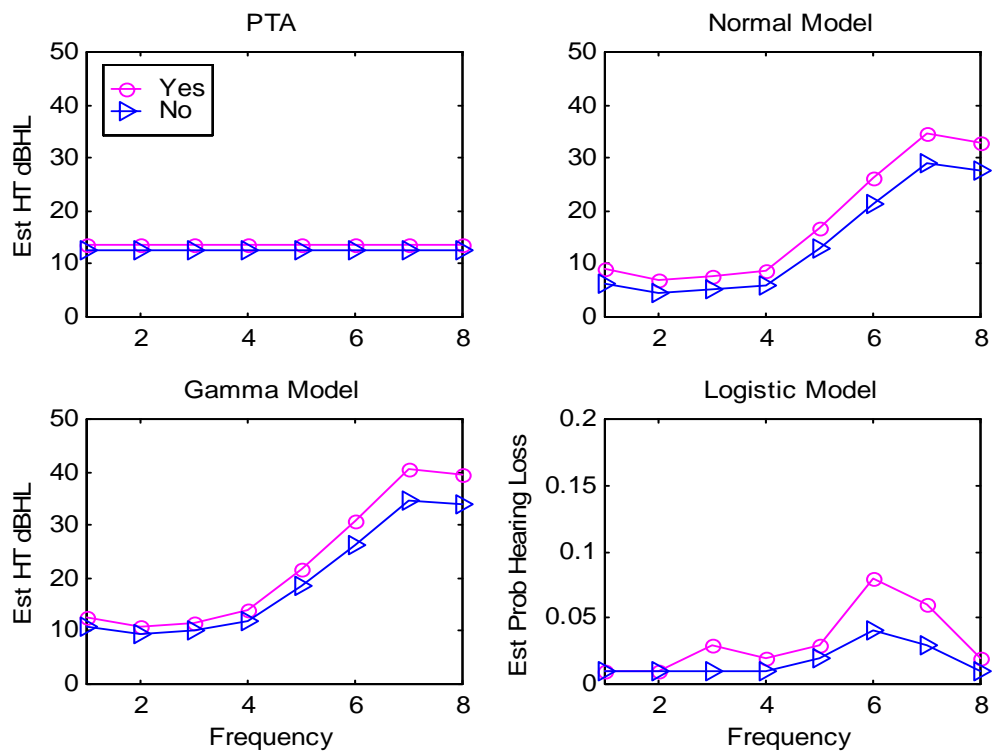**Figure 4 – Fitted Values by Frequency and Gender**

**Figure 5 – Fitted Values by Frequency and Industrial Noise**

**Figure 6 – Fitted Values by Frequency and Diabetes**

# Table 1 – Parameter Estimates

| Variable | PTA uncorrelated normal model $\widehat{\beta}$ | SE($\widehat{\beta}$) | Correlated normal model $\widehat{\beta}$ | SE($\widehat{\beta}$) | Gamma GEE $\widehat{\beta}$ | SE($\widehat{\beta}$) | Logistic GEE $\widehat{\beta}$ | SE($\widehat{\beta}$) |
|---|---|---|---|---|---|---|---|---|
| Intercept | 3.26 | 0.11 | 2.14 | 0.13 | 2.24 | 0.07 | -10.46 | 0.61 |
| **Frequency** | | | | | | | | |
| 250Hz | | | 0.36 | 0.09 | 0.13 | 0.04 | -0.21 | 0.43 |
| 500Hz (referent) | | | - | - | - | - | - | - |
| 1000Hz | | | 0.13 | 0.09 | 0.06 | 0.04 | 0.77 | 0.49 |
| 2000Hz | | | 0.29 | 0.13 | 0.23 | 0.05 | 0.53 | 0.58 |
| 3000Hz | | | 1.46 | 0.13 | 0.68 | 0.06 | 0.89 | 0.66 |
| 4000Hz | | | 2.49 | 0.13 | 1.04 | 0.06 | 1.86 | 0.71 |
| 6000Hz | | | 3.24 | 0.14 | 1.31 | 0.06 | 1.58 | 0.82 |
| 8000Hz | | | 3.10 | 0.15 | 1.28 | 0.06 | 0.49 | 0.86 |
| **Ear** | | | | | | | | |
| Better (referent) | | | - | - | - | - | - | - |
| Worse | | | 1.01 | 0.05 | 0.48 | 0.03 | 0.91 | 0.06 |
| **Age Group** | | | | | | | | |
| 90+ | 2.71 | 0.33 | 2.81 | 0.36 | 0.95 | 0.14 | 0.11* | 0.01 |
| 80-89 | 2.42 | 0.12 | 2.18 | 0.15 | 0.90 | 0.06 | ?? | |
| 70-79 | 1.53 | 0.10 | 1.28 | 0.12 | 0.58 | 0.06 | | |
| 60-69 | 0.63 | 0.10 | 0.45 | 0.12 | 0.24 | 0.06 | | |
| 50-59 (referent) | - | - | - | - | - | - | | |
| **Sex** | | | | | | | | |
| Male (referent) | - | - | - | - | - | - | - | - |
| Female | -0.19 | 0.07 | 0.51 | 0.09 | 0.20 | 0.04 | 0.51 | 0.13 |
| **Ear*Frequency** | | | | | | | | |
| Ear(worse)* 250Hz | | | 0.00 | 0.02 | -0.01 | 0.01 | 0.19 | 0.07 |
| 1000Hz | | | -0.04 | 0.03 | -0.03 | 0.01 | -0.20 | 0.06 |
| 2000Hz | | | 0.09 | 0.03 | -0.04 | 0.01 | -0.10 | 0.07 |
| 3000Hz | | | -0.04 | 0.03 | -0.11 | 0.01 | -0.10 | 0.08 |
| 4000Hz | | | -0.02 | 0.03 | -0.13 | 0.01 | 0.04 | 0.08 |
| 6000Hz | | | -0.05 | 0.03 | -0.17 | 0.01 | 0.24 | 0.10 |
| 8000Hz | | | 0.00 | 0.03 | -0.17 | 0.01 | 0.14 | 0.09 |
| **Freq*Age Group** | | | | | | | | |
| 250Hz* 90+ | | | -0.04 | 0.18 | -0.01 | 0.06 | 0.002* | 0.006 |
| 80-89 | | | -0.14 | 0.09 | -0.06 | 0.03 | | |
| 70-79 | | | -0.03 | 0.08 | -0.03 | 0.03 | | |
| 60-69 | | | -0.02 | 0.08 | -0.01 | 0.03 | | |
| 1000Hz* 90+ | | | -0.08 | 0.22 | 0.07 | 0.07 | 0.004* | 0.007 |
| 80-89 | | | 0.22 | 0.09 | 0.08 | 0.04 | | |
| 70-79 | | | 0.12 | 0.08 | 0.06 | 0.04 | | |
| 60-69 | | | 0.03 | 0.08 | 0.03 | 0.04 | | |
| 2000Hz *90+ | | | 0.72 | 0.25 | 0.26 | 0.09 | 0.01* | 0.01 |
| 80-89 | | | 0.85 | 0.13 | 0.21 | 0.05 | | |
| 70-79 | | | 0.67 | 0.12 | 0.20 | 0.05 | | |
| 60-69 | | | 0.39 | 0.12 | 0.12 | 0.05 | | |
| 3000Hz * 90+ | | | 0.81 | 0.28 | 0.21 | 0.10 | 0.03* | 0.01 |
| 80-89 | | | 0.80 | 0.14 | 0.10 | 0.06 | | |
| 70-79 | | | 0.70 | 0.13 | 0.15 | 0.05 | | |
| 60-69 | | | 0.46 | 0.13 | 0.12 | 0.05 | | |
| 4000Hz * 90+ | | | 0.51 | 0.31 | 0.05 | 0.11 | 0.02* | 0.01 |
| 80-89 | | | 0.66 | 0.15 | 0.00 | 0.06 | | |
| 70-79 | | | 0.62 | 0.13 | 0.07 | 0.05 | | |
| 60-69 | | | 0.41 | 0.13 | 0.07 | 0.06 | | |
| 6000Hz * 90+ | | | 0.48 | 0.30 | -0.07 | 0.10 | 0.04* | 0.01 |
| 80-89 | | | 0.50 | 0.15 | -0.13 | 0.06 | | |
| 70-79 | | | 0.53 | 0.13 | -0.03 | 0.05 | | |
| 60-69 | | | 0.30 | 0.13 | -0.01 | 0.06 | | |
| 8000Hz * 90+ | | | 1.01 | 0.30 | 0.01 | 0.10 | 0.06* | 0.01 |
| 80-89 | | | 0.99 | 0.16 | -0.05 | 0.06 | | |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 70-79 | | | 0.98 | 0.14 | 0.06 | 0.05 | | |
| 60-69 | | | 0.51 | 0.15 | 0.04 | 0.06 | | |
| **Sex*Frequency** | | | | | | | | |
| Sex(female) * 250Hz | | | -0.23 | 0.05 | -0.07 | 0.02 | -0.09 | 0.09 |
| 1000Hz | | | -0.23 | 0.06 | -0.10 | 0.02 | -0.23 | 0.10 |
| 2000Hz | | | -0.49 | 0.08 | -0.21 | 0.03 | -0.56 | 0.12 |
| 3000Hz | | | -1.16 | 0.09 | -0.42 | 0.03 | -1.36 | 0.14 |
| 4000Hz | | | -1.49 | 0.09 | -0.49 | 0.04 | -1.67 | 0.16 |
| 6000Hz | | | -1.15 | 0.09 | -0.37 | 0.04 | -1.36 | 0.17 |
| 8000Hz | | | -1.00 | 0.09 | -0.32 | 0.04 | -1.24 | 0.18 |
| **Ear* Age Group** | | | | | | | | |
| Ear(worse) * 90+ | | | -0.36 | 0.11 | -0.21 | 0.05 | | |
| 80-89 | | | -0.16 | 0.08 | -0.17 | 0.03 | | |
| 70-79 | | | -0.12 | 0.06 | -0.11 | 0.03 | | |
| 60-69 | | | -0.05 | 0.06 | -0.06 | 0.03 | | |
| **Industrial Noise** | 0.36 | 0.07 | 0.23 | 0.09 | 0.08 | 0.04 | 0.41 | 0.09 |
| **Ind. Noise*Freq** | | | | | | | | |
| IndNoise * 250Hz | | | -0.11 | 0.05 | -0.04 | 0.02 | | |
| 1000Hz | | | 0.04 | 0.06 | 0.02 | 0.02 | | |
| 2000Hz | | | 0.15 | 0.08 | 0.05 | 0.03 | | |
| 3000Hz | | | 0.30 | 0.09 | 0.07 | 0.04 | | |
| 4000Hz | | | 0.25 | 0.09 | 0.05 | 0.04 | | |
| 6000Hz | | | 0.12 | 0.09 | 0.02 | 0.04 | | |
| 8000Hz | | | 0.06 | 0.10 | 0.00 | 0.04 | | |
| **Family History** | 0.24 | 0.06 | 0.21 | 0.06 | 0.08 | 0.02 | 0.29 | 0.08 |
| **Stroke** | | | 0.17 | 0.11 | | | 0.26 | 0.16 |
| **Ear(worse) * Stroke** | | | -0.12 | 0.06 | | | -0.27 | 0.13 |
| **Diabetes** | 0.45 | 0.11 | 0.50 | 0.13 | 0.16 | 0.05 | 0.59 | 0.16 |
| **Ear Infection** | 0.18 | 0.08 | 0.19 | 0.08 | 0.08 | 0.03 | 0.34 | 0.10 |
| **Chicken Pox** | | | | | -0.06 | 0.02 | | |
| **Alcohol** | | | | | | | | |
| None (referent) | | | | | - | - | | |
| <8 drinks | | | | | -0.08 | 0.03 | | |
| 8-20 drinks | | | | | -0.04 | 0.03 | | |
| 20-40 drinks | | | | | 0.01 | 0.06 | | |
| >40 drinks | | | | | 0.10 | 0.10 | | |
| **Current Smoker** | | | | | | | 0.29 | 0.14 |

**Table 2 – Correlation Parameter Estimates**

| Model | Parameter | Estimate | p-value |
|---|---|---|---|
| Normal | $\sigma_1^2$ | 2.01 | <0.0001 |
| | $\sigma_{12}$ | 1.44 | <0.0001 |
| | $\sigma_2^2$ | 2.41 | <0.0001 |
| | $\rho$ | 0.73 | <0.0001 |
| Gamma (AR(1)) | $\rho$ | 0.82 | |
| Logistic (AR(1)) | $\rho$ | 0.53 | |

**Captions for illustrations**

**Figure 1 – Average Hearing Thresholds by Risk Factors for the Better Ear**

**Figure 2 – Histograms of Hearing Thresholds for the Better Ear**

**Figure 3 – Fitted Values by Frequency and Age Group**

**Figure 4 – Fitted Values by Frequency and Gender**

**Figure 5 – Fitted Values by Frequency and Industrial Noise**

**Figure 6 – Fitted Values by Frequency and Diabetes**

**Table 1 – Parameter Estimates**

**Table 2 – Correlation Parameter Estimates**