**Lab Report #1**                                                              **Jinxuan Lu**

## 1. Introduction

This report is about the investigation result of Sue's court case of suing her father for not being with Sue and her mother Kate 20 years ago on the 'Unsinkable' ship Titanic, which led to Kate's death. Based on multiple features of passengers included in the data and whether the passengers survived the disaster, I have created a model that can effectively predict the chances of survival of Sue and Kate, with and without Leonardo, respectively.

## 2. Data and Method

Based on initial exploratory analysis and visualization, several variables were excluded from the dataset in this study. The final dataset included 9 items: survival status, ticket class, sex, age, number of siblings/partners, number of parents/children, fare, cabin, and embarked port. The data was composed of 891 valid individuals aged five months to 80 years old. The average age was around 29, and among the samples, 35% were female.

A binomial logistic regression analysis was conducted to analyze the relationship between the survival status (as the outcome variable) and other predictors.

## 3. Results

The final logistic regression model included an outcome variable (survival status) and five predictors: ticket class, sex, age, number of siblings/spouses, and the port of embarkation. Fare and the number of parents/children were discarded from the model because they were both insignificant, which means that we cannot reject the null hypothesis that there was no relation between survival status and them. Cabin was also excluded since there were too many missing values in Cabin from the beginning, which I was not sure the data was simply missing, or it meant "no cabin". Besides, the intercept was insignificant with Cabin in the model.

The result shows that the final model had a significantly better fit than the null model ($Chi^2 = 390.77$, df = 7, $p < 0.001$). The deviance (-2LL) for this model was 795.89, which was much lower than the null model (1186.66), and it means that the amount of error left after accounting for all the variance explained by the predictors in our model was less than the null model. Besides, the AIC of my model was 811.89, which was also less than the null model (1188.66), more than 2 points (See **Table 1**).

As for the effectiveness, the model explained 32.9% of the variance (McFadden R^2 = 0.329). With 38% of passengers surviving in the original sample (342 out of 891 passengers), the final model correctly predicted 74.56% of the cases for those who actually survived the sink and 81.24% of the cases for those who actually died. In total, 78.68% of the survival status was correctly predicted by this model.

**Table 1** shows the coefficients and other statistics descriptions for each predictor in the model.
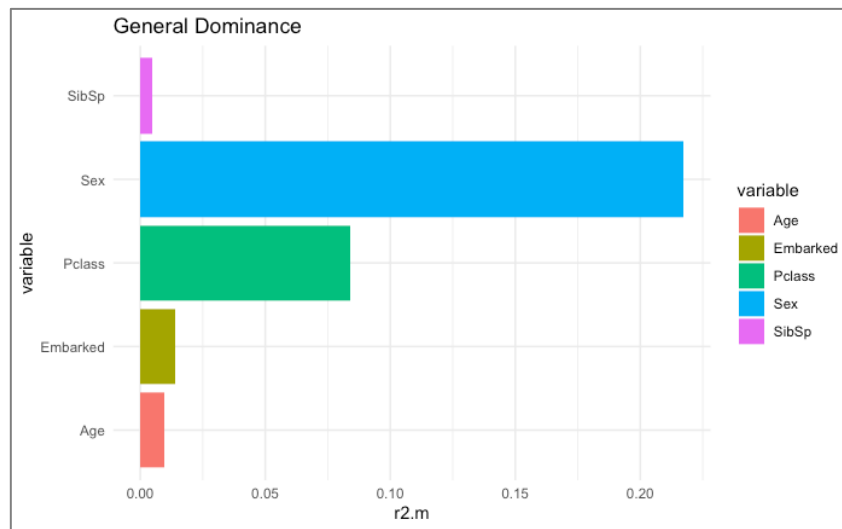
**Table 1.** Logistic regression statistics: significance and odd ratios

| Predictors | Survived (mod1) | | | | | Survived (mod_null) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | *Log(OR)* | *Odds Ratios* | *CI* | *Z value* | *p* | *Log(OR)* | *Odds Ratios* | *CI* | *p* |
| (Intercept) | 1.098 | 3.00 | 1.54 – 5.93 | 3.21 | **0.001** | -0.473 | 0.62 | 0.54 – 0.71 | **<0.001** |
| Pclass [2] | -0.959 | 0.38 | 0.23 – 0.64 | -3.60 | **<0.001** | | | | |
| Pclass [3] | -2.166 | 0.11 | 0.07 – 0.19 | -8.72 | **<0.001** | | | | |
| Sex [female] | 2.681 | 14.59 | 10.07 – 21.51 | 13.85 | **<0.001** | | | | |
| Age | -0.028 | 0.97 | 0.96 – 0.99 | -4.02 | **<0.001** | | | | |
| SibSp | -0.289 | 0.75 | 0.61 – 0.91 | -2.83 | **0.005** | | | | |
| Embarked [Q] | 0.007 | 1.01 | 0.48 – 2.11 | 0.02 | 0.985 | | | | |
| Embarked [S] | -0.463 | 0.63 | 0.40 – 1.00 | -1.98 | **0.047** | | | | |
| Observations | 891 | | | | | 891 | | | |
| Deviance | 795.886 | | | | | 1186.655 | | | |
| AIC | 811.886 | | | | | 1188.655 | | | |
| log-Likelihood | -397.943 | | | | | -593.328 | | | |

Note: Log(OR)=Regression Coefficient; OR = Odds Ratio; CI = 95% Confidence Intervals; AIC=Akaike Information Criterion; Deviance=-2LL.

To determine the relative contribution of predictors in the model, I computed the mean of each predictor's conditional measures and concluded that Sex had the highest value (0.217) and generally dominated all other predictors. In contrast, the number of siblings or partners had the lowest value (0.005) and contributed little to the model (see **Figure 1.**).

**Figure 1.** Average General Dominance of Predictors

The regression equation of the model was written as following: **Survived = 1.10** + 0 \* (**1st Class**) + (**-0.96**) \* (**2nd Class**) + (**-2.17**) \* (**3rd Class**) + 0 \* **Male** + **2.68** \* **Female** + (**-0.03**) \* **Age** + (**-0.29**) \* (**Number of SibSp**) + 0 \* (**C-port Embark**) + **0.01** \* (**Q-port Embark**) + (**-0.46**) \* (**S-port Embark**)

The initial survival rate for Kate was 63.41% without Leonardo on board, whereas it decreased to 56.46% when assuming Leonardo was on board with her. For Sue, whether Leonardo, as a father, was on board with his daughter, the survival rates weren't influenced, and it was always 73.69%.

## 4. Discussion

Based on the results, we can conclude that the presence of Leonardo would have no influence on Sue, but it would negatively affect Kate's survival chance. The best predictors of survival in the model, as I mentioned above, were sex and ticket class, and the presence of a spouse would negatively influence the final survival chance. And since the presence of parents was discarded from the model, it would not influence Sue's final result.

## 5. Reference

Azen, R. & Traxel, N. (2009). Using Dominance Analysis to Determine Predictor Importance in Logistic Regression. *Journal of Education Behavior Statistics, 34,* 319-347. https://doi.org/10.3102/1076998609332754

R Core Team. (2020). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL https://www.R-project.org/.

## 6. Appendices

Exploratory analysis and visualization: (analysis of these plots can be checked in code file)

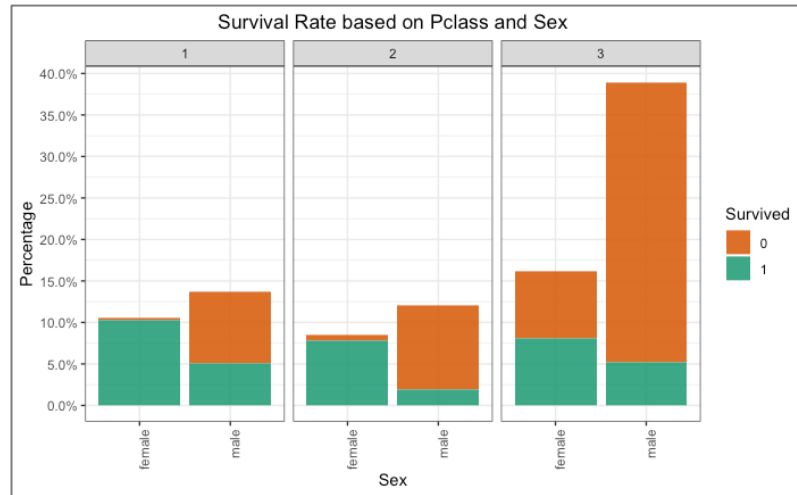**Figure 2.** Survival Rate based on Ticket Class & Sex



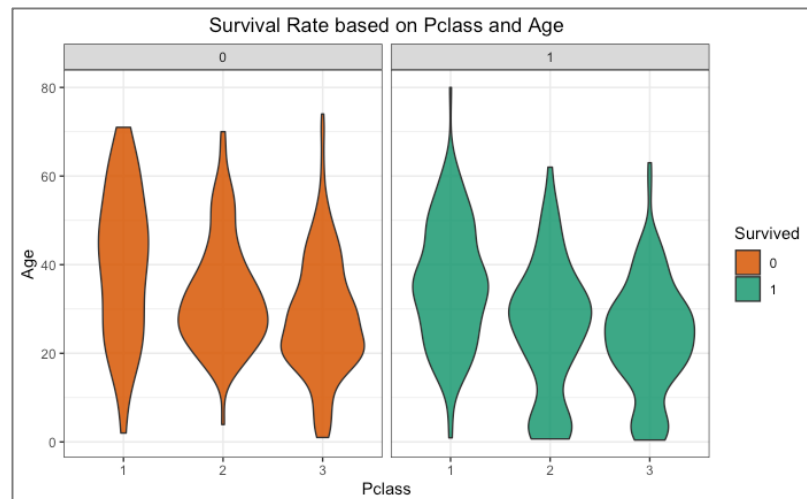**Figure 3.** Survival Rate based on Ticket Class & Age



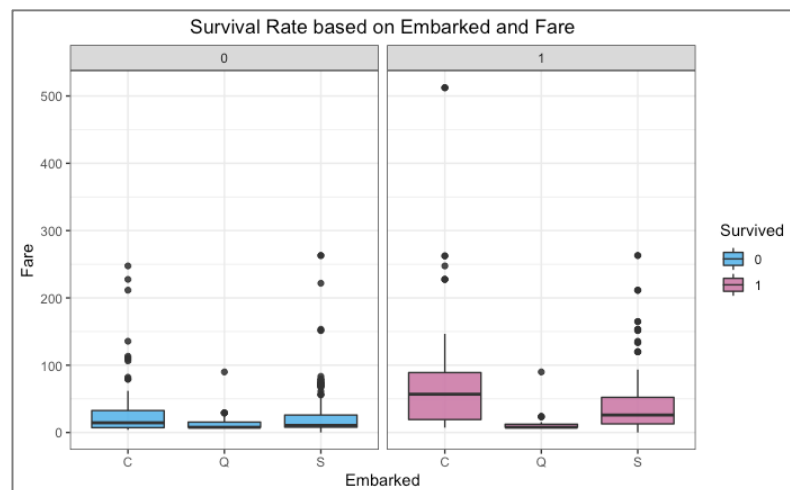**Figure 4.** Survival Rate based on Embarked Port & Fare

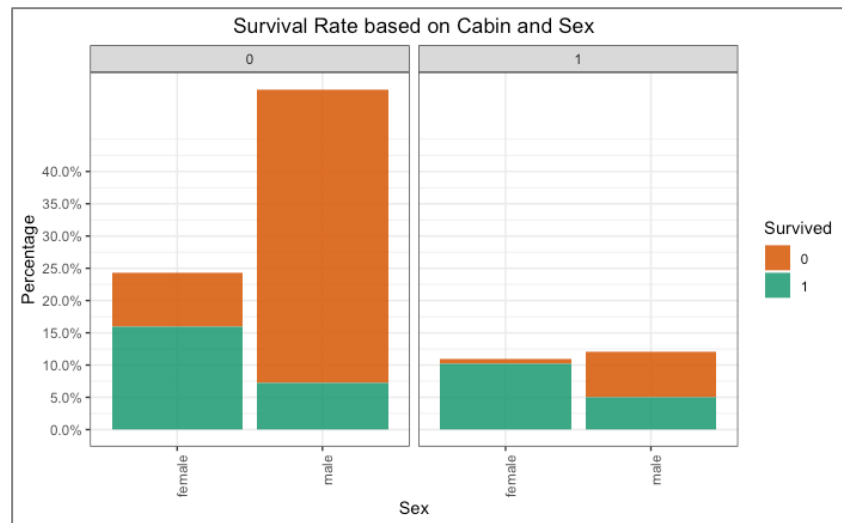**Figure 5.** Survival Rate based on Cabin & Sex



**Link to the codes of report #1:**

https://github.com/Dino-Lu/Lab-Assignments-for-SIMM61 (for all 4 reports)

https://github.com/Dino-Lu/Lab-Assignments-for-SIMM61/blob/main/Lab%20Report%231%20-%20Jinxuan%20Lu.Rmd

**\*\*Note\*\*:** The method I used to tackle **missing values** in the **Age** column was to take a random list of ages that maintains the original statistical summary values. When reproducing the code, the values of the Age column might change. Thus, I will also upload the data file that imputed missing values already for reference.

**Lab Report #2**                                                                    **Jinxuan Lu**

## 1.  Introduction

The aim of the analysis is to estimate patients' postoperative pain after the wisdom tooth surgery on the basis of the data we have collected. Given the clustering structure of the data (patients nested in different hospitals sites), it was considered necessary to build the model by inserting the fixed factors, such as age, sex, pain catastrophizing, serum cortisol, mindfulness, anxiety trait, and a random intercept relating to the different hospitals investigated in the survey.

## 2.  Data and Method

The two datasets used in this study were composed of 200 individuals from 10 different hospitals, respectively. The average age of patients in both datasets was 40, and while there were 48.5% female patients in dataset A, dataset B included a few more females, which accounted for 51% of total patients.

This project mainly included three stages. Firstly, in order to determine the influence of different parameters on the postoperative pain of wisdom tooth surgery, we mainly conducted a linear mixed model (LMM) analysis on dataset A, specifically, random intercept model, since no prior data or theory has shown that different hospital sites would influence the effects of certain parameters on the pain. In addition, we also used the model built on dataset A to predict pain in dataset B. And last, the most influential predictor in the first stage would be used to build both random intercept and slope models on dataset A, and the model fit would be compared based on the graph and other statistics. A description of all variables used in this project is shown in **Table 1**.

**Table 1.** Description and measures of variables in the project.

| Variables | Measuring questions |
|---|---|
| **Fixed effect predictors** | |
| sex | Female/Male |
| Age | Dataset A: 27-53; Dataset B: 26-52 |
| STAI_trait | The State Trait Anxiety Inventory measures trait anxiety on a scale of 20 to 80. The higher the score, the higher the anxiety. |
| pain_cat | The Pain Catastrophizing Scale measures the extent of pain catastrophizing on a scale of 0 to 52. The higher the score, the higher the catastrophizing level. |
| mindfulness | The mindfulness was measured by the Mindful Attention Awareness Scale (MAAS), ranging from 1 to 6, with higher scores representing higher dispositional mindfulness. |

| cortisol_serum | Serum cortisol was measured by collecting blood samples from participants in the waiting room 5 minutes before their operations. |
|---|---|
| **Outcome variable** | |
| pain | The level of pain was recorded using a numerical rating scale of 0 to 10, where 0 means "no pain" and 10 means "worst pain I can imagine". |
| **Random effect predictor** | |
| hospital | The 10 hospitals were coded from hospital_1 to hospital_10 |

Source: Adapted from the file "Lab assignment – Mixed linear models"

## 3. Results

The diagnostics considered and checked for this linear mixed model were the normal distribution of the residues, the verification of the linearity, and independence of all the variables (see graphs in Appendix).

The results showed that the random intercept model was significantly better than the null model, where the fixed effect predictors together explained 38.5% of the variance of post-operative pain of patients (marginal $R^2 = 0.385$ [95% CI = 0.301, 0.488]). And the conditional $R^2$ showed that the entire model explained 46.3% of the variance of post-operative pain (conditional $R^2 = 0.463$). Besides, the differences among different hospitals explained about 13% of the variance that is "leftover" after the variance explained by the fixed effects.

The model coefficients and the confidence intervals of the coefficients for each fixed effect predictor on dataset A are shown in **Table 2**. From the table, we can tell that the effect of **serum cortisol** is the strongest ($\beta = 0.51$, $p < 0.001$, $r^2 = 0.123$) and we can say so with a good degree of certainty (95% CI of $\beta = [0.33\text{-}0.69]$).

In addition, using the model coefficients obtained on dataset A to predict pain in dataset B showed that 36.6% of the variance could be explained by the model for dataset B ($R^2 = 0.3657$).

**Table 2.** Total Model statistics: significance and variance explained

| Predictors | Estimates | Std. Beta | CI | (Std. CI) | p |
|---|---|---|---|---|---|
| (Intercept) | 4.08 | 0.07 | 1.34 – 6.81 | (-0.17 – 0.32) | **0.004** |
| age | -0.06 | -0.19 | -0.10 – -0.02 | (-0.31 – -0.07) | **0.002** |
| sex [female] | -0.23 | -0.15 | -0.55 – 0.09 | (-0.37 – 0.06) | 0.159 |
| STAI trait | -0.02 | -0.08 | -0.06 – 0.02 | (-0.21 – 0.06) | 0.257 |

| | | | | | |
|---|---|---|---|---|---|
| pain cat | 0.08 | 0.26 | 0.04 – 0.13 | (0.12 – 0.41) | **<0.001** |
| mindfulness | -0.23 | -0.14 | -0.44 – -0.02 | (-0.26 – -0.01) | **0.033** |
| cortisol serum | 0.51 | 0.34 | 0.33 – 0.69 | (0.22 – 0.46) | **<0.001** |

**Random Effects**

| | |
|---|---|
| $\sigma^2$ | 1.20 |
| $\tau_{00\ hospital}$ | 0.17 |
| ICC | 0.13 |
| N $_{hospital}$ | 10 |
| Observations | 200 |
| Marginal $R^2$ / Conditional $R^2$ | **0.385 / 0.463** |

Note: $\sigma^2$ = residual; $\tau_{00\ hospital}$ = estimates of hospital; ICC = intraclass correlation coefficient.

The following figures displayed the separate fitted regression lines for the hospitals from the mixed model, including only the most influential predictor – the cortisol serum level – in the model. **Figure 1** included regression lines from the random intercept model and **Figure 2** from the random slope model.

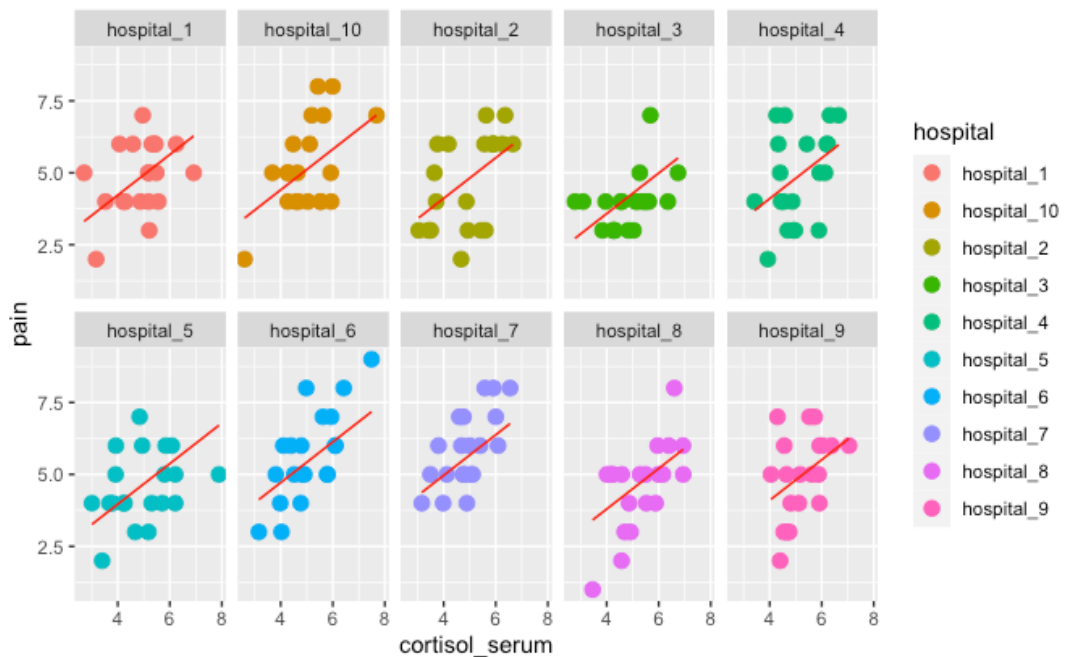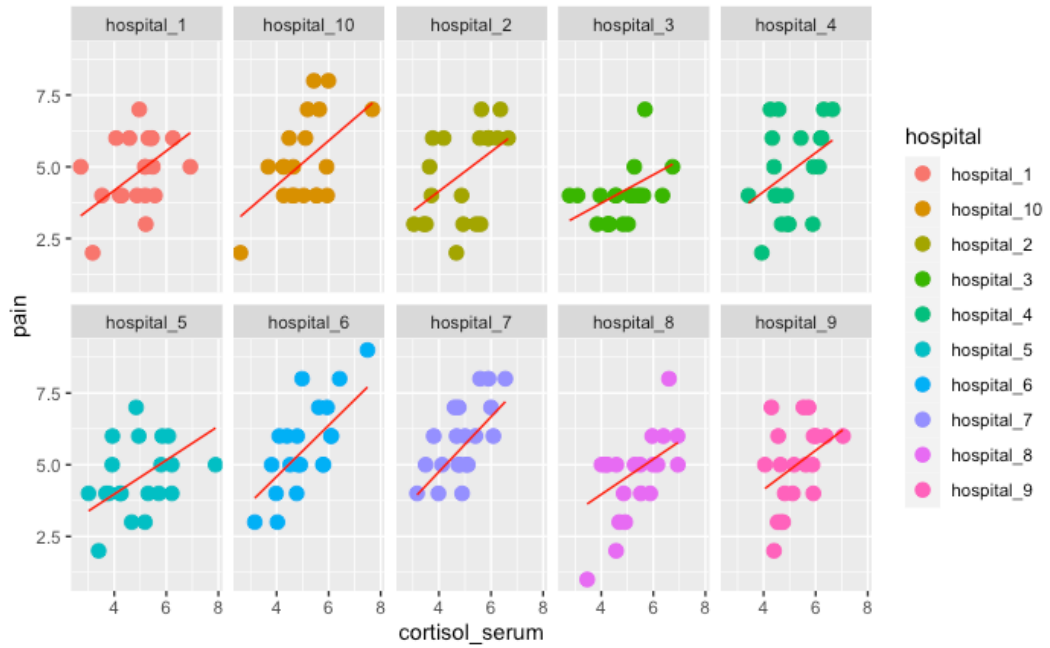**Figure 1.** The fitted regression lines for the hospitals from the random intercept model

**Figure 2.** The fitted regression lines for the hospitals from the random slope model

## 4. Discussion

The $R^2$ obtained on dataset B (0.366) was close to the marginal $R^2$ (0.385) on dataset A but smaller than the conditional $R^2$ (0.463). This is probably because when calculating the total sum of squared differences (TSS) of dataset B, I added the random effect into the calculation without fixed predictors (check the code), so the final result of $R^2$ of dataset B mainly explained the variance of fixed effects of the predictors on the data.

The difference between the predictions of the two models seemed unremarkable from the figures above, but the random intercept model produced a slightly better model fit according to the cAIC (cAIC intercept = 664.54, cAIC slope = 681.61). Thus, the random intercept model results would be presented as follows. The comparison table of the two models is listed in the Appendices.

**Table 3.** Statistics of the random intercept model with the most influential predictor

|  | b | 95%CI lb | 95%CI ub | Std.Beta | p-value |
|---|---|---|---|---|---|
| **(Intercept)** | 1.38 | 0.45 | 2.31 | 0 | .004 |
| **cortisol_serum** | 0.70 | 0.53 | 0.88 | 0.47 | <.001 |

## 5. Reference

Nakagawa, S., Johnson, P.C.D., & Schielzeth, H. (2017). The coefficient of determination R2 and intraclass correlation coefficient from generalized linear mixed-effects models revisited and expanded. *J. R. Soc. Interface 14*.
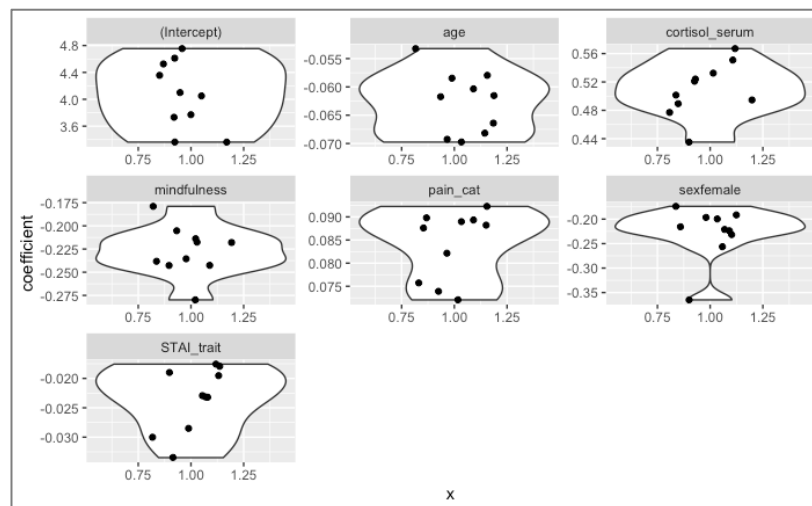
Nakagawa, S. & Schielzeth, H. (2013). A general and simple method for obtaining R2 from generalized linear mixed-effects models. *Methods Ecol Evol, 4*(2), 133-142.

R Core Team. (2020). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL https://www.R-project.org/.
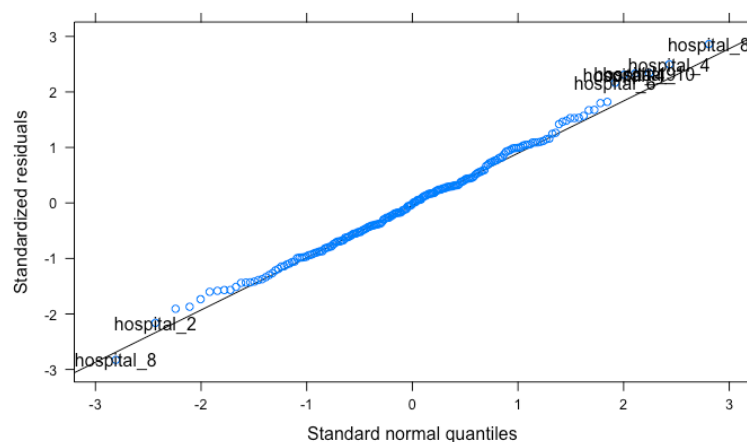
## 6. Appendices

6.1. Model diagnostics plots:

6.1.1. Checking influential outliers: the plots do not indicate extreme influential cases
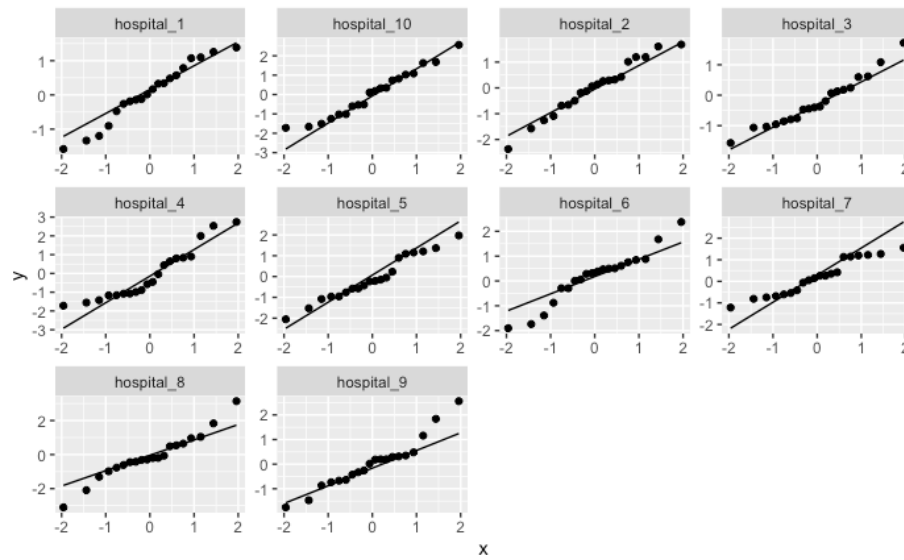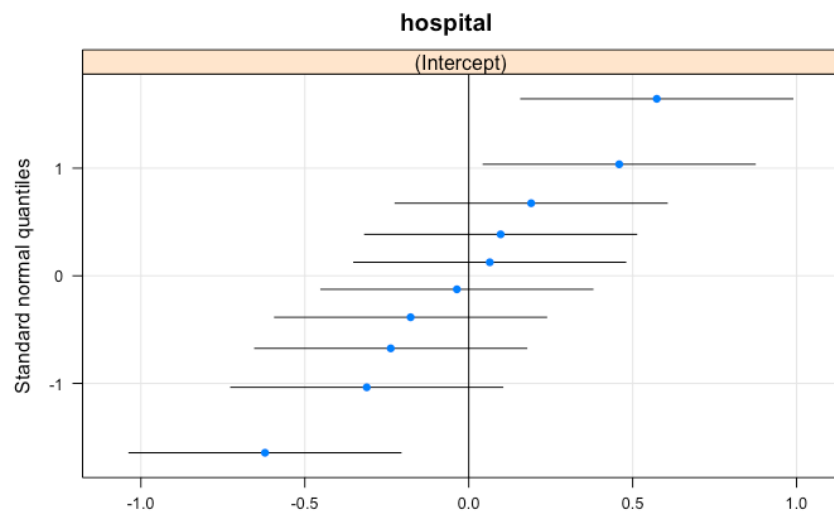


6.1.2. Checking normality

Normality of all residuals: points fall nicely onto the line – normally distributed
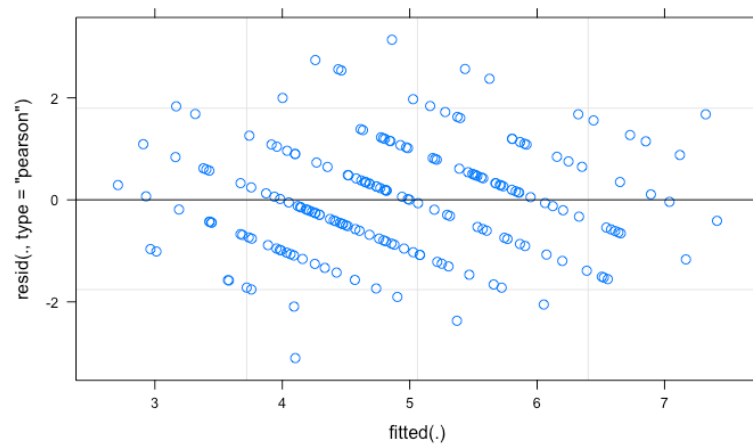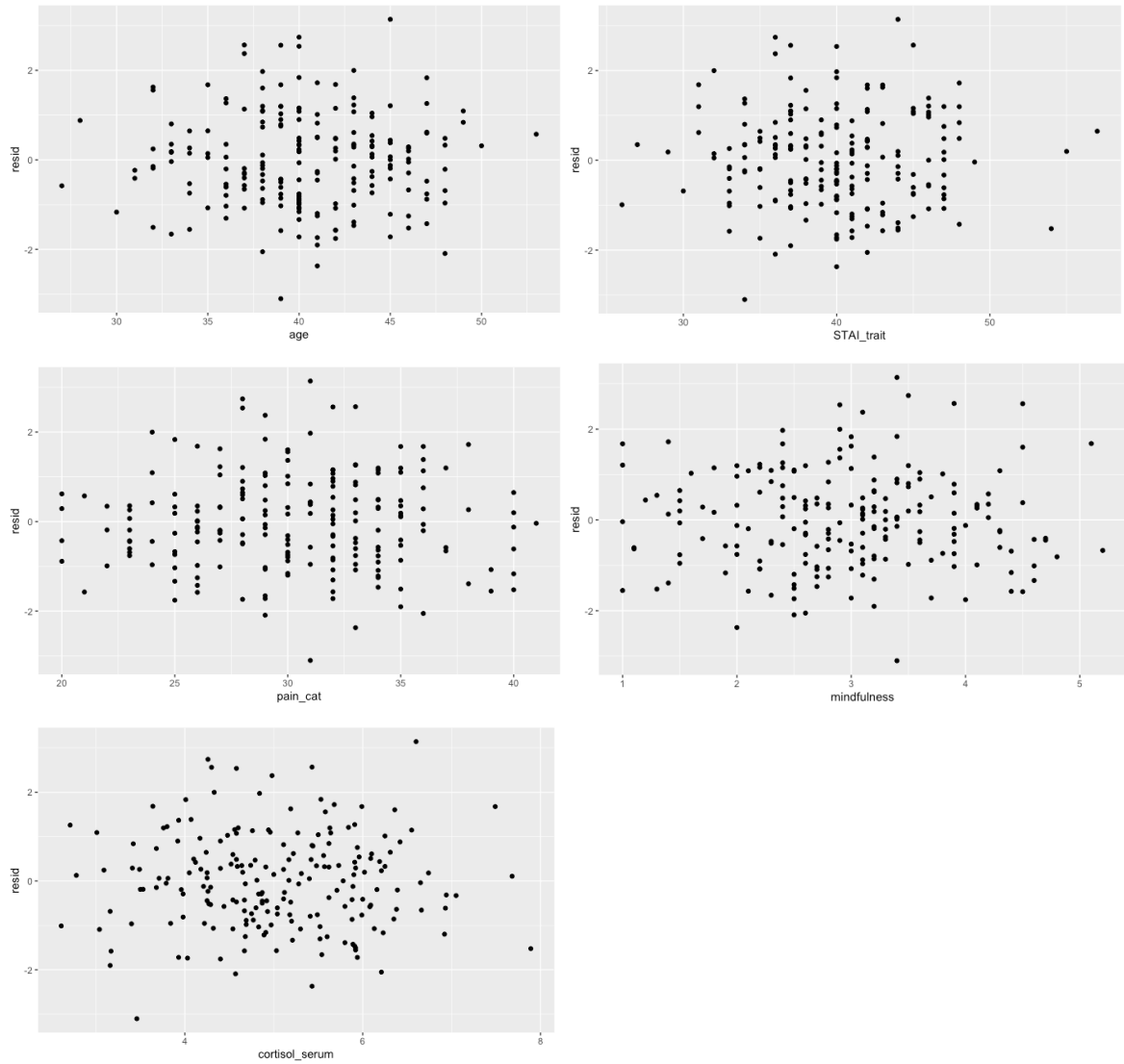
Normality of residuals within clusters:



Normality of random effects: the points roughly fit on a straight line



6.1.3. Checking linearity: the plot shows a non-linear relationship with the residuals.

6.2. Comparison between the random intercept model and the random slope model with only the most influential predictor.

**Table 4.** Statistics of the random intercept model with the most influential predictor

| Predictors | Pain (random intercept) | | | Pain (random slope) | | |
|---|---|---|---|---|---|---|
| | *Estimates* | *CI* | *p* | *Estimates* | *CI* | *p* |
| (Intercept) | 1.38 | 0.44 – 2.32 | **0.004** | 1.37 | 0.45 – 2.28 | **0.004** |
| cortisol serum | 0.70 | 0.53 – 0.88 | **<0.001** | 0.71 | 0.51 – 0.91 | **<0.001** |
| **Random Effects** | | | | | | |
| $\sigma^2$ | 1.53 | | | 1.50 | | |
| $\tau_{00}$ | 0.23 hospital | | | 0.13 hospital | | |
| $\tau_{11}$ | | | | 0.03 hospital.cortisol_serum | | |

| | | -0.92 $_{hospital}$ |
|---|---|---|
| $\rho_{01}$ | | |
| ICC | 0.13 | 0.16 |
| N | 10 $_{hospital}$ | 10 $_{hospital}$ |
| Observations | 200 | 200 |
| Marginal R$^2$ / Conditional R$^2$ | 0.221 / 0.321 | 0.220 / 0.342 |

Note: $\sigma^2$ = residual; $\tau_{00\ hospital}$ = estimates of hospital; ICC = intraclass correlation coefficient.

## Link to the codes of report #2:

https://github.com/Dino-Lu/Lab-Assignments-for-SIMM61

https://github.com/Dino-Lu/Lab-Assignments-for-SIMM61/blob/main/Lab%20Report%232%20-%20Jinxuan%20Lu.Rmd

## 1. Introduction

This study aims to explore the underlying factors that govern individuals' attitudes towards animal rights and animal research project and then to use the underlying factors to predict whether or not a particular person is conservative or liberal.

## 2. Data and Method

The data used in this study is the 28-item Animal Rights Scale (ARS, Wuensch, Jenkins & Poteat, 2002). The survey questionnaire comprised 28 items on similar scales based on several factors identified from the literature. After tackling missing values, 149 out of 154 valid individuals with 120 females and 29 males were finally included in the data.

Since the major goal of this study is to explore the latent factors underlying the response items rather than simply reduce dimensions into a smaller number of components, I used the Exploratory Factor Analysis (EFA) in the first part. For the second part, a linear regression analysis was conducted to predict the level of liberal based on factor scores.
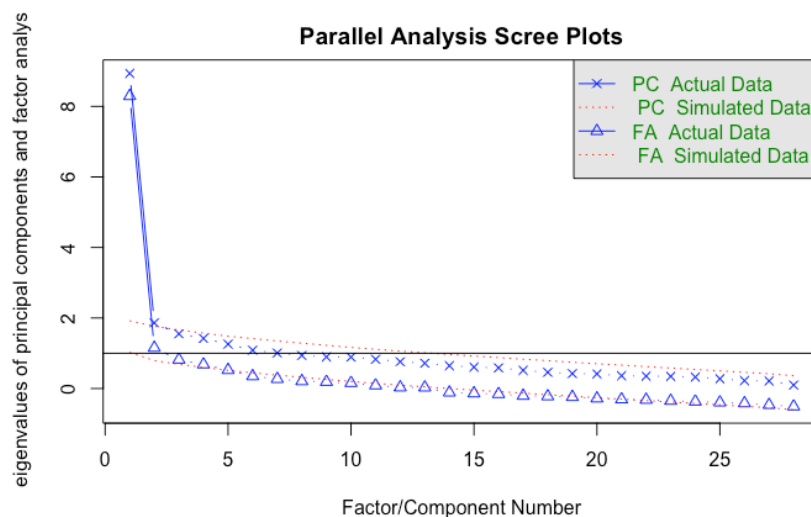
## 3. Results

According to the descriptive statistics result, there were 23% of missing values in the original dataset, which was high (>10%), but they concentratedly located across five rows, so I removed the five rows with missing values instead of using imputation methods to replace them. The items in the data all had acceptable skewness (-2 to +2) and kurtosis (-7 to +7) (Hair et al., 2010; Bryne, 2010). Using Mahalanobis squared distance to check outliers in the correlation matrix, 5 multivariate outliers were identified (Field, 2009). Then I conducted a sensitive analysis by comparing the results of two factor-analysis models with one excluding outlier rows. However, the factor loadings in the non-outlier model did not really make sense considering the theory behind items, so I kept the original dataset without removing the outliers.

The reliability of the correlation matrix was good, with Cronbach's alpha being 0.92, showing good internal consistency. Bartlett's sphericity test showed a small value (<0.05) of the significance level, which rejected the null hypothesis that the correlation matrix was an identity matrix. The total KMO was 0.87 (> 0.6), indicating that, based on this test, we could probably conduct a factor analysis (Kaiser, 1974). Besides, all the results of the Henze-Zirkler test and the multivariate skewedness and kurtosis tests were in a significant level (p<0.05),

which means the rejection of the null hypotheses, indicating violation of the multivariate normality assumption. Therefore, this study would use the principal axis factor extraction method instead of using maximum likelihood estimation.

Among different tests, except for the parallel analysis, which suggested 4 factors (the plot suggested 2 though), the Scree test, VSS, and MAP criterion all suggested 2 factors to retain, so I listened to this advice and followed the 2-factor structure. After I deleted several unneeded items, the tests on the final dataset also suggested 3 factors as a choice, so I checked 3-factor structure as well. However, the 2-factor structure made more sense theoretically than the 3-factor structure, even though the 3-factor structure might have a slightly higher average communality (0.424) than the 2-factor structure (0.380).

**Figure 1.** Parallel analysis scree plot



The post-extraction eigenvalues and variance explained were shown in **Table 1**.

**Table 1.** Eigenvalues, Variance Explained, and Factor Correlations

Eigenvalues, Variance Explained, and Factor Correlations for Rotated Factor Solution:

| Property | Factor_1 | Factor_2 |
|---|---|---|
| **SS loadings (eigenvalues)** | 4.223 | 4.136 |
| **Proportion Var** | 0.192 | 0.188 |
| **Cumulative Var** | 0.192 | **0.380** |
| **Proportion Explained** | 0.505 | 0.495 |

Eigenvalues, Variance Explained, and Factor Correlations for Rotated Factor Solution:

| Property | Factor_1 | Factor_2 |
|---|---|---|
| **Cumulative Proportion** | 0.505 | 1.000 |
| **Factor_1** | 1.000 | **0.650** |
| **Factor_2** | **0.650** | 1.000 |

This study used the oblique rotation method (Promax) rather than the orthogonal rotation, and the rotated results were examined for simple structure, following Kline's (2002, p. 65) criteria that each factor had several high loadings, with the rest being zero or less than ±0.10 (see **Table 2**). I used promax because factors in this study tended to be correlated since all the items were about animal rights. This was also evidenced by the resulting correlation matrix for the factors (see **Table 1**): the lowest correlation between factors was 0.65, which exceeded the Tabachnick and Fiddell threshold of 0.32, and this basically meant that there was more than 10% overlap in variance among factors, enough variance to warrant oblique rotation.

**Table 2.** Final structure of the factor analysis

| | Factor_1 (AnimalResearch _concern) | Factor_2 (AnimalRights _concern) | Communality ($h^2$) | Uniqueness ($1-h^2$) | Complexity |
|---|---|---|---|---|---|
| | **Factor analysis results** | | | | |
| ar6 | 0.924 | *-0.158* | 0.69 | 0.31 | 1.06 |
| ar27 | 0.726 | *-0.137* | 0.42 | 0.58 | 1.07 |
| ar2 | 0.685 | *0.071* | 0.54 | 0.46 | 1.02 |
| ar17 | 0.642 | *-0.144* | 0.31 | 0.69 | 1.10 |
| ar9 | 0.529 | *0.055* | 0.32 | 0.68 | 1.02 |
| ar18 | 0.478 | *-0.105* | 0.17 | 0.83 | 1.10 |
| ar20 | 0.471 | *0.162* | 0.35 | 0.65 | 1.23 |
| ar12 | 0.460 | *0.295* | 0.48 | 0.52 | 1.70 |
| ar15 | 0.428 | *0.324* | 0.47 | 0.53 | 1.86 |

| | Factor analysis results | | | | |
| --- | --- | --- | --- | --- | --- |
| | Factor_1 (AnimalResearch _concern) | Factor_2 (AnimalRights _concern) | Communality (h^2) | Uniqueness (1-h^2) | Complexity |
| ar19 | -0.410 | *-0.154* | 0.27 | 0.73 | 1.28 |
| ar21 | -0.387 | *-0.101* | 0.21 | 0.79 | 1.13 |
| ar10 | *-0.116* | 0.802 | 0.54 | 0.46 | 1.04 |
| ar7 | *-0.011* | 0.685 | 0.46 | 0.54 | 1.00 |
| ar5 | *0.147* | 0.672 | 0.60 | 0.40 | 1.10 |
| ar24 | *0.187* | -0.637 | 0.29 | 0.71 | 1.17 |
| ar4 | *0.030* | 0.588 | 0.37 | 0.63 | 1.01 |
| ar13 | *0.305* | 0.567 | 0.64 | 0.36 | 1.53 |
| ar26 | *0.167* | 0.517 | 0.41 | 0.59 | 1.21 |
| ar28 | *0.124* | -0.452 | 0.15 | 0.85 | 1.15 |
| ar23 | *0.231* | 0.418 | 0.35 | 0.65 | 1.56 |
| ar22 | *0.109* | 0.369 | 0.20 | 0.80 | 1.17 |
| ar3 | *0.006* | 0.365 | 0.14 | 0.86 | 1.00 |

Note: Complexity, specifically "Hoffman's index of complexity for each item," means how much an item reflects a single construct.

The final factor structure with communalities and loadings of each item was shown in order in Table 2, and the name of each factor was also listed above. Each factor included exactly 11 items, and 6 out of 28 items were excluded from the structure in the end due to the low loadings (<0.32) and cross-loading problems. Theoretically, ar1 and ar8 were both "indirect" questions, which we didn't see a clear construct that defined the two. And ar11 concerning insect pests was rather irrelevant to the main topic of this study. In addition, ar14 "I would rather see human die than to see animals used in research" was too extreme an item, and ar16 "God put animals on Earth for man to use" was too general and wasn't captured by any of the factors. And last, ar25 was excluded because of its low loading.

The summary of the simple linear regression and the moderated multiple regression results were shown below, respectively. Since I assumed that the two factors were intended to be correlated, it was reasonable to conduct a moderated multiple regression as well. Based on the results, the moderated multiple regression model was better ($R^2 = 0.103$). The factor of animal rights concern was the most influential predictor and had a positive effect on how liberal a person is ($\beta = 0.29$, $p < 0.01$). Besides, the result also showed a negative moderating effect of concerns about animal research on the relationship between general animal rights concerns and a person's liberal level ($\beta = -0.17$, $p < 0.01$).

**Table 3.** Simple linear regression results

| | liberal | | |
|---|---|---|---|
| *Predictors* | *Estimates* | *CI* | *p* |
| (Intercept) | 2.91 | 2.77 – 3.04 | **<0.001** |
| AnimalResearch concern | -0.09 | -0.27 – 0.08 | 0.296 |
| AnimalRights concern | 0.27 | 0.09 – 0.45 | **0.004** |
| Observations | 149 | | |
| $R^2$ / $R^2$ adjusted | 0.065 / 0.052 | | |

**Table 4.** Moderated multiple regression results

| | liberal | | |
|---|---|---|---|
| *Predictors* | *Estimates* | *CI* | *p* |
| (Intercept) | 3.02 | 2.87 – 3.17 | **<0.001** |
| AnimalRights concern | 0.29 | 0.12 – 0.47 | **0.001** |
| AnimalResearch concern | -0.06 | -0.23 – 0.12 | 0.520 |
| AnimalRights concern * AnimalResearch concern | -0.17 | -0.29 – -0.06 | **0.003** |
| Observations | 149 | | |
| $R^2$ / $R^2$ adjusted | 0.122 / 0.103 | | |

## 4. Discussion

One limitation of this study is that the average communality was only 0.38, still below the expected 0.6. And I need more practice on exploratory analysis, such as data cleaning and data management. It took me a lot of time at the current stage, and still, many things are unclear.

Moreover, Exploratory Factor Analysis is a complex topic for me, and I need more time to do research on it and practice to be more familiar with it. For example, I read many articles on the rotation method, and turned out there was no one commonly agreed method. And for the factor loadings, there was also no uniform standard on how good the loadings are for research.

**Reference**

Byrne, B.M. (2010). Structural Equation Modeling with AMOS: Basic Concepts, Applications, and Programming (2nd Edition). *New York: Taylor and Francis Group Publication.*

Field, A. P. (2009). Discovering statistics using SPSS: (And sex and drugs and rock 'n' roll). *Thousand Oaks, CA: Sage.*

Hair, J., Black, W.C., Babin, B. J., & Anderson, R.E. (2010). Multivariate Data Analysis (7th Edition). *NJ: Prentice-Hall Publication.*

R Core Team. (2020). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL https://www.R-project.org/.

Tabachnick, B. G. & Fidell, L. S. (2007). Using multivariate statistics (5th ed.). *Boston, MA: Allyn & Bacon*

Wuensch, K. L., Jenkins, K. W., & Poteat, G. M. (2002). Misanthropy, idealism, and attitudes towards animals. *Anthrozoös*, *15*, 139-149

**Link to the codes of report #3:**

https://github.com/Dino-Lu/Lab-Assignments-for-SIMM61

https://github.com/Dino-Lu/Lab-Assignments-for-SIMM61/blob/main/Lab%20Report%233%20Jinxuan%20Lu.Rmd
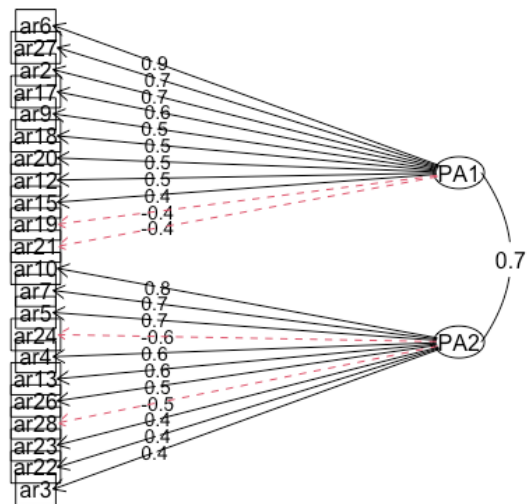
# 5. Appendices

**Figure 2.** Factor analysis diagram



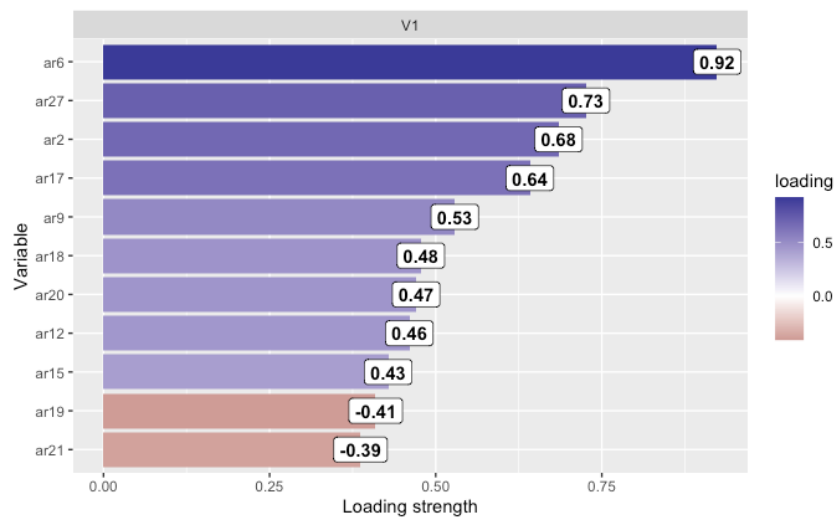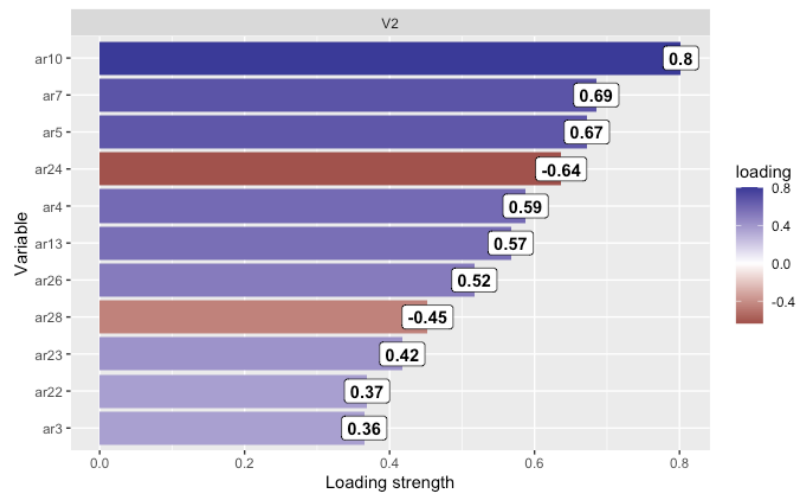**Figure 3.** Loadings diagram for Factor 1 – AnimalResearch_concern



**Figure 4.** Loadings diagram for Factor 1 – AnimalRights_concern

## 1. Introduction

The aim of this study is to assess the fit of a theoretical factor structure of a mental ability test and the relationships between different test scores. And the results of three tasks will be covered in this report, respectively.

## 2. Data and Method

The data used in this study consists of mental ability test scores of 301 seventh- and eighth-grade children from two different schools (Holzinger & Swineford, 1939). There are also other demographical variables, but only variables about different test scores will be used in this study.

As for the method, this study will use the structural equation modeling (SEM) analysis. While the first and second tasks include latent variables, the last task is primarily about path analysis, especially the mediation analysis.

## 3. Results

3.1. Model_A specification and analysis

Firstly, based on the theory, I specified a structural equation model (Model_A), which includes three latent factors: "Visper", "Verbal", and "Procespd". The result shows that the degree of freedom of Model_A is 32, which is over-identified, and the model is solvable.

By checking the kurtosis and skewness, I found a violation of the assumption of multivariate normality distribution ($p<0.05$). Thus, I chose to use bootstrapped ML estimator to fit the model since it can provide bootstrapped standard errors and p-values. And according to an article from Falk, with nonnormal data, the bootstrap method should be preferred than the robust likelihood-based approach and robust standard errors (Falk, 2018).

The fit statistics of Model_A are shown in **Table 1**. Based on Hu and Bentler's (1999) cutoff criteria, the model fit indices are acceptable (RMSEA = 0.076) or slightly less than the good fit values (p of Chi^2 < 0.05, CFI = 0.940, TLI = 0.916).

**Table 1.** Fit statistics of Model_A

|  | df | Chisq | p-value | TLI | CFI | RMSEA | 90%CI lb | 90%CI ub |
|---|---|---|---|---|---|---|---|---|
| **Model_A** | 32 | 87.964 | 0.000 | 0.916 | 0.940 | 0.076 | 0.057 | 0.095 |
| **cutoff** | > 0 |  | > 0.05 | > 0.95 | >= 0.95 | <= 0.05 (0.05-0.08: acceptable) | | |

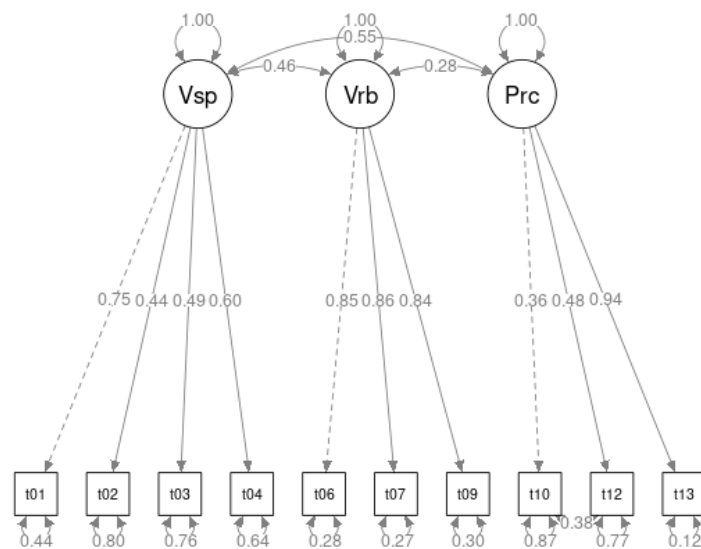3.2.Model_B specification & comparison with Model_A

Considering the correlation between t10_addition and t12_countdot not caused by "Procespd," I specified the Model_B including this new path. The AIC, BIC, and SABIC of Model_B are all smaller than Model_A, and can be seen as significantly different from Model_A (smaller more than 2 points). The Chi-squared difference between the two models is significant, and Model_B has the lower value. Besides, the difference in CFI is more than 0.01 between the two models, and Model_B has the higher value of CFI, which is more than the cutoff (0.95). Therefore, based on the above test statistics, we can conclude that there is a significant difference between the two models, and Model_B has a better fit (see **Table 2**).

**Table 2.** Comparison between Model_A & Model_B

|  | df | AIC | BIC | SABIC | CFI | Chisq | Chisq diff | p |
|---|---|---|---|---|---|---|---|---|
| **Model_A** | 32 | 8296.9 | 8382.1 | 8309.2 | 0.940 | 87.964 | 31.204 | <0.001 |
| **Model_B** | 31 | 8267.7 | 8356.6 | 8280.5 | 0.972 | 56.759 | | |

Among the manifest variables t01, t02, t03, and t04, we can tell from **Figure 1** that the t02 is the least influenced variable by Visual perception ability ("Visper"). Figure 1 shows the factor loadings (standardized estimates) of the latent factors on each of their manifest variables, and it's clear that the factor loading of Visper is strongest for t01_addition (0.75), and weakest for t02_cubes (0.44).

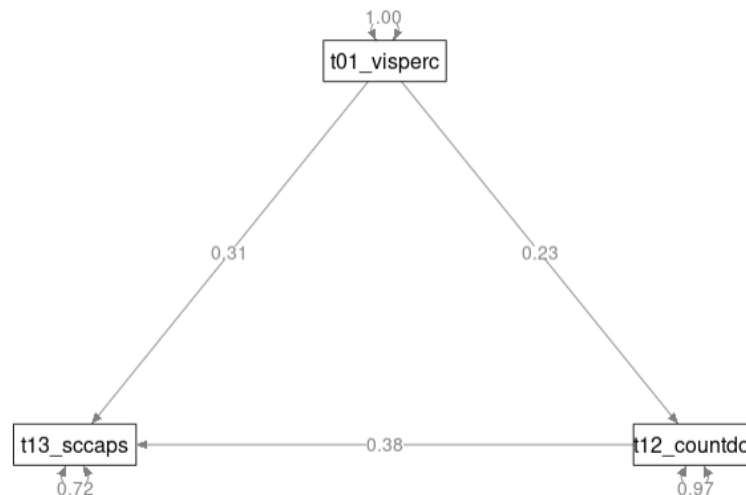**Figure 1.** Path diagram of Model_B

This result can also be verified by the squared multiple correlations (R^2) for endogenous variables (t01-t13), which shows that only 19.5% of the variance of t02_cubes is explained in the model (R^2 = 0.195).

3.3. The mediation model – Path Analysis

I reproduced the mediation model and got exactly the same graph as the following. The graph shows an indirect effect of t01_visperc, via a mediating variable t12_countdot, on t13_sccaps. The direct effect of t01 on t13 is 0.31, and the indirect effect of t01 on t13 can be calculated as: 0.23 * 0.38, which is 0.087. Thus, if the independent variable t01 increases by 1 unit, the dependent variable t13 would be expected to increase by 0.397 units (0.23*0.38 + 0.31). Also, the analysis results indicate that the indirect effect of t01 on t13 mediated through r12 is significant, so the model supports a mediation effect.

**Figure 2.** Path diagram of the mediation analysis model
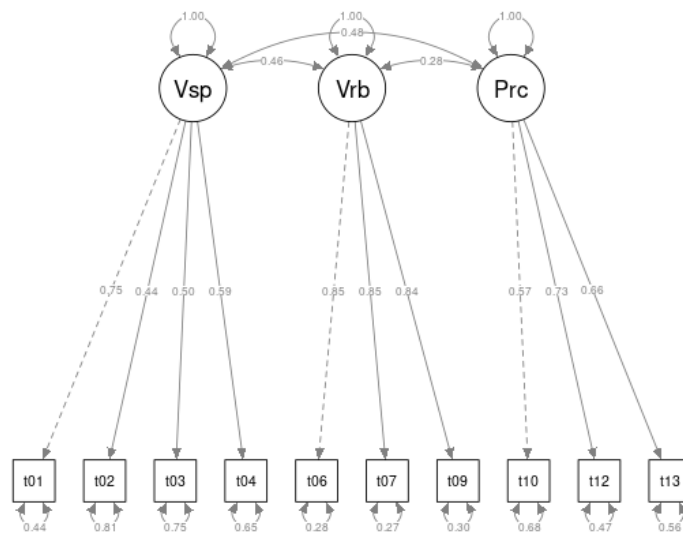


## 4. Reference

Falk, Carl F. (2018). Are Robust Standard Errors the Best Approach for Interval Estimation with Nonnormal Data in Structural Equation Modeling? *Structural Equation Modeling: A Multidisciplinary Journal, 25*(2), 244-266. DOI: 10.1080/10705511.2017.1367254

Li-tze, Hu & Peter M. Bentler. (1999). Cutoff criteria for fit indexes in covariance structure analysis: Conventional criteria versus new alternatives. *Structural Equation Modeling: A Multidisciplinary Journal, 6*(1), 1-55. DOI: 10.1080/10705519909540118

R Core Team. (2020). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL https://www.R-project.org/.

# 5. Appendices

**Figure 3.** Path diagram of Model_A



**Link to the codes of report #4:**

https://github.com/Dino-Lu/Lab-Assignments-for-SIMM61

https://github.com/Dino-Lu/Lab-Assignments-for-

SIMM61/blob/main/Lab%20Report%20%234%20-%20Jinxuan%20Lu.Rmd