# Examining the Efficacy of Beat-the-Blues Treatment on Depression Utilizing GEE Modeling and Inverse Probability Weighting

## CHL5222 Correlated Data Analysis - Group C

Victoria Chui | Faculty of Information | victoria.chui@mail.utoronto.ca
Alexander Gebreamlak | Dalla Lana School of Public Health | alex.gebreamlak@mail.utoronto.ca
Yuchen Jiang | Dalla Lana School of Public Health | ycc.jiang@mail.utoronto.ca
Wenhao Wang | Dalla Lana School of Public Health | wenhao.wang@mail.utoronto.ca

April 12, 2024

## Contents

# Abstract

Cognitive-behavioural therapy (CBT) is a widely recommended treatment for those suffering from depression and/or anxiety. However, CBT requires both time and emotional commitment from those who regularly pursue this technique. In this study, we statistically compare Beat-the-Blues - an electronic CBT program - to "treatment as usual", or standard depression/anxiety treatments, over a two-month-long treatment phase and a subsequent six-month-long follow-up phase. We had data from 48 participants who were assigned to the control (or standard treatment) group and 52 participants who completed two months of Beat-the-Blues. The severity of depression/anxiety of participants was measured using the Beck Depression Inventory. Using general estimating equations, and accommodating missing data using inverse probability weighting, the population trajectories of the control and treatment groups were compared, with the goal of determining if a treatment and time interaction effect was significant. The rate ratio between the treatment and control groups for average BDI was 0.77, with a 95% confidence interval between [0.52,0.99], indicating a significantly smaller population BDI for those completing Beat-the-Blues. As well, the treatment and time interaction effect was significant at 0.05 level, with a decrease in BDI of 0.88 for the Beat the Blues group with each month's progression. Limitations include frequently missing data points and a broad definition of "treatment as usual" for participants in the control group.

GitHub link: https://github.com/victoriachui/chl5222

# Introduction

## 1.1 Beat-the-Blues

Cognitive-behavioral therapy (CBT) is an evidence based psychotherapy that aims to help individuals develop skills for maintaining health by changing inaccurate beliefs and focusing on how our perceptions shape our emotions.[1] It is widely used and often as effective as medication, benefiting from format flexibility and shorter duration, yet it requires substantial patient effort and may not be ideal for individuals with complex needs or in addressing root causes.[2,3] Research shows mixed results on effectiveness and cost-efficiency of electronic versus face-to-face therapy. With some studies suggesting eCBT might outperform face-to-face CBT in symptom reduction, whereas others highlight the superior clinical effectiveness of traditional face-to-face CBT.[4]

Beating the Blues (BtheB) is a self-help eCBT program featuring eight weekly customized sessions lasting 50 minutes each, where after initial supervised session, interaction with the research coordinator is minimal.[5] In a pragmatic trial comparing the CCBT program to previous treatments for depression or anxiety, 80% of users providing feedback stated that Beating the Blues was as good as or better than their prior treatment experiences.[6]

## 1.2 The Beck Depression Inventory

Introduced in 1961, the Beck Depression Inventory (BDI) serves as a critical tool in both clinical and research settings to assess depression severity.[7] This 21-item self-report questionnaire, which takes about 5 minutes to complete either via computer or paper and pencil, is designed for both clinical patients and non-clinical individuals, aiming to reflect the intensity of depression with higher scores indicating greater depression.[7] Specific cut-offs vary but generally, scores ranging from 0–9 suggest no depression, 10–18 signify mild to moderate depression, 19–29 point to moderate to severe depression, and 30–63 denote severe depression.[8] Regarding the psychometric properties of the Beck Depression Inventory, based on four studies, the sensitivity was 0.85 (95% CI, 0.79 - 0.90) and the specificity was 0.83 (95% CI, 0.70 - 0.91).[9]

### 1.3 Research Question

A randomized control trial was conducted comparing Beating the Blues psychotherapy to the standard treatment as usual (TAU) in patients suffering from anxiety and/or depression, with the primary response being Beck Depression Inventory scores. Therefore, the aim of this paper is to examine the efficacy of BtheB (in both the active treatment phase and the post-treatment follow-up) measured by BDI, controlling for potential confounding variables.

# Methods

## 2.1 Exploratory Data Analysis

Dataset for depression study contains 100 observations in total. To be specific, 100 patients suffering from anxiety or depression were invited to join this randomized clinical trial and their Beck Depression Inventory(BDI) levels were recorded at 5 different time occasions before and after the active treatment phase. Here, BDI level is the primary response variable, and it belongs to one measurement index in evaluating patients' depression level. Baseline BDI level of patients was measured at the time of study entry. Then, the second BDI level of patients was measured at the end of treatment phase. Since the duration of treatment including BtheB therapy and TAU therapy lasts for 2 months, clearly, the second BDI level of patients was measured 2 months after the time of study entry. For the rest of three measurements of BDI level, they were recorded at 4 months, 6 months, and 8 months after the time of study entry during the post-treatment period. In terms of possible predictors which might be used in the regression model, they are drug, length as well as treatment. Drug is a binary variable, and it has a value of 1 when a patient was prescribed concomitant drug therapy along with current treatment, and 0 otherwise. Length is a binary variable as well. It has a value of 1 when the current lasting time of depression at the study entry is longer than 6 months, and 0 otherwise. Treatment is a categorical variable, and it represents the treatment assigned to each patient in this depression study.

One thing needs to be emphasized here is that there are quite a lot of missing values in the outcome of interests(i.e. BDI level). In other words, patients enrolled in the clinical trial have different numbers of measurement records for BDI level over time. Statistically, 3% of patients quit this trial after the first measurement of BDI, 24% of patients quit this trial after the second measurement of BDI, 15% of patients quit this trial after the third measurement of BDI, 6% of patients quit this trial after the fourth measurement of BDI and only 52% of patients complete whole trial for depression study. That is, nearly half of patients have incomplete measurement records for BDI level and this will greatly influence the choice of modeling approaches in the following part.

## 2.2 Assumption Checks

### 2.2.1 Check for BDI Distributions

**Distribution of Baseline BDI level**



**Distribution of BDI level 2M**



**Distribution of BDI level 4M**



**Distribution of BDI level 6M**
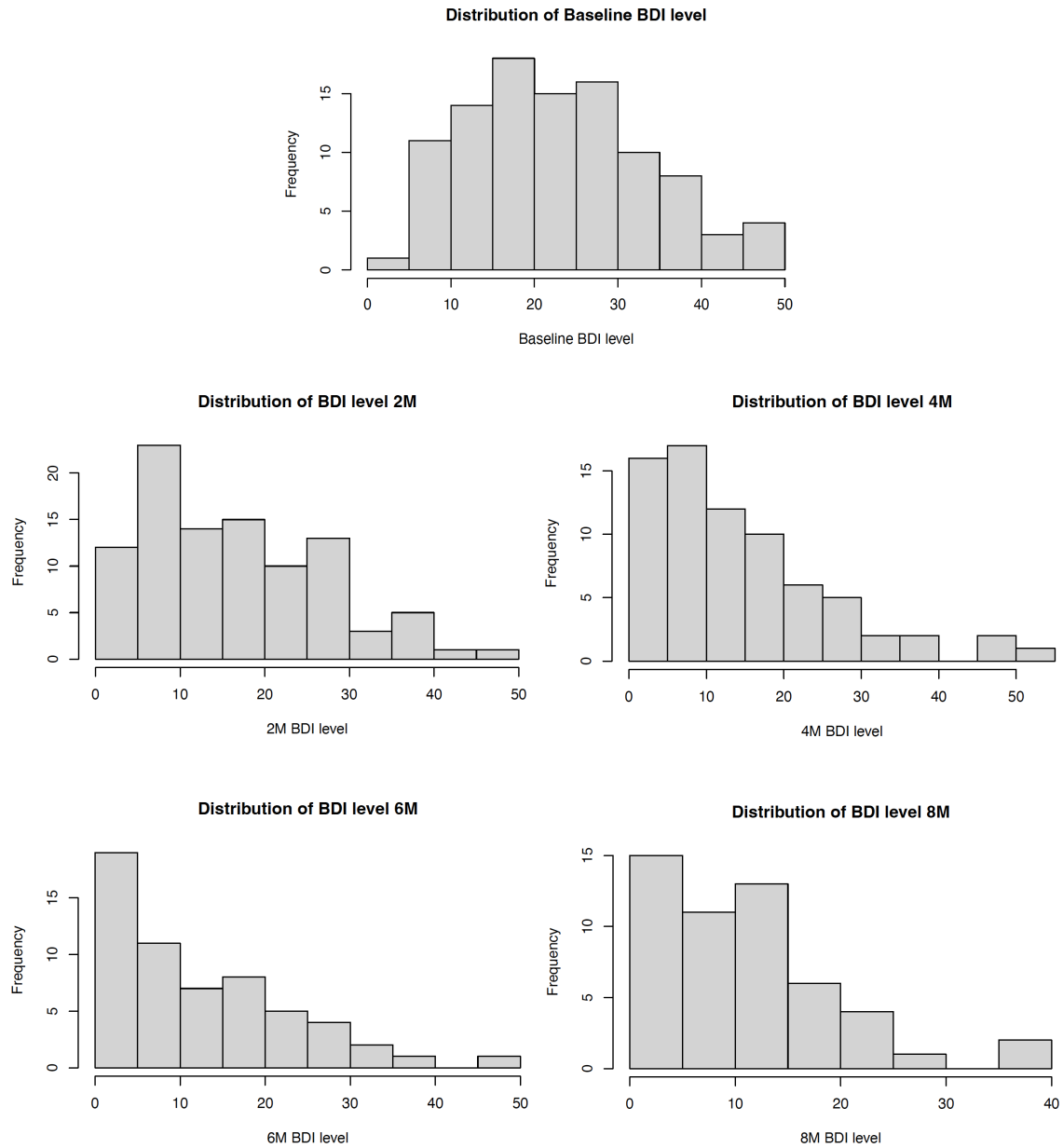


**Distribution of BDI level 8M**



Figure 1: BDI Distributions

These figures provide the distribution of BDI measurements for different months of visits. Baseline BDI levels appear to be approaching a normal distribution. Post-intervention, the BDI becomes more right-skewed (decreasing) over time.
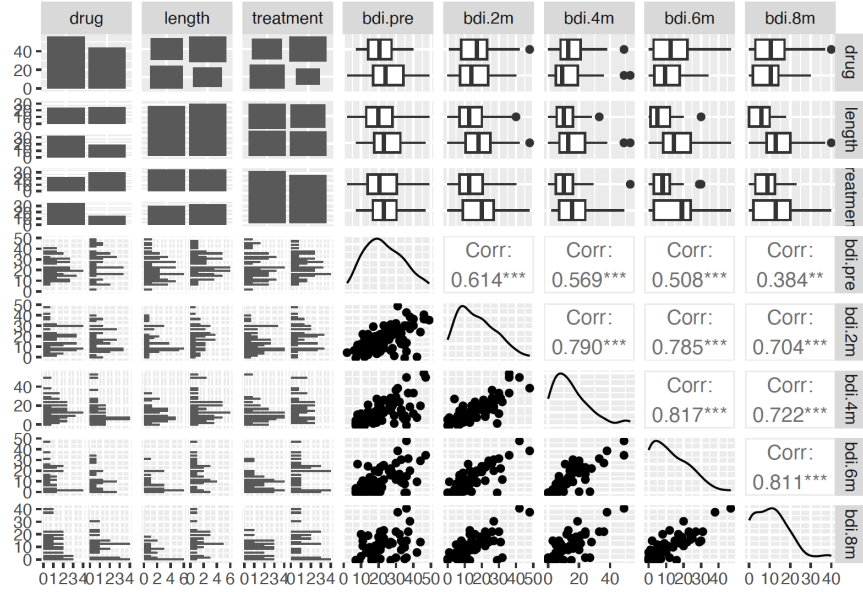
## 2.2.2 Check for Independence



Figure 2: Correlation Plot

Upon examination of the mosaic plots, located in the upper left corner of the matrix, no apparent variation in the size of the boxes is observed, which represent the permutations of categories across the drug, length, and treatment variables. This uniformity suggests a balanced distribution of observations across each subgroup, with no evidence of significant association between the two variables examined. Variables assume to be independent with each other.

Next, the scatter plots situated in the lower right corner of the matrix reveal the dynamics of the Beck Depression Inventory (BDI) measures over the study. There is a demonstrable linear relationship between baseline BDI levels and those measured at 2 months, which weakens over time. The relationship between baseline BDI levels and those at 8 months post-baseline is comparatively weaker, indicating a decrease in the strength of the relationship as the separation between assessments increases.

The correlation coefficients displayed in the central right area of the matrix further substantiate these findings. Correlations approaching 1 indicates a strong positive linear relationship, whereas those closing -1 denote a strong negative linear relationship.

Consistently, we observe a diminishing in correlation coefficients as the time interval between assessments increase, suggesting that the strength of the correlation between the BDI levels of different visits decreases as the time interval extends.

## 2.3 Missing Data Visualization

Approximately 50% of observations are missing Beck Depression Inventory (BDI) measurements across four visits. Understanding the nature of these missing values is critical for selecting the appropriate data handling techniques and for determining the robustness of subsequent analyses.
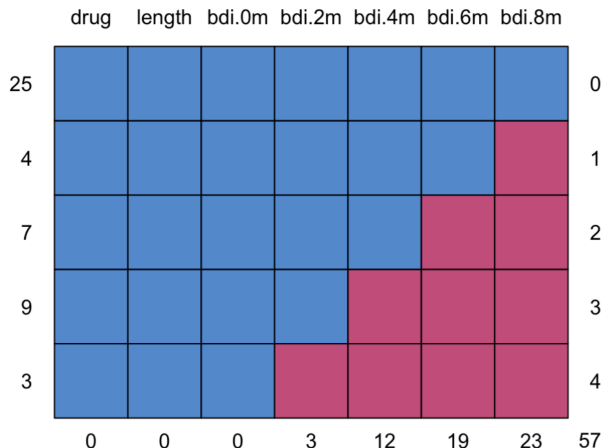


Figure 3: Monotone Missing Data Pattern

The pattern of missingness, as evidenced in the figure, appears to exhibit a monotonic nature. This indicates that once an observation lacks a measurement at a given time point, it continues to be missing for all subsequent time points for that observation, suggesting permanent dropout from the study.
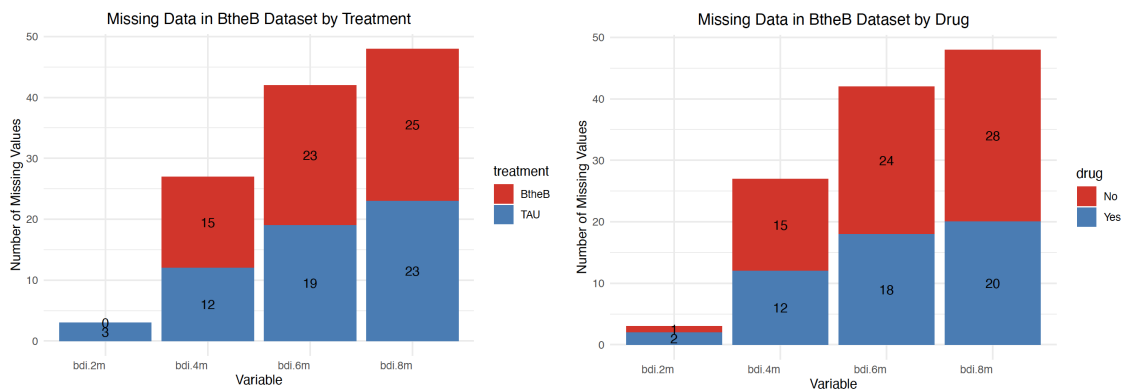


Figure 4: Missingness by Treatment and Drug

There are generally three types of missing data: Missing Completely at Random (MCAR), Missing at Random (MAR), and Missing Not at Random (MNAR). MCAR occurs when the probability of missingness is independent of both observed and unobserved data. MAR occurs when missingness is systematically related to observed variables but not the missing data itself. MNAR occurs when the probability of missingness is related to the unobserved data, in other words, the missingness is related to why the value is missing.

From the above two figures, the missing values show an increase over time across the arms and drug groups. It is noticeable that the missing values are equally distributed across the treatment groups (BtheB and TAU) and drug groups (with and without drug), it supports the assumption that data may be MAR, given that the MCAR assumption is often not realistic. A more detailed investigation would be required to

confirm the MAR assumption. However, such an analysis exceeds the scope of this report and will be further elaborated upon in the discussion section.

## 2.4 More Linear Models

### 2.4.1 Linear Spline Model

After examining the plot for the trajectory of mean BDI level between treatment groups, we found the trend of mean BDI level is piecewise linear between any two consecutive measurement occasions. Linear spline models come to mind in this case. However, due to the existence of missing values in BDI level, the linear spline model is not appropriate to be used here since it is highly sensitive to missing values and lacks corresponding mechanisms to deal with these missing values. In detail, if missing values appear in the measurement time point where the knot is put, then this model is likely to produce biased and unstable estimates for slope before and after that time point. Alternatively, a linear spline model assumes the whole dataset is complete and this implies those missing values need to be imputed at first before such a model can be fitted successfully. Since our strategy in dealing with missing values in response variables is omitting them directly, the assumption of using a linear spline model cannot be met in this case.

2.4.2 Linear Mixed Effects Model

Linear mixed models (LMM), also known as mixed-effects models, account for both fixed effects, which are consistent and predictable factors that you expect to influence the response variable across all subjects or groups, and random effects, which are random variations that occur at the group or subject level. LMMs allow you to model random effects, which capture the variations between clusters or subjects that are not explained by the fixed effects. The fixed effects estimates show the expected change in BDI score associated with each variable, holding the other variables constant. The random effects portion of the model indicates that there is a participant-specific random intercept.

## 2.5 General Estimating Equations Modeling

By exploring marginal models, a continuous response can be analyzed with population-level considerations. Since our objective is to examine the efficacy of the BtheB treatment compared to treatment as usual, a population-average oriented modeling strategy is preferable as it allows us to examine trends more broadly than individual-specific linear mixed effects modeling. By excluding random effects and the effect of previous responses, solely the estimates of our covariates will construct and define the mean trend of the treatment groups.

Applying general estimate equation (GEE) modeling to this data, we examine those factors contributing to the mean trajectory changes of BDI between the BtheB and treatment as usual groups. As defined in section 2.1, our data consists of a baseline and end of active phase BDI measurement, as well as three follow-up measurements each two months apart. There are binary indicators for if the participant is taking a drug medication for their depression/anxiety and for if the participant has had depression/anxiety either longer than or less than 6 months prior to the baseline measurement date.

We perform GEE modeling of the BDI measurements from all five timepoints using main effects for Time, Treatment, and Knot. The Time factor is centered at the end of the active phase (two months into the study), where the baseline is Time -2, and the follow-ups are Time 2, 4, and 6. Our Knot exists at the end of the active phase (two months into the study), as we aimed to look for a significant difference in BDI trajectory when comparing active phase to follow-up. The Treatment groups are pre-defined as either BtheB or TAU (treatment as usual). The equation used for GEE modeling showing the main and interaction effects used is found below.

$$\log E(Y_{ij}) = \beta_1 + \beta_2 \text{knot}_{ij} + \beta_3 \text{treatment}_i + \beta_4 \text{time}_{ij} + \beta_5 \text{knot}_{ij} \times \text{time}_{ij} + \beta_6 \text{treatment}_i \times \text{time}_{ij}$$

## 2.6 Inverse Probability Weighting

Since we only have five measurements (all BDI values in the study) from 52 participants, we explored multiple options for handling missing data. Alongside GEE modeling on complete case and all available data (results in Table X), inverse probability weighting (IPW) was used in an attempt to weight the contribution of individuals more coherently to balance the missingness. IPW was chosen as a preferable strategy as it handles monotone, missing at random data well.10 (Seaman & White, 2013) This strategy aims to estimate the probability of a participant leaving the study and reweigh their BDI values appropriately.

To execute IPW, a "dropout model" must be specified according to the covariates available, and a following model must use the "independence" correlation structure. For our data, the "dropout model" was developed using logistic regression and factors including a binary Treatment indicator, alongside previous BDI lag measurements, and an interaction between Treatment and previous BDI indicator. The GEE model of the equation shown above was fit using the reweighted data, as our final conclusive model for this data.

# Results

## 3.1 Model Output

| Coefficients | Estimate | SE | Wald | Pr( $\|W\|$ ) |
|---|---|---|---|---|
| Intercept | 3.00 | 0.105 | 823.0 | 2e-16 |
| Ageknot | 0.0326 | 0.0528 | 0.4 | 0.537 |
| Treatment | -0.331 | 0.163 | 4.2 | 0.042 |
| Time | -0.0906 | 0.0420 | 4.7 | 0.031 |
| Ageknot*Treatment | 0.109 | 0.0889 | 1.5 | 0.220 |
| Treatment*Time | -0.130 | 0.0662 | 3.9 | 0.049 |

Figure 5: GEE Model Summary

The GEE model fitted to the IPW data yielded the coefficients in the above table. At the level of alpha equals to 0.05, the main effects for Time and Treatment are significant, as is their interaction.

| Variable | Complete Case | All Available | IPW |
|---|---|---|---|
| Treatment | -0.713 | -0.286 | -0.331 |
|  | 1.87 | 0.124 | 0.163 |
| Time | -0.0856 | -0.102 | -0.0906 |
|  | 0.0465 | 0.0325 | 0.0420 |
| Treatment*Time | -0.265 | -0.111 | -0.130 |
|  | 0.225 | 0.0498 | 0.0662 |
| QIC | -12732 | -23716 | -23608 |

Figure 6: Comparison of Model Coefficients

Table II compares the significant coefficients of this model to those models fitted to the complete case and all available data. The standard error for each coefficient is listed below. Fitting the same GEE to the complete case dataset sees Treatment and Treatment×Time coefficient values are less consistent with those produced by the GEE from the all available and IPW data. As well, the standard errors for the main Treatment and Treatment×Time effects are larger when fit to the complete case data.

## 3.2 Rate Ratio

From Table I, we can calculate the expected BDI for both treatment groups at each time occasion. At the end of the active phase (two months), the expected BDI for the control group is 20.1 and for the treatment

group is 14.4. Both BDI means decrease by the four month mark, to 17.9 for the control group and 12.3 for the treatment.

The rate ratio of predicted BDI between the end and beginning of the active phase (baseline) for the control group is 0.91 (95% CI, 0.84, 0.99), and for the treatment group is 0.88 (95% CI, 0.77, 1.00). The rate ratio of predicted BDI between the treatment and control groups across the study is 0.72 (95% CI, 0.52, 0.99).

# Discussion

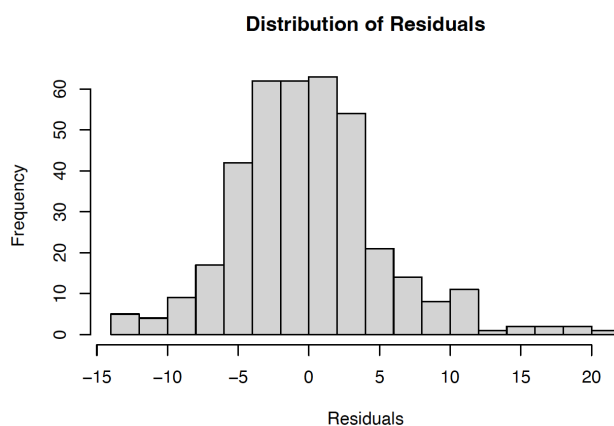## 4.1 Model Diagnostics for LMM

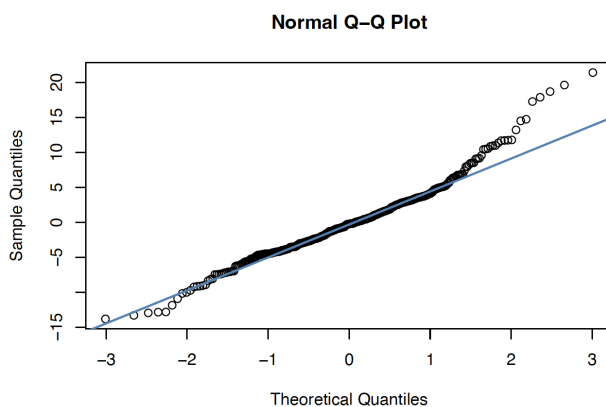### 4.1.1 Normality



Figure 7: Residual Distribution



Figure 8: Residual QQ Plot

By the distribution and QQ plot, residuals exhibit a good normality.

### 4.1.2 Homoscedasticity

The residuals do not appear to have a pattern across the range of fitted values. However, there seems to be a slight "funnel" shape, where the spread of residuals increases as the fitted values increase. This could
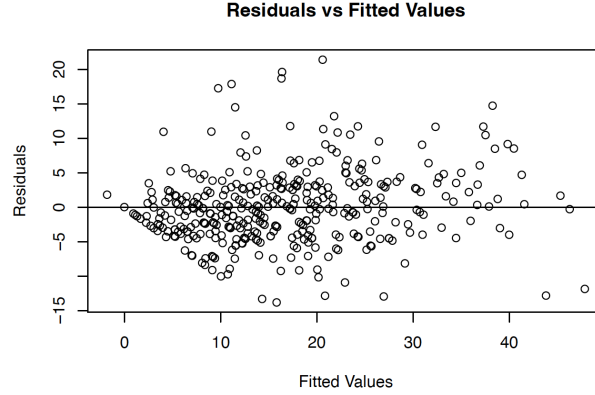
Figure 9: Plotting Residuals and Fitted Values

be an indication of heteroscedasticity. That is, the variance of the residuals might not be constant across the range of fitted values. It slightly violates the assumption of homoscedasticity.

### 4.1.3 Multicollinearity

```
##                   GVIF Df GVIF^(1/(2*Df))
## month       1.001649  4         1.000206
## treatment   1.001316  1         1.000658
## length      1.002573  1         1.001286
```

Figure 10: Multicollinearity Statistics Output

The VIF of each predictor is nearly to 1, indicating no correlation with the other variables, suggesting that there is no multicollinearity between these predictors; In other words, each predictor in the model is independent of the other predictors. All GVIF^ (1/(2*DF)) values are close to 1, indicating that there are no multicollinearity concerns.

## 4.2 Missing Data Handling

A more nuanced interpretation of the missing data's nature requires access to in-depth information regarding the study design, the protocols for participant engagement over time, and any complications encountered during the data collection phase. Such information would enable a more precise inference and inform the choice of strategy for handling the missing data in the analysis. Advanced statistical techniques may have to be employed to determine whether the missingness is systematic or not, and to select the most suitable method for addressing the missing data.

# Conclusion

## 5.1 Conclusion

Since the rate ratio between the treatment and control groups is 0.88, with a confidence interval between 0.77 and 1, we conclude that the Beat-the-Blues treatment sees a significant reduction in BDI. We can also see a significant reduction in BDI level throughout the active phase for both groups, indicating effective decrease in BDI for both the treatment as usual (control) and Beat-the-Blues active regimen. Applying the IPW method to the dataset in order to accommodate the missing data saw a smaller standard error for the coefficient estimates and a decrease in QIC from the complete case dataset. Using IPW and GEE to fit the data saw a significant reduction in BDI when undergoing Beat-the-Blues treatment, and for six months follow-up.

## 5.2 Limitations

Due to the missing data in the study, IPW had to be applied in order to re-weight the observations respective to their probability of remaining in the study. Future study re-design should consider alternative ways to ensure participants attend follow-up BDI measurements at the chosen two-month intervals. While this is difficult for participants with severe depression and anxiety, meaning those remaining in the study may have a lower BDI as they are more likely to return for BDI re-evaluation, honorariums or other incentives could be attempted to promote return to the study.

As well, other demographic information such as age, gender, or a more granular measurement of length of time with anxiety/depression could be collected to develop a more nuanced model around these features. Since "treatment as usual" can describe a broad range of treatment strategies (such as medication, counseling, weekly therapy, or a combination of these efforts), the control group encompasses a wide range of participants. This study would benefit from further distinction of "treatment as usual", potentially into multiple subgroups, to additionally compare whether frequency of regular treatment also impacts BDI in a similar manner to Beat-the-Blues.

# References

1. Driessen, E., & Hollon, S. D. (2010). Cognitive behavioral therapy for mood disorders: efficacy, moderators and mediators. The Psychiatric clinics of North America, 33(3), 537–555. https://doi.org/10.1016/j.psc.2010.04.005

2. Cuijpers, P., Miguel, C., Harrer, M., Plessen, C. Y., Ciharova, M., Ebert, D., & Karyotaki, E. (2023). Cognitive behavior therapy vs. control conditions, other psychotherapies, pharmacotherapies and combined treatment for depression: a comprehensive meta-analysis including 409 trials with 52,702 patients. World psychiatry : official journal of the World Psychiatric Association (WPA), 22(1), 105–115. https://doi.org/10.1002/wps.21069

3. David, D., Cristea, I., & Hofmann, S. G. (2018). Why Cognitive Behavioral Therapy Is the Current Gold Standard of Psychotherapy. Frontiers in psychiatry, 9, 4. https://doi.org/10.3389/fpsyt.2018.00004

4. Luo, C., Sanger, N., Singhal, N., Pattrick, K., Shams, I., Shahid, H., Hoang, P., Schmidt, J., Lee, J., Haber, S., Puckering, M., Buchanan, N., Lee, P., Ng, K., Sun, S., Kheyson, S., Chung, D. C., Sanger, S., Thabane, L., & Samaan, Z. (2020). A comparison of electronically-delivered and face to face cognitive behavioural therapies in depressive disorders: A systematic review and meta-analysis. EClinicalMedicine, 24, 100442. https://doi.org/10.1016/j.eclinm.2020.100442

5. Cavanagh, K., Herbeck Belnap, B., Rothenberger, S. D., Abebe, K. Z., & Rollman, B. L. (2017). My care manager, my computer therapy and me: The relationship triangle in computerized cognitive behavioural therapy. Internet interventions, 11, 11–19. https://doi.org/10.1016/j.invent.2017.10.005

6. Cavanagh, K., Shapiro, D. A., Van Den Berg, S., Swain, S., Barkham, M., & Proudfoot, J. (2009). The acceptability of computer-aided cognitive behavioural therapy: a pragmatic study. Cognitive behaviour therapy, 38(4), 235–246. https://doi.org/10.1080/16506070802561256

7. Jackson-Koku G. (2016). Beck Depression Inventory. Occupational medicine (Oxford, England), 66(2), 174–175. https://doi.org/10.1093/occmed/kqv087

8. Beck, A. T., Ward, C. H., Mendelson, M., Mock, J., & Erbaugh, J. (1961). An inventory for measuring depression. Archives of general psychiatry, 4, 561–571. https://doi.org/10.1001/archpsyc.1961.01710120031004

9. Gaynes, B. N., Asher, G., Gartlehner, G., Hoffman, V., Green, J., Boland, E., Lux, L., Weber, R. P., Randolph, C., Bann, C., Coker-Schwimmer, E., Viswanathan, M., & Lohr, K. N. (2018). Definition of Treatment-Resistant Depression in the Medicare Population. Agency for Healthcare Research and Quality (US).

10. Seaman, S. R., & White, I. R. (2013). Review of inverse probability weighting for dealing with missing data. Statistical Methods in Medical Research, 22(3), 278–295. doi:10.1177/0962280210395740

# Appendix

GitHub link: https://github.com/victoriachui/chl5222

```r
#Set-Up
library(here)
library(tidyverse)
library(tinytex)
library(knitr)
library(kableExtra)
library(gee)
library(geepack)
library(lubridate)
library(texreg)
library(MetBrewer)
library(lme4)
library(nlme)
library(broom)
library(ggplot2)
library(glmnet)
library(geepack)
library(naniar)
library(vcd)
library(GGally)
library(lmerTest)
library(car)
knitr::opts_chunk$set(eval = FALSE)
knitr::opts_knit$set(root.dir = here())
knitr::opts_chunk$set(warning = FALSE, message = FALSE)
theme_set(theme_bw(base_size = 15))
print("Set Up Complete")
```

```
## [1] "Set Up Complete"
```

---

```r
#Importing/cleaning data:
data <- read.table("/Users/yuchenjiang/Desktop/group_project/btheb.txt", header=TRUE)
data <- tibble::rowid_to_column(data, "ID")
print(dim(data))
head(data,5)
```

---

```r
# Exploratory Data Analysis

#Renaming columns:
names(data)[names(data) == 'bdi.pre'] <- 'bdi.0m' #referring to baseline as "0 month"
names(data)[names(data) == 'treatment'] <- 'Treatment'

#Pivoting data:
data_long <-
```

14

```r
  pivot_longer(data,  #wide data
    cols = starts_with("bdi"), #columns for BDI measures observations
    names_to = "Month", #new column name
    values_to = "BDI" #BDI measurements
  )

#Trajectory of BDI by treatment:
data_long %>%
  mutate(Treatment = as.factor(Treatment)) %>% #change to factor
  ggplot(aes(y = BDI, x = Month, group=Treatment,color = Treatment)) +
  stat_summary(aes(group = Treatment), fun = mean, geom="line") +
  labs(title = "Trajectory of BDI per Treatment", y="BDI",
       caption='Figure 1. Mean trajectory of BDI, grouped by Treatment') +
  theme(plot.caption = element_text(hjust=0))

#Baseline BDI distributions
ggplot(data, aes(x=bdi.0m, fill=Treatment)) +
  geom_boxplot()


## Missingness
# Examining monotone missingness pattern
library(mice)

data_miss <- data %>%
  dplyr::select(-ID,-Treatment) %>%
  md.pattern()
```

---

```r
#Complete Case
complete <- data[!is.na(data$bdi.8m),]

#Pivoting data:
complete_long <-
  pivot_longer(complete,  #wide data
    cols = starts_with("bdi"), #columns for BDI measures observations
    names_to = "Month", #new column name
    values_to = "BDI" #BDI measurements
  )

complete_model <- complete_long %>% #adding Time for month difference from 2 months
  mutate(Time=case_when(
    Month == 'bdi.0m' ~ -2,
    Month == 'bdi.2m' ~ 0,
    Month == 'bdi.4m' ~ 2,
    Month == 'bdi.6m' ~ 4,
    Month == 'bdi.8m' ~ 6,
    ))

#adding pre-/post-treatment
complete_model <- complete_model %>%
  mutate(knot = Time * I(Time >= 0)) #adding Time0 for before/after knot (2 months)
```

```r
complete_model <- complete_model %>%
  mutate(Treatment=as.factor(Treatment))

#setting control as ref
complete_model <- within(complete_model, Treatment <- relevel(Treatment, ref = 'TAU'))

#with knotxTreatment and TimexTreatment
GEE_model_complete <- geeglm(BDI ~ knot*Treatment+Time*Treatment,
                 id = ID,
                 family = poisson(link = "log"),
                 corstr = "unstructured",
                 scale.fix = TRUE,
                 data = complete_model)
summary(GEE_model_complete)

#no predictors are significant at alpha = 0.05

QIC(GEE_model_complete)["QIC"]


## All Available Data

#Removing Na --> all available data:
data_long_na = data_long[complete.cases(data_long), ]

all_available <- data_long_na %>% #adding Time for month difference from 2 months
  mutate(Time=case_when(
    Month == 'bdi.0m' ~ -2,
    Month == 'bdi.2m' ~ 0,
    Month == 'bdi.4m' ~ 2,
    Month == 'bdi.6m' ~ 4,
    Month == 'bdi.8m' ~ 6,
    ))

#adding pre-/post-treatment
all_available <- all_available %>%
  mutate(knot = Time * I(Time >= 0)) #adding Time0 for before/after knot (2 months)

all_available <- all_available %>%
  mutate(Treatment=as.factor(Treatment))

#setting control as ref
all_available <- within(all_available, Treatment <- relevel(Treatment, ref = 'TAU'))

#with knotxTreatment and TimexTreatment
GEE_all_available <- geeglm(BDI ~ knot*Treatment+Time*Treatment,
                 id = ID,
                 family = poisson(link = "log"),
                 corstr = "unstructured",
                 scale.fix = TRUE,
                 data = all_available)
summary(GEE_all_available)

#main effects for Treatment and Time are significant, as are interactions
```

```r
QIC(GEE_all_available)["QIC"] #smaller than complete case


## Inverse Probability Weighting

#Drop Out Model
ipwdat <- data_long %>%
  group_by(ID) %>%
  mutate(prevy = lag(BDI)) %>%
  ungroup() %>%
  mutate(r = ifelse(is.na(BDI), 0, 1),
         t2 = ifelse(Month == 'bdi.2m', 1, 0),
         t4 = ifelse(Month == 'bdi.4m', 1, 0),
         t6 = ifelse(Month == 'bdi.6m', 1, 0),
         t8 = ifelse(Month == 'bdi.8m', 1, 0),
         trt = ifelse(Treatment=='BtheB',1,0),
         trt.prevy = trt * prevy) %>%
  filter(!is.na(BDI)|!is.na(prevy))

#Fitting Model
r <- ipwdat$r
xmat <- as.matrix(
  cbind(rep(1, length(r)), ipwdat[,c("t2", "t4", "t6", "trt", "prevy", "trt.prevy")])
  )
rmod <- glm(r ~ xmat, family = poisson("log"))
summary(rmod)$coef

dropcoef <- summary(rmod)$coef[,1]


ipwdat <- ipwdat %>%
  mutate(logit = xmat %*% dropcoef,
         logitp = as.numeric(logit),
         phat = ifelse(Month == 'bdi.0m', 1, exp(logitp) / (1 + exp(logitp)))) %>%
  group_by(ID) %>%
  mutate(cumprob = cumprod(phat),
         ipw = 1/cumprob) %>%
  ungroup()

#Comparing Weights by Month
ipwdat %>%
  filter(r == 1) %>%
  ggplot(aes(y = ipw, x = as.factor(Month))) +
  geom_boxplot() +
  labs(y = "IPW", x = "Occasion")


#GEE modeling on this dataset
data_model <- ipwdat %>%
  mutate(Time=case_when(
    Month == 'bdi.0m' ~ -2,
    Month == 'bdi.2m' ~ 0,
    Month == 'bdi.4m' ~ 2,
    Month == 'bdi.6m' ~ 4,
    Month == 'bdi.8m' ~ 6,
```

```r
  ))

knot <- 0
#adding knot
data_model <- data_model %>%
  mutate(knot = (Time-knot)* I(Time >= knot))

data_model <- data_model %>%
  mutate(Treatment=as.factor(Treatment))

#setting control as ref
data_model <- within(data_model, Treatment <- relevel(Treatment, ref = 'TAU'))

# ipw-gee
ipwgee <- geeglm(BDI ~ knot*Treatment+Time*Treatment,
                 family = poisson("log"),
                 id = ID, scale.fix = TRUE,
                 corstr = "independence",
                 weights = ipw,
                 data = data_model)
summary(ipwgee)

QIC(ipwgee)['QIC'] #similar to all available


## Missing values by treatment

BtheB <- read.table("btheb.txt",
                    header = TRUE)
selected_columns <- BtheB[, c("bdi.2m", "bdi.4m", "bdi.6m", "bdi.8m", "treatment")]
# calculate the number of missing values for bdi by treatment group
missing_data_frame <- selected_columns %>%
  gather(key = "variable", value = "value", -treatment) %>%
  # Exclude treatment column from gathering
  group_by(variable, treatment) %>%
  summarize(num_missing = sum(is.na(value)))

# colored by treatment
ggplot(missing_data_frame, aes(x = variable, y = num_missing, fill = treatment)) +
  geom_bar(stat = "identity") +
  geom_text(aes(label = num_missing), position = position_stack(vjust = 0.5), size = 4) +
  theme_minimal() +
  labs(x = "Variable", y = "Number of Missing Values",
       title = "Missing Data in BtheB Dataset by Treatment") +
  theme(plot.title = element_text(hjust = 0.5),
        axis.text.x = element_text(hjust = 1)) +
  scale_fill_brewer(palette = "Set1")


## Missing values by drug

selected_columns <- BtheB[, c("bdi.2m", "bdi.4m", "bdi.6m", "bdi.8m", "drug")]
# calculate the number of missing values for bdi by drug group
missing_data_frame <- selected_columns %>%
  gather(key = "variable", value = "value", -drug) %>%
```

```r
  group_by(variable, drug) %>%
  summarize(num_missing = sum(is.na(value)))

# colored by drug
ggplot(missing_data_frame, aes(x = variable, y = num_missing, fill = drug)) +
  geom_bar(stat = "identity") +
  geom_text(aes(label = num_missing),
            position = position_stack(vjust = 0.5), size = 4) +
  theme_minimal() +
  labs(x = "Variable", y = "Number of Missing Values",
       title = "Missing Data in BtheB Dataset by Drug") +
  theme(plot.title = element_text(hjust = 0.5),
        axis.text.x = element_text(hjust = 1)) +
  scale_fill_brewer(palette = "Set1")


# Check BDI Distribution
# Create a histogram of baseline BDI levels
baseline <- BtheB[["bdi.pre"]]
hist(baseline, main = "Distribution of Baseline BDI level", xlab = "Baseline BDI level")

bdi_2m <- BtheB[["bdi.2m"]]
hist(bdi_2m, main = "Distribution of BDI level 2M", xlab ="2M BDI level")

bdi_4m <- BtheB[["bdi.4m"]]
hist(bdi_4m, main = "Distribution of BDI level 4M", xlab = "4M BDI level")

bdi_6m <- BtheB[["bdi.6m"]]
hist(bdi_6m, main = "Distribution of BDI level 6M", xlab = "6M BDI level")

bdi_8m <- BtheB[["bdi.8m"]]
hist(bdi_8m, main = "Distribution of BDI level 8M", xlab = "8M BDI level")


# Check for Independence
# Correlation Plot
ggpairs(BtheB)


# Data Processing

# convert to long format
BtheB_long <- BtheB %>%
  mutate(id = row_number()) %>%
  pivot_longer(cols = starts_with("bdi"),
               names_to = "month",
               values_to = "bdi") %>%
  mutate(month = factor(month,
                        levels = c("bdi.pre", "bdi.2m", "bdi.4m", "bdi.6m", "bdi.8m"),
                        labels = c("Baseline", "2m", "4m", "6m", "8m"))) %>%
  mutate(treatment = as.factor(ifelse(treatment == "BtheB", 1, 0)),
         length = as.factor(ifelse(length == ">6m", 1, 0)))

# Linear Mixed Effect Model
model <- lmer(bdi ~ month + treatment + length + (1 | id), data = BtheB_long)
# summary(model)
```

```r
# test the overall significance of the fixed effects
# anova(model)

# check residuals distribution
residuals <- residuals(model)
hist(residuals, breaks = 'Scott', main = "Distribution of Residuals", xlab = "Residuals")

qqnorm(residuals)
qqline(residuals, col = "steelblue", lwd = 2)

# Shapiro-Wilk test for normality
shapiro.test(residuals)

# check for Homoscedasticity
plot(fitted(model), residuals,
     xlab = "Fitted Values",
     ylab = "Residuals",
     main = "Residuals vs Fitted Values")
abline(h = 0, col = "black", lwd = 1)

# Calculate VIF
vif_result <- vif(model)
print(vif_result)
```