

Bayesian Analysis of Randomized Controlled Trials

Julian Bautista[†], Alex Pavlakis[†], Advait Rajagopal[†]

April 29, 2018

Abstract

Objective :

Method :

Results :

Discussion :

Keywords – binge-eating disorder, bulimia nervosa, guided self-help, smartphone, Bayesian inference, randomized control trial, RCT, multilevel, hierarchical Poisson models

[†]The New School for Social Research, 6 E 16 St, New York, NY, 10003.
Correspondence : rajaa598@newschool.edu

1 Introduction

mention (Grotzinger, Hildebrandt, & Yu, 2015) here

Bayesian methods are gaining popularity in medical, pharmaceutical, and social-science research. Bayesian data analysis combines mathematical assumptions with data collected, to produce a robust estimate of a data generating process. The purpose of this paper is to motivate the use of Bayesian data analysis in applied research, especially in analyzing results of Randomized Control Trials (RCTs). While regression models are well-suited for analyzing RCTs because they allow researchers to estimate causal effects more precisely by controlling for all available pre-treatment covariates [Gelman and Hill, 2007]; Bayesian regression methods have the added benefit of allowing the inclusion of prior information that can increase predictive accuracy, and estimating varying treatment effects among subgroups of the population. These improvements in the analysis of treatment effects make Bayesian methods especially valuable for clinical researchers in the field of psychiatry.

The analysis of RCTs is also part of standard practice for assessing efficacy of treatments for binge eating and/or purging for binge eating disorder (BED), bulimia nervosa (BN), or subthreshold variants. In this paper, we use the data from recently published RCT examining the efficacy of cognitive behavioral therapy version of guided self-help (CBT-GSH) with the addition of a smart-phone app designed for the intervention [Hildebrandt et al., 2017]. We aim for this paper to be an introduction to the method with a focus on providing guidelines for executing this method in a rigorous way and with enough information to facilitate replication of these methods, especially in the context of treatment of eating disorders.

The rest of the paper is organized as follows, section 2 explains the process and steps of Bayesian data analysis and inference. Section 3 discusses the advantages of Bayesian analysis, section 4 explains the experiment, modeling and results. Section 5 concludes and the Appendix contains relevant code for easy reproduction of results.

2 Methods

Bayesian Data Analysis

Bayesian data analysis (BDA) is a method for building statistical models to describe data. Researchers begin with explicit assumptions about the data generating process based on past research and the scientific nature of the problem. Based on these beliefs and assumptions, they then collect data from designed experiments or observational studies. Using Bayes' rule they combine their assumptions or *prior information* with the actual data into a comprehensive probability model. This model contains information about all known (observed data) and unknown (unobserved parameters) quantities related to the data generating process. There are three main steps to BDA which are listed below and explained in detail in the subsequent sections:

1. Specify a model based on scientific knowledge of the data generating process.
2. Estimate model parameters based on observed data.
3. Evaluate the model's accuracy and expand or alter the model.

2.1 Probability notation and Bayes' rule

In Bayesian Data Analysis, assumptions are specified mathematically as *prior distributions*. Data is represented through a *likelihood model*. Bayes' Rule combines prior distribution and data likelihood into a *posterior distribution*. A formal expression of the Bayes' Rule is as follows;

$$p(\theta|y, x) = \frac{p(\theta)p(y|\theta, x)}{p(y)} \quad (1)$$

θ is the estimated parameter. In the context of an RCT, θ is the effect of the treatment x on dependent variable y . $p(\theta)$ is the prior distribution of the treatment effect, which captures the researcher's beliefs about the model parameter prior to any analysis. $p(y|\theta, x)$ is the likelihood function and is the probability of the observed data given the parameter. $p(y)$ is a normalizing constant with respect to θ that ensures the left hand side $p(\theta|y, x)$ is a proper probability distribution that integrates to 1. For a full treatment of proper distributions and normalizing constants see Gelman et al. (2014) .

For our purpose, we can ignore this denominator and rewrite the expression for the posterior $p(\theta|y)$ as;

$$p(\theta|y, x) \propto p(\theta)p(y|\theta, x) \quad (2)$$

The posterior distribution is proportional to the product of the prior and the likelihood. Ultimately the goal of modeling is to learn the posterior distribution $p(\theta|y, x)$ and summarize it accurately (Gelman et al., 2014). Based on this we can make inferences and predictions.

2.2 Model development

Model development involves specifying a model that accounts for all observed data and unobserved parameters. The model should include all knowledge of the experiment or data collection process and should be logically consistent with scientific nature of the problem. We approach model development in three steps.

2.2.1 Exploratory data analysis

Summarizing data through visualizations and summary statistics is a critical model building step. We look at summary statistics of our relevant variables to understand their distributions. We explore data through univariate analyses, such as histograms to visualize distributions, and bivariate analyses, such as scatter plots to see if there are meaningful linear or non linear relationships. This exploratory process informs us of the distribution of individual variables, their correlations, and other relevant information for modeling.

Exploratory analysis helps us to choose an appropriate likelihood model to describe the data. For example, if our variables are normally distributed and linearly correlated, we might choose a linear regression model. If the relationship between variables is not straightforward we might choose more complicated models. It allows researchers to explore and solidify intuition about the problem at hand and the nature of the data. An introductory treatment of exploratory data analysis can be found in Tukey (1977). A more advanced explanation of exploratory data analysis in the context of Bayesian statistics is available in (Gabry et al., 2017) and Betancourt (2018).

2.2.2 Setting up a likelihood model

The likelihood function is often analogous to a traditional regression equation. The researcher must select independent variables that represent important determinants of the outcome variable. The likelihood represents the distribution of the outcome variable given the independent variables and model parameters. To choose the right likelihood function, it is critical to know the type of data. Data can be binary, categorical, ordinal, count, or continuous and each of these types of data requires a different kind of model. In the example in this paper, the outcome variable is

non-negative and discrete so we choose a Poisson likelihood function (see section 4.2).

2.2.3 Choosing a prior distribution

The prior distribution is a mathematical encoding of researchers' assumptions. A prior distribution serves three functions. First, it makes assumptions about the underlying scientific nature of the problem explicit. Second, it regularizes or constrains the parameter space by specifying likely ranges for parameter values. Third, it facilitates the calculation of a posterior distribution and makes it possible to generate simulations from that distribution.

Prior choice depends on the parameter or coefficient of interest. We could assign a completely "noninformative" or flat prior to our coefficient by specifying a uniform distribution as the prior. This is equivalent to saying that our parameter is equally likely to assume any value from negative infinity to positive infinity and we have no more information about it. Using a noninformative uniform prior is the same as carrying out a maximum likelihood estimation of the parameter of interest.

Researchers rarely know nothing about the problem or relevant parameters. They often possess valuable information about the parameter which statisticians can incorporate as informative prior information. For instance, if a coefficient of interest is a proportion, then it *must* be between 0 and 1, and we can assign a *Beta* prior to that coefficient. If we believe that a coefficient is close to zero, but may be positive or negative, we could assign an informative *Normal*(0, 1) prior distribution to it. More information about prior choice for the same can be found on the GitHub page for Stan developers (Stan Development Team, 2015). We explain the Bayesian probabilistic programming language, Stan, in section 2.3.

2.3 Model estimation

Once the likelihood model and prior distribution are specified, the posterior distribution of an outcome variable can be estimated. As discussed in section 2.1, the posterior distribution is obtained by multiplying the prior and likelihood. The distribution obtained by this process is proportional to the true posterior because we can ignore the normalizing constant (the denominator in equation 1). We use an approximation of the posterior as shown in equation 2 because calculating the true posterior analytically may be practically impossible. A standard practice is to use Markov Chain Monte Carlo (MCMC) sampling methods to approximate the posterior up to a normalizing

constant and sample from it (Gelman et al., 2014). There are other approaches to calculating the posterior distribution, but those are beyond the scope of this paper. Analyses in this paper are carried out with the Bayesian probabilistic modeling language Stan. The R interface of the language can be understood at Stan Development Team (2018) and Carpenter et al. (2017) present a clear conceptual overview of the language. Stan uses a Hamiltonian Monte Carlo sampling algorithm (from a broader class of MCMC sampling methods) to approximate the posterior distribution. Betancourt (2017) has a clear exposition of how the algorithm works. Stan returns the full posterior distribution of the desired parameters. Stan can also predict values based on the specified model which can be used for model checking, validation, and expansion.

2.4 Model checking and expansion

We evaluate whether our model explains the data by investigating parameter distributions and posterior predictive checking. Visualizations of parameter distributions, such as histograms, enable us to summarize estimates and the uncertainty around them. Posterior predictive checks involve “simulating replicated data under the fitted model and then comparing these to the observed data” (Gelman & Hill, 2006, p. 158). While we may not expect our model to generate our data exactly, it should recover important patterns. Once we have estimated a model and studied its properties, we can *expand* it through reparameterizations, adding parameters, or changing prior distributions. After we have fit multiple models, we may want to compare their performances. There is wide literature on the most effective ways to compare and assess Bayesian models. These include cross validation methods, Bayes factors, and information criterion. One popular Bayesian approach to model checking, is *leave-one-out-cross validation (loo-cv)*. The process compares models by fitting them on all data points except one, then evaluating how they predict the remaining data point. This process is repeated until all data points have been left out once. There are also several information criterion that can be used. An exhaustive summary of these “practical” model checking methods can be found in Vehtari et al. (2017). Some model checking methods that are closer to the null-hypothesis and significance testing framework are the ROPE (region of practical equivalence)(Kruschke, 2014) and Bayes factor (Rouder et al., 2009) approaches.

Model checking is a critical step in Bayesian Data Analysis; however, there is no one-size-fits-all approach. We recommend posterior predictive checks because it is a direct way to assess the model fit to various aspects of data. By using posterior predictive checks we neither “accept” nor “reject” models but aim to understand their limitations in realistic replications (Gelman et al., 2014) .

3 Advantages of BDA for RCT

In this section we discuss some of the major benefits of using a Bayesian approach to the analysis of Randomized Control Trials.

3.1 Heterogenous Treatment Effects

Among the advantages of Bayesian modeling is its ability to capture heterogeneous treatment effects. For example in our study, we control for variation in demographic variables and variation in the treatment across the time. This approach is commonly known as a varying slope - varying intercept model (Gelman & Hill, 2006). Our approach of hierarchical modeling and partial pooling allows us to accomplish this naturally.

Hierarchical modeling has two closely related meanings (Feller & Gelman, 2015). Hierarchies can explain a hierarchical data structure like spatial or temporal variation. In the RCT we consider in this paper, each treatment period is treated as a level and so we obtain different treatment effects for each period. Hierarchies also describe how parameters are modeled. We model our treatment effects for each level or category in the data, to come from a common underlying prior distribution (Gelman et al., 2014). This is the idea of *partially pooled* estimates of treatment effects. Partial pooling in a Bayesian context is different than traditional complete pooling or no pooling approaches.

In non-Bayesian methods, researchers have two options: complete pooling and no pooling. With complete pooling, the categories in data are completely interchangeable, thus ignoring the uniqueness of categories. With no pooling, each category is treated as independent from the others, ignoring the interrelatedness of each of the categories. For example, if an RCT was conducted in multiple locations by different staff, the researcher must assume that the treatments are entirely identical in a complete pooling model, ignoring the differences in execution that may have taken place by differing staff. In the no pooling model, the RCT would assume that each treatment would be entirely different across locations, despite having close similarities in how the treatments were applied. Setting up a prior distribution allows researchers to solve this problem by saying that the treatment effects across the locations have a common mean, but vary based on location. Thus partial pooling is often a better representation of data as it takes into account both the uniqueness and interrelatedness of categories within a hierarchical model.

3.2 Making uncertainty explicit

Frequentist (Classical or non Bayesian) inference techniques like Null Hypothesis Significance Testing (NHST) and Confidence Interval (CI) estimation are often misinterpreted in different ways. Hoekstra et al. (2014) have a very clear exposition of the common ways in which NHST and CI's are often misused or misinterpreted. The reason we recommend summarizing uncertainty in parameter estimates with posterior intervals is because of the heavy dependence of confidence intervals on hypothetical replications and asymptotics (Gelman, 2013). In the frequentist paradigm, the source of uncertainty in the parameter estimate is from sampling; as the number of hypothetical replications tends towards infinity and with a completely non informative prior, 95% (typically used in the social sciences) of the estimates will fall within the 95% CI. Interpreting the width of a CI as a measure of the precision of the estimate is not correct (Morey et al., 2016). In a CI approach, error around the parameter estimate represents sampling error associated with a test statistic and not genuine uncertainty in the model.

Bayesians on the other hand, start out with certain assumptions about the estimated parameters and the uncertainty that they inherently contain. These assumptions are encapsulated in the prior, then iteratively updated using the new information contained with each observation of the data. The assumptions can include common sense knowledge such as understanding that estimates on ratios will be between 0 and 1, or it could be an expected estimate based on previous studies done. These models are more useful because they don't constrain the estimates to move asymptotically towards a pre-defined distribution of uncertainty. Instead, it allows the assumptions and data to meet together to form a distribution that captures the uniqueness of the specific scientific problem. In the Bayesian approach, we focus on the posterior distribution of parameters, eliminating the need for hypothetical replications or the issue of multiple comparisons (Gelman et al., 2014).

In the Frequentist (Classical or non-Bayesian) paradigm, the source of uncertainty is from sampling; as the number of replications tends towards infinity, 95% of estimates will fall within the 95% interval. Frequentist results rely heavily on asymptotics. Error in the models represent sampling error associated with a test statistic, not genuine uncertainty about the data generating process. This is problematic because modeling events as complex as human reaction always contains non-sampling error as well. However, due to this reliance on asymptotics, researchers using Frequentist methods are forced to assume that any non-sampling error disappears in the case of arbitrarily large replications. On the other hand, Bayesians start out

with certain assumptions about the estimated parameters and the uncertainty that they inherently contain. These assumptions are encapsulated in the prior, then iteratively updated using the new information contained with each observation of the data. This type of model is more useful because it does not constrain the estimates to move asymptotically towards a specific, pre-defined distribution of uncertainty. Instead, it allows the assumptions and data to meet together to form a distribution that captures the uniqueness of the specific scientific problem. This is often a more accurate representation of uncertainty.

Furthermore, models that use priors have access to information beyond what was collected within the study. Data does not live only on the spreadsheet. It is contained in the researcher’s knowledge of the literature, the data generating process, and the specific scientific problem. This can include common sense knowledge such as the understanding that estimates on ratios will be between 0 and 1, or it could be an expected estimate based on previous studies done. Regardless of the source of the information, knowledge about an estimator can be translated into the prior to improve accuracy.

By embracing the uncertainty inherent in statistical estimators, the Bayesian approach leaves little need for the p-values that are often used as thresholds in classical statistics. Intuitively, p-values represent the probability of the data conditional on the null hypothesis, under hypothetical replications. However, since researchers have many opportunities to choose how data is coded, which data is included, and what comparisons are made, it is possible to obtain “statistically significant” p-values even if there are no “true” effects (Gelman, 2013) In the Bayesian approach, we focus on the posterior distribution of parameters, leaving no need for worrying about hypothetical replications or the issue of multiple comparisons (Gelman et al., 2014).

Results

4 Impact of Smartphone App on Eating Behavior

Hildebrandt et al. (2017) conducted an experiment to test whether the Noom Monitor, a smartphone application, could augment the effect of in-person therapeutic treatment on binge eating behavior. The treatment, known as *guided self-help treatments based on cognitive-behavior therapy*

(CBT-GSH), had been shown in previous research to reduce binge eating behavior by 10-50%. The Noom Monitor application was designed to facilitate CBT-GSH. For this example, we consider two research questions from the experiment:

1. Is CBT-GSH more effective at reducing binge eating behavior when facilitated by the Noom Monitor?
2. Does the effect of the Noom Monitor vary over time?

4.1 Experimental design

66 men and women with Bulimia Nervosa (BN) or Binge Eating Disorder (BED) were randomly assigned into two treatment conditions: CBT-GSH (N= 33) or CBT-GSH + Noom (N=33). Therapy lasted for 12 weeks. Assessments were conducted at weeks 0, 4, 8, 12, 24, and 36. The primary outcome was Objective Bulimic Episodes (OBE). For more information about the experiment, choice of dependent (outcome) variable, historical context and past research in this area, see Hildebrandt et al. (2017). There is a discussion of Bulimia Nervosa and Binge Eating as well as an explanation of the choice of outcome variable OBE.

4.2 Model development

Exploratory data analysis

We start the analysis by plotting the outcome variable OBE. Figure 1 displays OBEs per week for each individual in both treatment conditions. A few aspects of the data immediately stand out, which suggest that any model should account for individual-level effects and time-level effects, and should let treatment effects vary over time.

- The number of OBEs decreases over the course of the treatment for almost all subjects.
- The biggest decreases in OBEs appear to occur in the early stages of treatment.
- The primary sources of variation in OBE appear to be *between people* and *over time*.

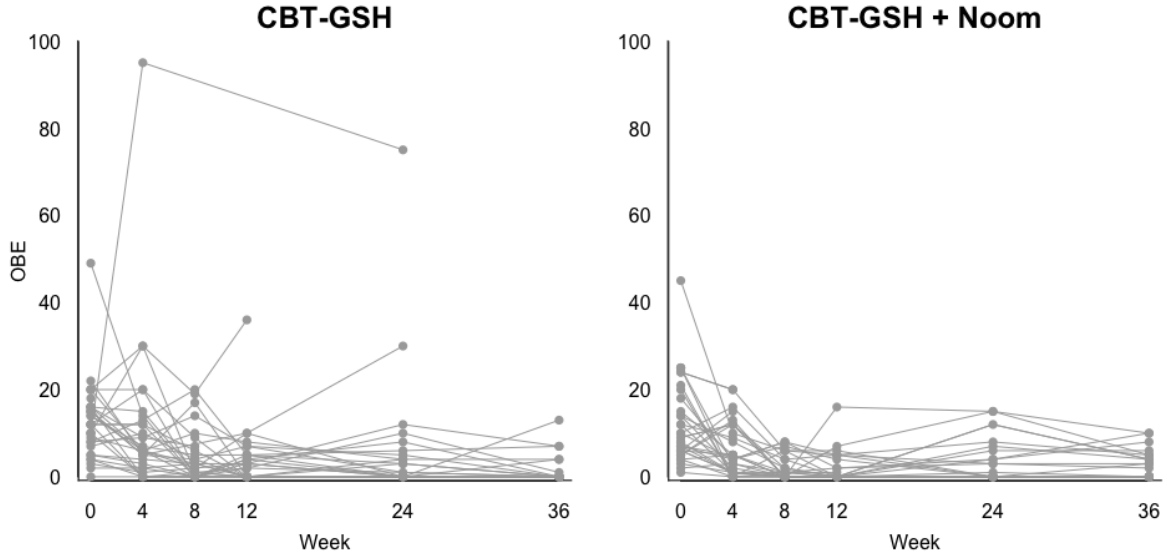


Figure 1: *Fig.1 (a) and (b) show the OBE measurements for individuals in the CBT-GSH group and CBT-GSH + Noom group respectively. The horizontal axis shows the week and the vertical axis shows the instances of OBE. The gray dots represent OBE readings over time for each individual.*

Figure 2 displays the distribution of OBEs in each condition in each week, aggregating across individuals. We notice three characteristics of the data from these histograms.

1. The distributions appear to condense around zero for both conditions over time
2. The distributions in the CBT-GSH condition appear to have longer tails than those in the CBT-GSH+Noom condition
3. OBEs are count data; they must be nonnegative integers.

These three characteristics suggest that the appropriate model for OBEs is the Poisson distribution, because it is restricted to nonnegative integers and can concentrate its density around low numbers with a long tail.

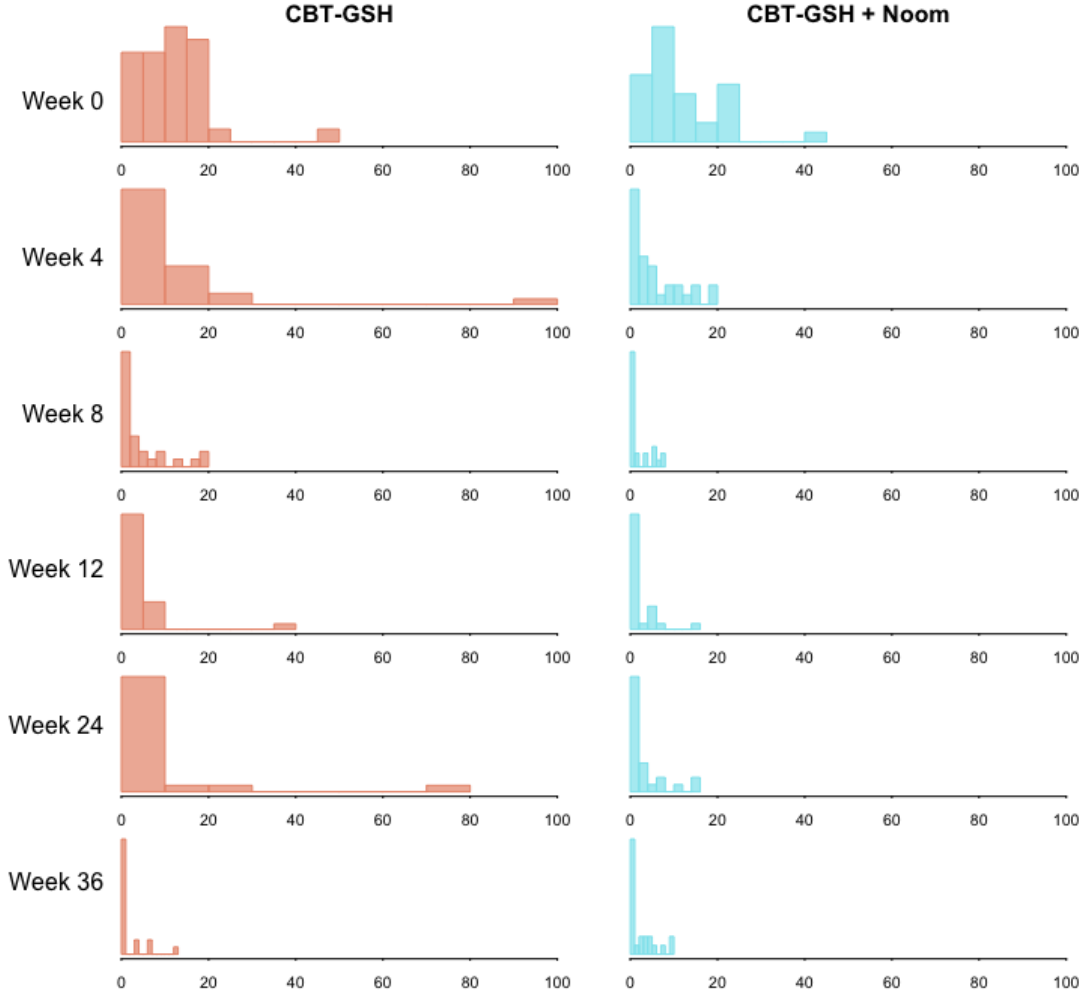


Figure 2: *Histograms display the distribution of OBEs in each condition in each stage of the treatment. The orange and light blue histograms show the distribution of OBEs for the CBT-GSH and CBT-GSH + Noom group respectively.*

Setting up a likelihood

We analyze RCTs by modeling the outcome of interest (in this case OBE) as a function of the treatment and all available pre-treatment covariates. The coefficients associated with the treatment are estimates of average treatment effects. Inclusion of all available pre-treatment covariates accounts for variation in the outcome variable, decreasing uncertainty around treatment effects and providing the model with more predictive power. We conduct *intent-to-treat* analysis, meaning that our inferences will be based on initial treatment assignment, and will not account for mid-

experiment dropouts.

The outcome variable is restricted to be nonnegative integers, so we fit a Poisson regression model, partially pooling across individuals, time periods, and treatment effects. For each individual in each time period, the number of OBEs follows a Poisson distribution, with a mean dependent on the characteristics of the individual and the time period.

$$OBE_{i,t} \sim \text{Poisson}(\lambda_{i,t}) \quad (3)$$

$$\lambda_{i,t} = \exp(\alpha_i + \beta_t + \gamma_t T_i + X_i \theta) \quad (4)$$

$$T_i = \begin{cases} 0, & \text{if } CBT - GSH \\ 1, & \text{if } CBT - GSH + Noom \end{cases} \quad (5)$$

α is an individual-specific intercept, β is a time-specific intercept, γ is a time-specific treatment effect, T is a treatment indicator, X is a matrix of individual level covariates (age, sex, race, etc), and θ is a vector of effects. Subscripts $i = 1, \dots, 66$ indicate individuals and subscripts $t = 0, 4, 8, 12, 24, 36$ indicate time periods.

Choosing a prior distribution

Table 1 has a list of different sources from which prior information has been obtained for this experiment. It aims to summarize the various methods which a researcher can use to incorporate prior information into the modeling process. There is a rich literature on binge-eating disorders and bulimia nervosa studies, implying a large amount of prior information. This makes analysis of similar RCTs amenable to Bayesian methods. For more examples of priors see the Appendix or check the “Prior Choice Recommendations” GitHub page (Stan Development Team, 2015).

Source of Prior Information	
Experimental Design	Outcome variable is nonnegative integers
Literature	Treatment effect size is small
	A large number of zeros in OBE data due to remission
Exploratory Data Analysis	There is variation in OBEs at the individual level
	There is variation in OBEs over time
	Treatment effects may vary over time

Table 1: *Sources of prior information.*

We believe that individual-level intercepts are simultaneously unique to the individual and common to the population; that is, each individual has their own baseline predilection to engage in eating disorder behavior, but those baseline predilections are not drastically different from each other. We operationalize this concept by modeling all individual-level intercepts as coming from a common distribution, with *hyperparameters* μ_α and τ_α . In Bayesian statistics, hyperparameters are parameters of prior distributions. In hierarchical models, we model hyperparameters explicitly.

$$\alpha_i \sim \text{Normal}(\mu_\alpha, \tau_\alpha) \quad \forall i \in 1, \dots, 66 \quad (6)$$

Similarly, we believe that time-specific treatment effects may be unique to each period but similar over time. We operationalize this concept by modeling all time-specific treatment effects γ as coming from a common distribution, with *hyperparameters* μ_γ and τ_γ .

$$\gamma_t \sim \text{Normal}(\mu_\gamma, \tau_\gamma) \quad \forall t \in 0, 4, 8, 12, 24, 36 \quad (7)$$

μ_γ is the *grand mean*, the overall treatment effect; τ_γ is the variation in treatment effects over time; and each individual γ_t is a time-period specific treatment effect. This approach has a natural smoothing effect: any extreme estimates of γ_t will be partially-pooled back toward the grand mean μ_γ .

We assign the following prior and hyperprior distributions:

$$\mu_\alpha \sim \text{Normal}(5, 5) \quad (8)$$

$$\tau_\alpha \sim \text{Cauchy}^+(0, 30) \quad (9)$$

$$\mu_\gamma \sim \text{Normal}(0, 5) \quad (10)$$

$$\tau_\gamma \sim \text{Cauchy}^+(0, 30) \quad (11)$$

$$\theta \sim \text{Normal}(0, 1) \quad (12)$$

The normal distributions around the individual and treatment effects allow us to guide the model to the appropriate range of parameter values, but with wide enough variance (5 in each case) to let the model find its own way in that range. Half cauchy priors on the variance parameters are weakly informative, with much of their mass around zero but gentle slopes in their tails, which have been shown to be effective prior distributions for variance parameters (Gelman, 2006).

4.3 Model estimation and results

We estimate this model with *Hamiltonian Monte-Carlo* in Stan. Model code is appended to this document. This is a particular algorithm from a larger class of Markov Chain Monte Carlo algo-

rithms, for more examples see Gelman et al. (2014).

Model results are displayed in Table 2. The table displays the mean of each of the parameter distributions along with the 50% posterior interval. Results suggest that using the Noom Monitor smartphone application during CBT-GSH may slightly decrease OBEs. There is evidence that the treatment effect varies over time, with the Noom effect being slightly more pronounced during stages 4, 8, 12 and 24 of therapy but decreasing by week 36.

	mean	25%	50%	75%
γ_0	0.18	-0.45	0.15	0.78
γ_4	-0.43	-1.05	-0.46	0.16
γ_8	-0.70	-1.33	-0.71	-0.10
γ_{12}	-0.65	-1.28	-0.68	-0.04
γ_{24}	-0.72	-1.34	-0.75	-0.11
γ_{36}	0.21	-0.42	0.19	0.82
μ_γ	-0.34	-0.98	-0.36	0.26
τ_γ	0.64	0.43	0.56	0.77

Table 2: *Table displays model results for Noom effects in all six time periods and grand mean and variance parameters.*

Discussion

4.4 Model checking, comparison, and expansion

Before using our model to make inferences about time-specific treatment effects, we check its fit by comparing model-simulated OBE to data OBE. If model simulations do not track the data well, we may want to revisit our model’s assumptions before trusting its inferences. If the model’s simulations recover patterns in the data, we are more inclined to trust it.

Figure 3 displays OBEs in each period for each individual in each treatment group, from raw data (upper plots) and model simulations (lower plots). Black lines display means for each period. This suggests that the model is able to pick up on the key variables that determine OBE over time for the duration of this experiment.

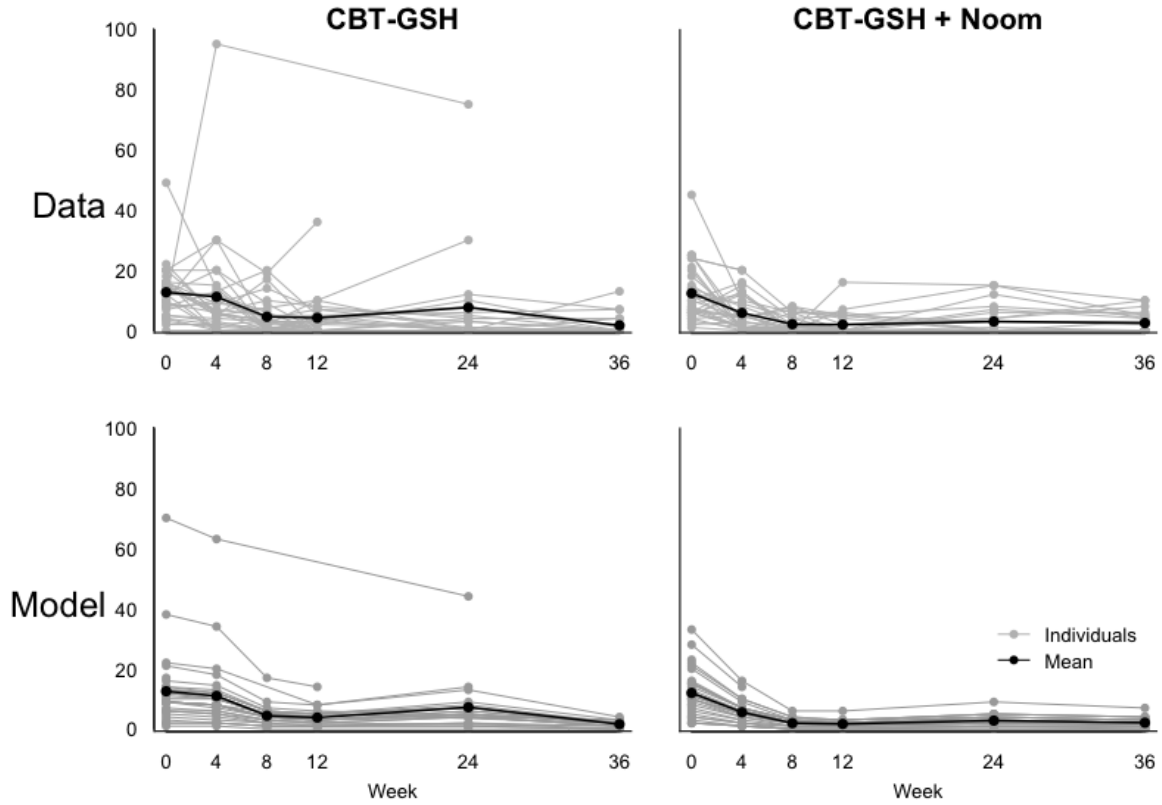


Figure 3: Fig.3 (a) and (b) in the upper row show OBEs in each period for each individual in CBT - GSH and CBT-GSH + Noom respectively. Black dots represent mean estimate for each period. Fig.3 (c) and (d) in the lower row show the simulated OBEs for each individual in CBT - GSH and CBT-GSH + Noom respectively. Black dots represent means from the simulated data for each period. The horizontal axis shows the week and the vertical axis shows the instances of OBE.

Another way to check the fit of the model is by comparing simulated data directly against the raw data. This is called posterior predictive checking and is our preferred form of model checking. Figure 4 shows this for both treatment conditions. Simulated data for the Noom condition appears to better track the raw data than simulated data for the no Noom condition. This is unsurprising, since the no Noom condition tended to have more outliers, which we would not expect (or want) our model to pick up perfectly from such a small sample.

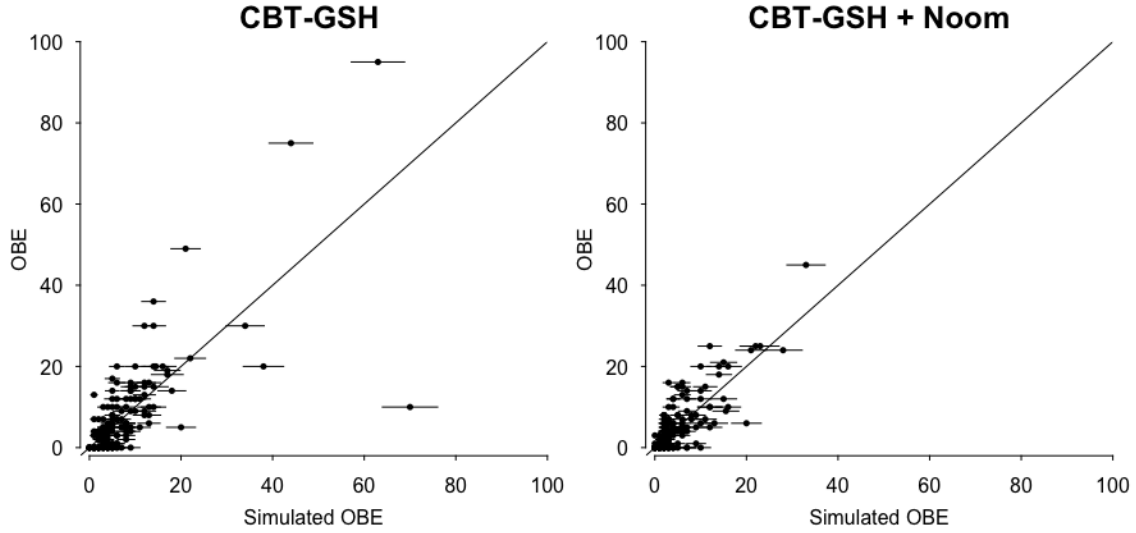


Figure 4: *Fig. 4(a) and (b) show the predicted OBE vs. the actual OBE data for CBT - GSH and CBT- GSH + Noom respectively. The horizontal axis shows the simulated OBE and the vertical axis shows the actual OBE. The black dots are the points that represent this and the lines around them show 50% intervals around the predictions. The upward sloping line is the 45 degree line.*

We compare the distributions of OBE for each condition in each time period by plotting density curves over the histograms in figure 2, displayed in figure 5. Figure 5 shows that our model is able to broadly pick up on the patterns in the data over time and between treatment conditions. We see that the density curves clearly peak at lower values close to zero and have long tails, correctly capturing the pattern in the OBE data.

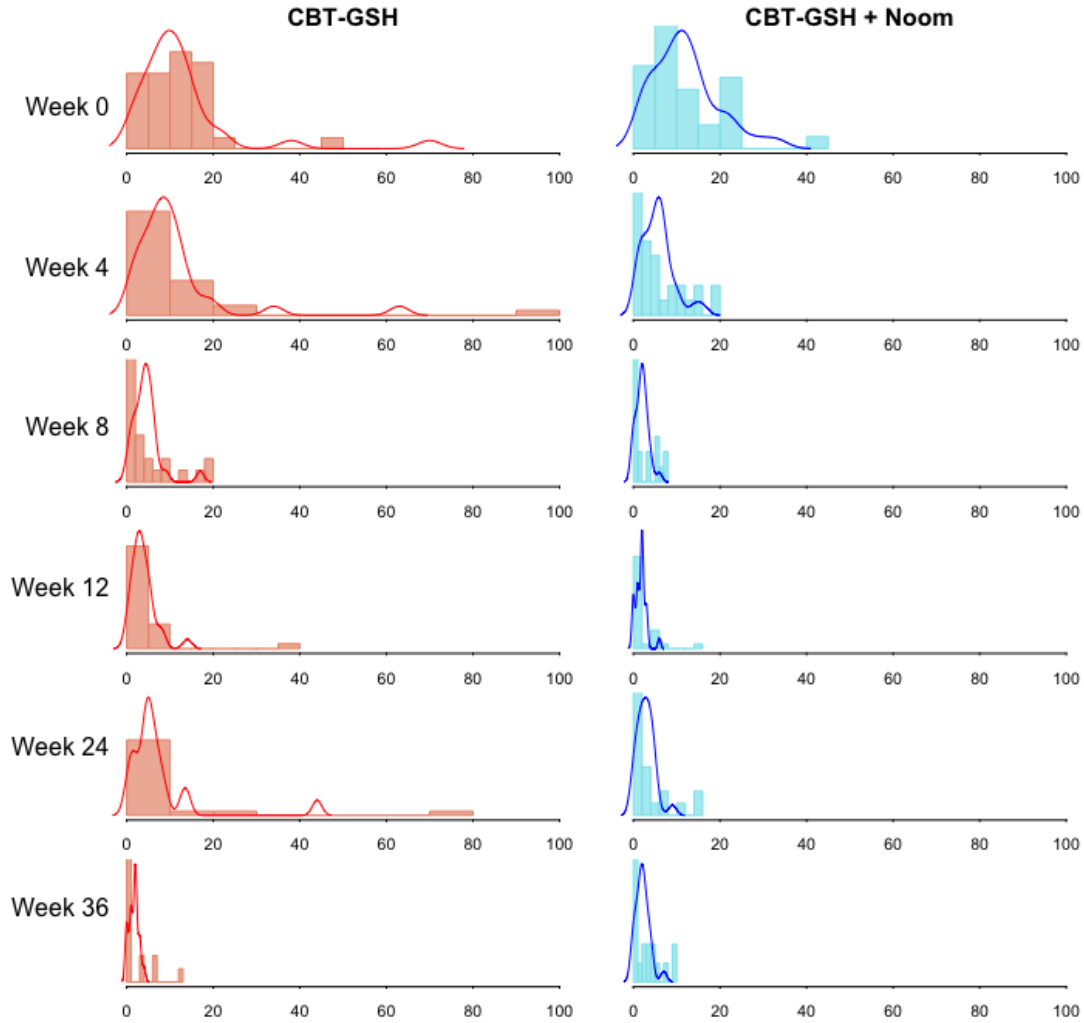


Figure 5: Histograms display the distribution of OBEs in each condition in each stage of the treatment. The orange and light blue histograms show the distribution of OBEs for the CBT-GSH and CBT-GSH + Noom group respectively. The red and blue lines show the predicted density curves obtained from the simulations.

Figure 6 displays the simulated OBE for both treatment groups (upper plot) and smoothed treatment effects (lower plot). In each measurement period, simulated OBE are higher for the CBT - GSH condition than for the CBT - GSH + Noom condition, with some of the difference likely attributable to use of the Noom Monitor smartphone app. This shows that the app has an effect on lowering episodes of binge eating.

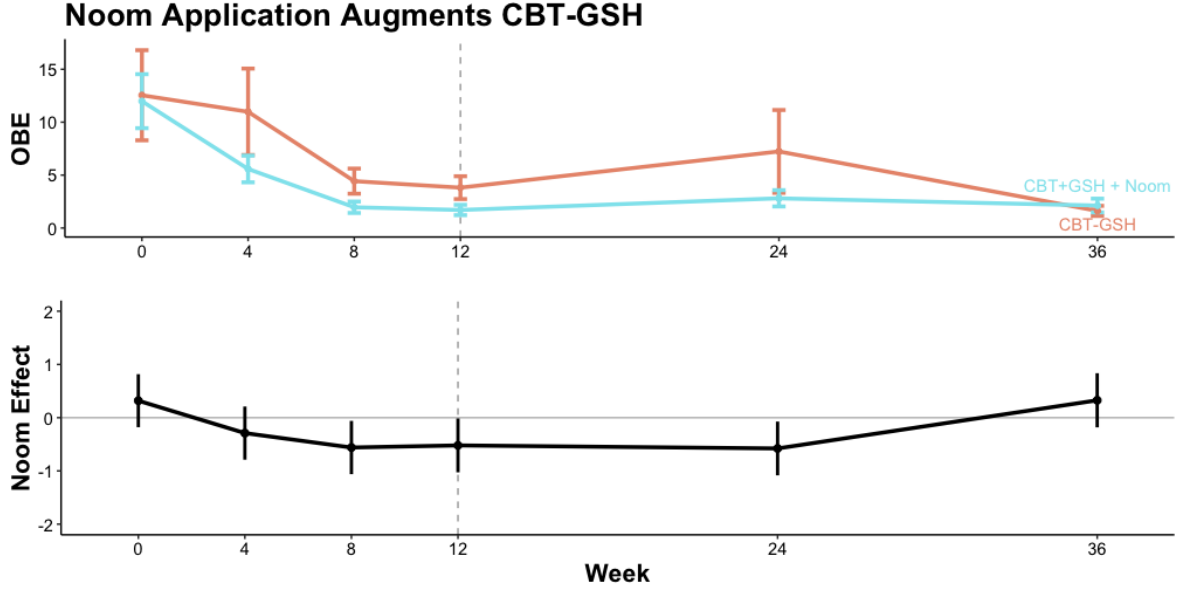


Figure 6: Fig. 6(a) the upper plot shows the simulated OBE in each time period. The CBT - GSH condition is in orange and the CBT - GSH + Noom condition is in orange. The bars around each point show the 95% interval. Fig. 6(b) the lower plot shows the treatment effect for each period with the bars showing the 50% interval. The horizontal axis shows the time period and the vertical axis for Fig.6 (a) shows the instances of OBE and for Fig.6 (b) shows the treatment effect of the Noom app.

5 Conclusion

Bayesian Data Analysis is a powerful tool for incorporating data and prior information into models. It is well-suited to analyze RCTs, where effects are small and data has complicated hierarchical structures. In this paper, we have explained the steps of Bayesian Data Analysis and shown how they can be used to analyze an RCT that evaluates treatments for binge-eating disorder (BED) and bulimia nervosa (BN). We argue that Bayesian methodologies are well suited to analyzing RCTs on eating disorder behavior because they allow researchers to fit models that portray uncertainty accurately and also taking into account heterogenous treatment effects. We demonstrate this by analyzing the impact of a smartphone app on binge eating behavior by fitting a hierarchical poisson model with individual level and time level effects, and time-varying treatment effects. These unique representations of assumptions can be useful to researchers, and we have shown that their use will be beneficial for analysis of this data.

6 Appendix

Example priors

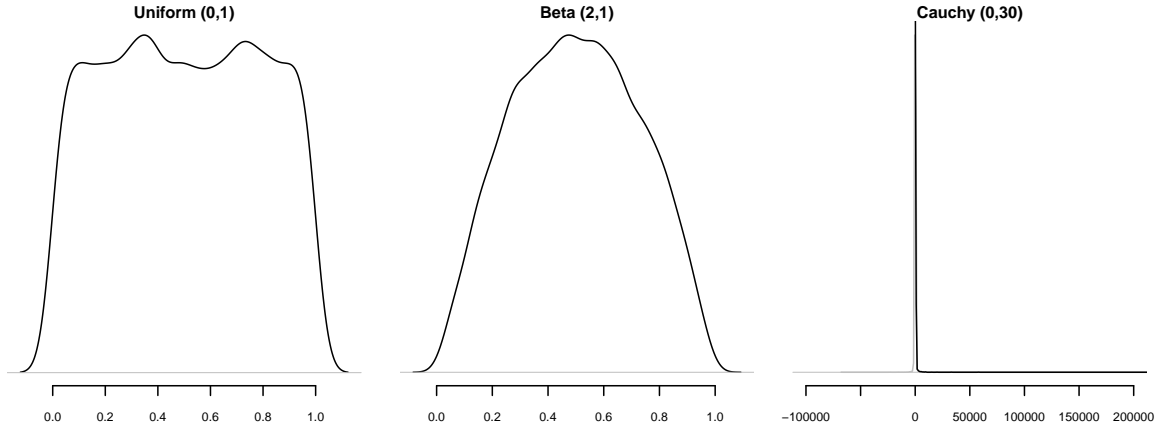


Figure 7: *Fig. 7(a) shows a Uniform (0,1) distribution. Fig. 7(b) shows a Beta (2,1) distribution. Fig. 7(c) shows a Cauchy distribution, where the black line shows the positive half and the gray line shows the negative half.*

In Figure 7, we show some sample prior distributions. The Uniform (0,1) distribution (Fig. 7(a)) can be used when we know the quantity of interest is constrained to be between 0 and 1 but believe that all values between 0 and 1 are equally likely. If our quantity of interest is between 0 and 1 but unlikely to take extreme values we can use a Beta (2,1) (Fig. 7(b)) prior as shown. $\text{Cauchy}^+(0,30)$ is the half-Cauchy distribution with location 0 and scale 30 (Fig. 7(c)). This is a good prior for variances because along with restricting the distribution to the positive real line it places most of the mass at 0 but allows for long smooth tails that the HMC algorithm can explore.

Stan Model

Listing 1: Stan code

```
data{\\ Declare data
  int N;
  int N_time;
  int N_ppl;
  int OBE[N];
  vector[N] tmt;
  int person_id[N];
  int time_id[N];
  int sex[N];
  int age[N];
  int black[N];
  int other_race[N];
  int hisp[N];
  int bn[N];
}

parameters{\\ Declare parameters for estimation
  real alpha[N_ppl];
  real mu_alpha;
  real<lower = 0> tau_alpha;
  real alpha_sex;
  real alpha_age;
  real alpha_black;
  real alpha_other;
  real alpha_hisp;
  real alpha_bn;
  real beta[N_time];
  real gamma[N_time];
  real mu_gamma;
  real<lower = 0> tau_gamma;
}
```

Listing 2: Stan code contd.

```
model{
  \\ Define likelihood
  for(i in 1:N) {
    OBE[i] ~ poisson_log(alpha[person_id[i]] + alpha_sex*sex[i]
                        + alpha_black*black[i]
                        + alpha_other*other_race[i]
                        + alpha_hisp*hisp[i] + alpha_bn*bn[i]
                        + beta[time_id[i]] + alpha_age*age[i]
                        + gamma[time_id[i]]*tmt[i]);
  }
  \\ Define priors
  alpha ~ normal(mu_alpha, tau_alpha);
  mu_alpha ~ normal(5, 2);
  tau_alpha ~ cauchy(0, 50);
  alpha_black ~ normal(0, 1);
  alpha_sex ~ normal(0, 1);
  alpha_other ~ normal(0, 1);
  alpha_hisp ~ normal(0, 1);
  alpha_bn ~ normal(0, 1);
  gamma ~ normal(mu_gamma, tau_gamma);
  mu_gamma ~ normal(0, 2);
  tau_gamma ~ cauchy(0, 30);
}
```

Listing 3: Stan code contd.

```
generated quantities{
  \\ Simulate data through the model
  real OBE_pred[N];
  for(i in 1:N) {
    OBE_pred[i] = poisson_log_rng(alpha[person_id[i]]
                                   + alpha_sex*sex[i]
                                   + alpha_black*black[i]
                                   + alpha_other*other_race[i]
                                   + alpha_hisp*hisp[i]
                                   + alpha_bn*bn[i]
                                   + beta[time_id[i]]
                                   + alpha_age*age[i]
                                   + gamma[time_id[i]]*tmt[i]);
  }
}
```

References

- Betancourt, M. (2017). A conceptual introduction to hamiltonian monte carlo. *arXiv preprint arXiv:1701.02434*.
- Betancourt, M. (2018). Calibrating model-based inferences and decisions. *arXiv preprint arXiv:1803.08393*.
- Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., ... Riddell, A. (2017). Stan: A probabilistic programming language. *Journal of statistical software*, 76(1).
- Feller, A., & Gelman, A. (2015). Hierarchical models for causal effects. *Emerging Trends in the Social and Behavioral Sciences: An interdisciplinary, searchable, and linkable resource*.

- Gabry, J., Simpson, D., Vehtari, A., Betancourt, M., & Gelman, A. (2017). Visualization in bayesian workflow. *arXiv preprint arXiv:1709.01449*.
- Gelman, A. (2006). Prior distributions for variance parameters in hierarchical models (comment on article by browne and draper). *Bayesian analysis*, 1(3), 515–534.
- Gelman, A. (2013). Commentary: P values and statistical practice. *Epidemiology*, 24(1), 69–72.
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (2014). *Bayesian data analysis* (Vol. 2). CRC press Boca Raton, FL.
- Gelman, A., & Hill, J. (2006). *Data analysis using regression and multilevel/hierarchical models*. Cambridge university press.
- Grotzinger, A., Hildebrandt, T., & Yu, J. (2015). The benefits of using semi-continuous and continuous models to analyze binge eating data: A monte carlo investigation. *International Journal of Eating Disorders*, 48(6), 746–758.
- Hildebrandt, T., Michaelides, A., Mackinnon, D., Greif, R., DeBar, L., & Sysko, R. (2017). Randomized controlled trial comparing smartphone assisted versus traditional guided self-help for adults with binge eating. *International Journal of Eating Disorders*, 50(11), 1313–1322.
- Hoekstra, R., Morey, R. D., Rouder, J. N., & Wagenmakers, E.-J. (2014). Robust misinterpretation of confidence intervals. *Psychonomic bulletin & review*, 21(5), 1157–1164.
- Kruschke, J. (2014). *Doing bayesian data analysis: A tutorial with r, jags, and stan*. Academic Press.
- Morey, R. D., Hoekstra, R., Rouder, J. N., Lee, M. D., & Wagenmakers, E.-J. (2016). The fallacy of placing confidence in confidence intervals. *Psychonomic bulletin & review*, 23(1), 103–123.
- Rouder, J. N., Speckman, P. L., Sun, D., Morey, R. D., & Iverson, G. (2009). Bayesian t tests for accepting and rejecting the null hypothesis. *Psychonomic bulletin & review*, 16(2), 225–237.
- Stan Development Team. (2015). *Prior choice recommendations*. <https://github.com/stan-dev/stan/wiki/Prior-Choice-Recommendations>. (Accessed: 2018-04-27)
- Stan Development Team. (2018). *RStan: the R interface to Stan*. <http://mc-stan.org/>.

(R package version 2.17.3)

Tukey, J. W. (1977). *Exploratory data analysis* (Vol. 2). Reading, Mass.

Vehtari, A., Gelman, A., & Gabry, J. (2017). Practical bayesian model evaluation using leave-one-out cross-validation and waic. *Statistics and Computing*, 27(5), 1413–1432.