

RESEARCH

Open Access

# Statistical analysis of two arm randomized pre-post designs with one post-treatment measurement



Fei Wan

## Abstract

**Background:** Randomized pre-post designs, with outcomes measured at baseline and after treatment, have been commonly used to compare the clinical effectiveness of two competing treatments. There are vast, but often conflicting, amount of information in current literature about the best analytic methods for pre-post designs. It is challenging for applied researchers to make an informed choice.

**Methods:** We discuss six methods commonly used in literature: one way analysis of variance (“ANOVA”), analysis of covariance main effect and interaction models on the post-treatment score (“ANCOVAI” and “ANCOVAII”), ANOVA on the change score between the baseline and post-treatment scores (“ANOVA-Change”), repeated measures (“RM”) and constrained repeated measures (“cRM”) models on the baseline and post-treatment scores as joint outcomes. We review a number of study endpoints in randomized pre-post designs and identify the mean difference in the post-treatment score as the common treatment effect that all six methods target. We delineate the underlying differences and connections between these competing methods in homogeneous and heterogeneous study populations.

**Results:** ANCOVA and cRM outperform other alternative methods because their treatment effect estimators have the smallest variances. cRM has comparable performance to ANCOVAI in the homogeneous scenario and to ANCOVAII in the heterogeneous scenario. In spite of that, ANCOVA has several advantages over cRM: i) the baseline score is adjusted as covariate because it is not an outcome by definition; ii) it is very convenient to incorporate other baseline variables and easy to handle complex heteroscedasticity patterns in a linear regression framework.

**Conclusions:** ANCOVA is a simple and the most efficient approach for analyzing pre-post randomized designs.

**Keywords:** Pre-post design, ANCOVA, ANOVA, Repeated measures, Change score, Treatment effect

## Background

Two arm parallel randomized trials have been widely used to compare the clinical effectiveness of competing treatments in improving patients’ health outcomes. In these trials, continuous outcomes of interest were routinely measured at baseline (defined

as “baseline score”) and one post treatment time point (defined as “post-treatment score”). The primary purpose of designing a pre-post randomized study is to answer the scientific question of interest: is treatment *A* more effective than treatment *B*? To assess the difference in the treatment effectiveness between two treatments, we need to select a study endpoint and quantify a treatment effect. Common study endpoints include the post treatment score, the change score from baseline to post treatment, a percentage

Correspondence: [wan.fe@wustl.edu](mailto:wan.fe@wustl.edu)

Department of Surgery, Division of Public Health Sciences, Washington University School of Medicine, Campus Box 8100, 660 S. Euclid Ave, St. Louis, MO, USA



© The Author(s). 2021 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

change from baseline, and rate of change from baseline. The difference between two arms on selected study endpoints is defined as the treatment effect. Few studies have investigated the links between these different metrics of treatment effect in a randomized pre-post trial. These underlying connections are critical in understanding the equivalence among some statistical methods that may appear to be very different at the first sight. We need to be certain about the type of treatment effect each method targets and select the one that yields an unbiased and the most efficient estimator of the treatment effect of our interest.

There are a number of statistical methods commonly used in analyzing pre-post trials. We can analyze the post-treatment score using one way analysis of variance model (*ANOVA*) [1, 2], analysis of covariance model adjusting for the baseline score (*ANCOVAI*) [2–7], and ANCOVA including a baseline score by treatment interaction (*ANCOVAII*) [3, 4, 8–10]. We can also analyze the change score using *ANOVA (ANOVA-Change)* [11]. Alternatively, we can model the baseline and post-treatment scores jointly using repeated measures models (*RM*) and constrained repeated measures models (*cRM*) [10, 12–14]. Despite of the simplicity and wide application of randomized pre-post designs, which method is the best analytic approach has been a debated topic and many methodological studies have been performed to compare different statistical methods for past decades [1–13]. However, it is challenging for applied researchers to evaluate this vast, but often conflicting, amount of information in current literature and make an informed choice.

In this study we aim to review *ANOVA, ANCOVAI, ANCOVAII, ANOVA-Change, RM, and cRM* from a practical standpoint, with the focus on delineating the differences and underlying connections between them. In section **Methods**, we first provide notations and assumptions for a typical pre-post design, define homogeneous and heterogeneous study populations, and discuss some common study endpoints and the associated metrics of treatment effects. We next analytically assess differences and connections between these competing models in the homogeneous and heterogeneous scenarios by first describing each model using the same set of population mean, variance, and covariance parameters. In section **Results**, we compare the relative efficiency of these competing methods theoretically using three simulated weight loss trial examples (homogeneous data, heterogeneous data with balanced design, heterogeneous data with unbalanced design). In the last two sections, we discuss the results and give recommendation on the best analytical approach in a randomized pre-post design.

## Methods

### A hypothetical weight loss trial and metrics of treatment effects

#### Notations

In a hypothetical two arm parallel weight loss trial comparing the effect of a new drug (“treatment”) and a placebo (“control”) in reducing participants’ body weights, we use  $Y_{ijt}$  to denote body weight of the  $i$  th subject ( $i = 1, 2, 3, \dots, n_j$ ) in the  $j$ th treatment arm ( $j = 0, 1$ ) at the  $t$  th time ( $t = t_0, t_1$ ).  $n_0$  and  $n_1$  are the number of subjects in the control and treatment arms.

We denote the mean baseline weights for the treatment and control arms by  $\mu_{1t_0}$  and  $\mu_{0t_0}$ , respectively. Random allocation guarantees  $\mu_{1t_0} = \mu_{0t_0}$  and we let  $\mu_{t_0}$  denote the overall mean baseline weight. The mean weights of the treatment and control arms at time  $t_1$  are denoted by  $\mu_{1t_1}$  and  $\mu_{0t_1}$ , respectively (Fig. 1). We define homogeneous and heterogeneous study populations as follows:

- i) **The homogeneous scenario:** every participant has the same pattern of variance and covariance structure for their baseline and post-treatment weights, which is parameterized as below:

$$\Sigma = \begin{bmatrix} \sigma_0^2 & \rho\sigma_0\sigma_1 \\ \rho\sigma_0\sigma_1 & \sigma_1^2 \end{bmatrix},$$

where  $\sigma_0^2$  and  $\sigma_1^2$  are the variances of the baseline and post-treatment weights,  $\rho$  is the correlation coefficient between the baseline and post-treatment weights.

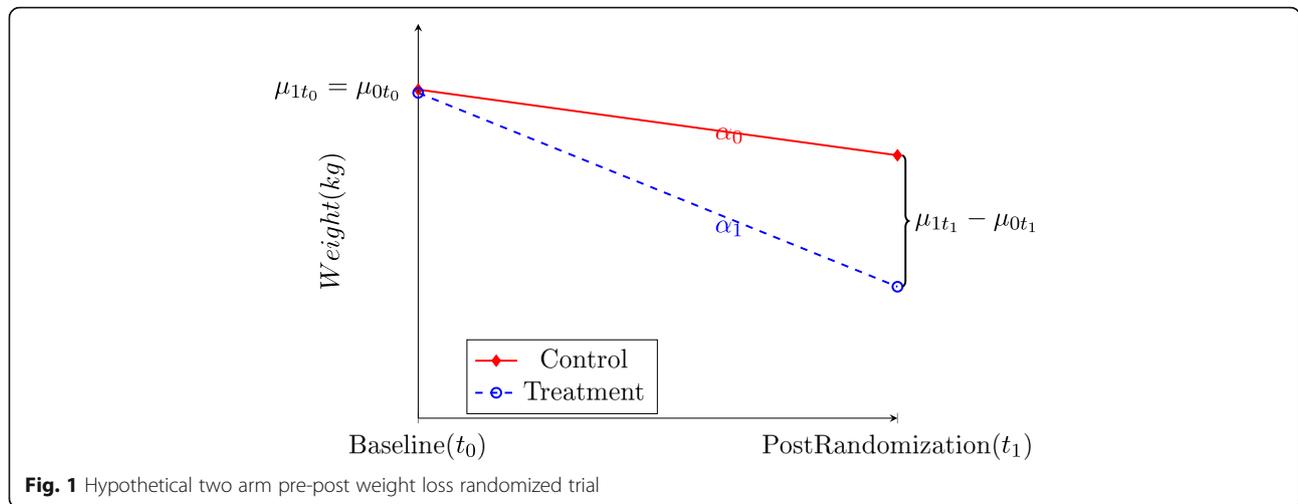
- ii) **The heterogeneous scenario:** variance and covariance structures of the baseline and post-treatment weights differ between the treatment and control arms. Formally, we have

$$\Sigma_0 = \begin{bmatrix} \sigma_0^2 & \rho_0\sigma_0\sigma_{01} \\ \rho_0\sigma_0\sigma_{01} & \sigma_{01}^2 \end{bmatrix},$$

and

$$\Sigma_1 = \begin{bmatrix} \sigma_0^2 & \rho_1\sigma_0\sigma_{11} \\ \rho_1\sigma_0\sigma_{11} & \sigma_{11}^2 \end{bmatrix},$$

where  $\sigma_0^2$  is the common variance of the baseline body weight in the control and treatment arms. Randomization guarantees that the variances of the baseline weights in both arms are equal to  $\sigma_0^2$ .  $\sigma_{01}^2$  and  $\sigma_{11}^2$  are the variances of the post-treatment weight in the control and treatment arms.  $\rho_0$  and  $\rho_1$  are the correlation coefficients between the baseline and post-treatment weights in the control and treatment arms,



respectively. In practice, participants may respond to the treatment more differently so that variability of the post-treatment weight tends to be larger in the treatment arm than in the control arm and the correlation between pre- and post-treatment weights are usually stronger in the control arm than in the treatment arm. i.e.,  $\rho_0 > \rho_1$  and  $\sigma_{11}^2 > \sigma_{01}^2$ .

**Metrics of treatment effect**

We discuss the following three metrics of treatment effect commonly reported in pre-post trials:

- i) The primary endpoint is the post-treatment weight measured at  $t_1$ . The difference in the mean post-treatment weights of two arms is defined as a treatment effect, which is parameterized as follows:

$$\tau = \mu_{1t_1} - \mu_{0t_1}$$

For example, if  $\tau = -10$ , we can interpret the results as “at the end of the trial, the mean weight was 10 pounds lower in the treatment group than in the control group.”

- ii) The primary endpoint is the change score calculated by subtracting the baseline weight from the post-treatment weight. i.e.,  $\Delta_{ij} = Y_{ijt_1} - Y_{ijt_0}$ . The difference in the mean change scores of two arms is a treatment effect. Formally, we have:

$$\tilde{\tau} = (\mu_{1t_1} - \mu_{1t_0}) - (\mu_{0t_1} - \mu_{0t_0})$$

e.g. if  $\tilde{\tau} = -10$ , this difference is usually interpreted as “weight reductions were 10 pounds greater in the treatment group than in the control group”. Since

randomization ensures  $\mu_{0t_1} = \mu_{0t_0}$ , it follows directly  $\tilde{\tau} = \tau$ . When we code “0” for  $t_0$  and “1” for  $t_1$ , the mean change score for each arm can also be interpreted as the mean change rate per unit time for each arm, represented by slopes in Fig. 1. Thus, the difference in slopes, denoted by  $\tilde{\tau} = \alpha_1 - \alpha_0$ , is also equivalent to  $\tau$ . As shown in previous section, ANOVA and ANCOVA target  $\tau$ , ANOVA-CHANGE targets  $\tilde{\tau}$ , and RM targets  $\tilde{\tau}$ . However, we can compare these statistical methods targeting seemingly very different types of treatment effects in a meaningful way because of the equivalence between  $\tau$ ,  $\tilde{\tau}$ , and  $\tilde{\tau}$  in randomized pre-post designs.

- iii) The primary endpoint is the percent change from baseline weight, denoted by  $\phi_{ij} = \frac{(Y_{ijt_1} - Y_{ijt_0})}{Y_{ijt_0}}$ . The mean difference in the percent change between two arms is defined as a treatment effect and parameterized as follows:

$$\tau^* = \bar{\phi}_1 - \bar{\phi}_0,$$

where  $\bar{\phi}_1$  and  $\bar{\phi}_0$  are the mean percent changes of the treatment and control arms. Although the percent change is popular among clinical researchers, this metric has several drawbacks [1, 15, 16]: i) the percent change is a function of ratio  $\frac{Y_{ijt_1}}{Y_{ijt_0}}$ . The distribution of the percent change is highly skewed. Analyzing it with normal-theory based statistical methods is not justified and non-parametric statistical methods are generally less powerful; ii) the percent change is not a symmetric measure. For example, the mean weight of adults over 20 in US is 197.8 pound for men and 170.5 pound for women. The mean difference is 27.3 pound between men and women. Men weight 16% (i.e.,  $100 \times ((197.8 - 170.5) / 170.5)$ ) more

than women, whereas women weight 13.8% (i.e.,  $100 \times ((197.8 - 170.5) / 197.8)$ ) less than men. The differences could be different depending on which sex is used as divisor; iii) the percent change is not an additive measure. For example, if a participant’s weight increases by 10% in first 6 months and fall by 10% for the next 6 months, the 2 % changes do not cancel out. The participant’s weight at the end would be only 99% of the participant’s starting weight.

**Statistical models**

In this section, we focus on six methods that estimate  $\tau$ . We describe each statistical model using the same set of population mean, variance, and covariance parameters defined in section [Methods](#) for homogeneous and heterogeneous scenarios, separately. For each method, we present the closed-form expressions of the point estimator of treatment effect and its variance. It often goes unnoticed in practice that different statistical methods have different types of variances (i.e., conditional vs. unconditional variances) associated with their treatment effect estimators. For example, the OLS model-based variances for ANCOVA are conditional because OLS assumes the baseline weight is fixed. Generally speaking, the baseline weight is random because we rarely enroll participants into randomized trials based on predetermined values of the baseline weight. Thus, the unconditional variance and the corresponding unconditional inference is of greater interest because we want the findings derived from the current sample to be generalizable to the population of interest. We will discuss in details whether the OLS model-based conditional inference (i.e., test statistics and  $p$ -values from standard statistical softwares) for ANCOVA is still valid for unconditional hypothesis testing and the potential fixes that we can use to draw valid unconditional inference if the usual OLS model-based inference is biased.

**When the study population is homogeneous**

**Method 1: ANOVA modeling post treatment measure (“ANOVA-Post”).** We model the post-treatment body weight  $Y_{ijt_1}$  using the binary treatment indicator  $G_{ij}$  (1 if in the treatment arm; 0 if in the control arm) as follows:

$$Y_{ijt_1} = \beta_0^{(1)} + \beta_1^{(1)} G_{ij} + e_{ij}^{(1)}, i = 1, 2, \dots, n_j; j = 0, 1; \tag{1}$$

$$e_{ij}^{(1)} \sim N(0, \sigma_1^2),$$

where  $\beta_0^{(1)} = \mu_{0t_1}, \beta_1^{(1)} = \mu_{1t_1} - \mu_{0t_1} = \tau$ , and  $e_{ij}^{(1)}$  is independently and identically distributed (i.i.d) random error.  $\beta_1^{(1)}$  represents the treatment effect. Model (1) is homoscedastic with a constant residual variance  $\sigma_1^2$ .

We can fit an ordinary least squares (OLS) regression to estimate the coefficients and standard errors of model (1). The closed-form expressions of the OLS estimator  $\hat{\beta}_{1,ols}^{(1)}$  and its “unconditional” variance, denoted by  $var(\hat{\beta}_{1,ols}^{(1)})$ , are presented in [Table 1](#).  $\hat{\beta}_{1,ols}^{(1)}$  is estimated by the sample group mean difference in the post-treatment weight between two arms.  $\hat{\beta}_{1,ols}^{(1)}$  is unbiased for  $\tau$ . The OLS model-based variance of  $\hat{\beta}_{1,ols}^{(1)}$  assuming known  $\sigma_1^2$  is:

$$var_{ols}(\hat{\beta}_{1,ols}^{(1)}) = \frac{\sigma_1^2}{\sum_{j=0}^1 \sum_{i=1}^{n_j} (G_{ij} - G_{..})^2},$$

where  $G_{..} = \frac{\sum_{j=0}^1 \sum_{i=1}^{n_j} G_{ij}}{n_0 + n_1} = \frac{n_1}{n_0 + n_1}$ .  $\sigma_1^2$  is estimated by

$$\hat{\sigma}_1^2 = \frac{\sum_{j=0}^1 \sum_{i=1}^{n_j} (y_{ijt_1} - \hat{y}_{ijt_1}^{(1)})^2}{(n_0 + n_1 - 2)},$$

where  $\hat{y}_{ijt_1}^{(1)} = \hat{\beta}_{0,ols}^{(1)} + \hat{\beta}_{1,ols}^{(1)} G_{ij}$  is the predicted value from model (1). We let  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(1)})$  denote the OLS model-based variance estimator with  $\hat{\sigma}_1^2$  substituted for  $\sigma_1^2$ , which is output by standard statistical softwares ([Table 1](#)). Since  $\sum_{j=0}^1 \sum_{i=1}^{n_j} (G_{ij} - G_{..})^2 = \frac{n_0 n_1}{n_0 + n_1}$ , it follows that  $var_{ols}(\hat{\beta}_{1,ols}^{(1)}) = var(\hat{\beta}_{1,ols}^{(1)})$ . It is well established that  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(1)})$  is unbiased for  $var_{ols}(\hat{\beta}_{1,ols}^{(1)})$ . Thus,  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(1)})$  is unbiased for  $var(\hat{\beta}_{1,ols}^{(1)})$ . The usual OLS model-based inference (i.e., test statistics  $t = \frac{\hat{\beta}_{1,ols}^{(1)}}{\sqrt{\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(1)})}}$  and the associated  $p$ -value) is valid for testing  $H_0: \tau = 0$  unconditionally.

**Method 2: ANCOVA modeling post treatment measure (“ANCOVAI”).** We model the post-treatment weight  $Y_{ijt_1}$  using the binary treatment indicator  $G_{ij}$  and the baseline weight  $Y_{ijt_0}$ .

$$Y_{ijt_1} = \beta_0^{(2)} + \beta_1^{(2)} G_{ij} + \beta_2^{(2)} Y_{ijt_0} + e_{ij}^{(2)}, i = 1, 2, \dots, n_j; j = 0, 1; \tag{2}$$

$$e_{ij}^{(2)} \sim N(0, \sigma_{e^{(2)}}^2) \text{ and } \sigma_{e^{(2)}}^2 = (1 - \rho^2) \sigma_1^2$$

, where  $\beta_0^{(2)} = \mu_{0t_1} - \rho \frac{\sigma_1}{\sigma_0} \mu_{t_0}$ ,  $\beta_1^{(2)} = \tau$ ,  $\beta_2^{(2)} = \rho \frac{\sigma_1}{\sigma_0}$ , and  $e_{ij}^{(2)}$  is i.i.d random error.  $\beta_1^{(2)}$  measures the treatment effect  $\tau$  and  $\beta_2^{(2)}$  represents the slope of the pre-post association between  $Y_{ijt_1}$  and  $Y_{ijt_0}$ . Model (2) has a common residual variance  $\sigma_{e^{(2)}}^2$  and implicitly assumes that two arms share the common baseline mean  $\mu_{t_0}$ .

**Table 1** Estimators of treatment effect and variance estimators in a homogeneous study population

Model	Estimator of treatment effect ( $\tau$ )	Type <sup>a</sup>	True variance of treatment effect estimator	OLS model based variance estimator
ANOVA-Post	$\hat{\beta}_{1,obs}^{(1)} = \bar{Y}_{.1t_1} - \bar{Y}_{.0t_1}$	U	$var(\hat{\beta}_{1,obs}^{(1)}) = \frac{\sigma_1^2}{n_0} + \frac{\sigma_1^2}{n_1}$	$\widehat{var}_{obs}(\hat{\beta}_{1,obs}^{(1)}) = \frac{\hat{\sigma}_1^2}{\sum_{j=0}^1 \sum_{i=1}^{n_j} (G_j - G_i)^2}$ $\hat{\sigma}_1^2 = \frac{\sum_{j=0}^1 \sum_{i=1}^{n_j} (Y_{jt_1} - \bar{Y}_{.jt_1})^2}{(n_0 + n_1 - 2)}$
ANCOVA-Post I	$\hat{\beta}_{1,obs}^{(2)} = (\bar{Y}_{.1t_1} - \bar{Y}_{.0t_1}) - \hat{\rho}_{2,obs}^{(2)} (\bar{Y}_{.1t_0} - \bar{Y}_{.0t_0})$	C	$var(\hat{\beta}_{1,obs}^{(2)}   Y_{jt_0}) = (\frac{1}{n_0} + \frac{1}{n_1} + \frac{\bar{Y}_{.1t_0} - \bar{Y}_{.0t_0}}{\sum_{j=0}^1 \sum_{i=1}^{n_j} (Y_{jt_0} - \bar{Y}_{.jt_0})^2}) \sigma_{\epsilon^{(2)}}^2$ $\sigma_{\epsilon^{(2)}}^2 = (1 - \rho^2) \sigma_1^2$	$\widehat{var}_{obs}(\hat{\beta}_{1,obs}^{(2)}   Y_{jt_0}) = (\frac{1}{n_0} + \frac{1}{n_1} + \frac{\bar{Y}_{.1t_0} - \bar{Y}_{.0t_0}}{\sum_{j=0}^1 \sum_{i=1}^{n_j} (Y_{jt_0} - \bar{Y}_{.jt_0})^2}) \hat{\sigma}_{\epsilon^{(2)}}^2$ $\hat{\sigma}_{\epsilon^{(2)}}^2 = \frac{\sum_{j=0}^1 \sum_{i=1}^{n_j} (Y_{jt_1} - \bar{Y}_{.jt_1})^2}{(n_0 + n_1 - 4)}$
RM	$\hat{\gamma}_{3, gts}^{(3)} = (\bar{Y}_{.1t_1} - \bar{Y}_{.1t_0}) - (\bar{Y}_{.0t_1} - \bar{Y}_{.0t_0})$	U	$var(\hat{\beta}_{1,obs}^{(2)}) = (\frac{1}{n_0} + \frac{1}{n_1})(1 - \rho^2) \sigma_1^2$	
cRM	$\hat{\gamma}_{3, gts}^{(4)} = (\bar{Y}_{.1t_1} - \bar{Y}_{.0t_1}) - \frac{\rho_0 \sigma_1}{\sigma_0^2} (\bar{Y}_{.1t_0} - \bar{Y}_{.0t_0})$	U	$var(\hat{\gamma}_{3, gts}^{(3)}) = (\frac{1}{n_0} + \frac{1}{n_1})(\sigma_1^2 + \sigma_0^2 - 2\rho_0 \sigma_0 \sigma_1)$ $var(\hat{\gamma}_{3, gts}^{(4)}) = (\frac{1}{n_0} + \frac{1}{n_1})(1 - \rho^2) \sigma_1^2$	
ANOVA-Change	$\hat{\beta}_{1,obs}^{(5)} = (\bar{Y}_{.1t_1} - \bar{Y}_{.1t_0}) - (\bar{Y}_{.0t_1} - \bar{Y}_{.0t_0})$	U	$var(\hat{\beta}_{1,obs}^{(5)}) = (\frac{1}{n_0} + \frac{1}{n_1})(\sigma_1^2 + \sigma_0^2 - 2\rho_0 \sigma_0 \sigma_1)$	$\widehat{var}_{obs}(\hat{\beta}_{1,obs}^{(5)}) = \frac{\hat{\sigma}_{\epsilon^{(5)}}^2}{\sum_{j=0}^1 \sum_{i=1}^{n_j} (G_j - G_i)^2}$ $\hat{\sigma}_{\epsilon^{(5)}}^2 = \frac{\sum_{j=0}^1 \sum_{i=1}^{n_j} (Y_{jt_1} - \Delta_{jt_1})^2}{(n_0 + n_1 - 2)}$

<sup>a</sup>U- unconditional variance; C- conditional variance

The coefficients and standard errors of model (2) are also estimated using an OLS regression. The OLS estimator  $\hat{\beta}_{1,ols}^{(2)}$  is derived as the sample mean difference in the post-treatment weight adjusting for the sample mean difference in the baseline weight between two arms. The group mean difference in the baseline weight can be seen as chance imbalance in a randomized trial.  $\hat{\beta}_{1,ols}^{(2)}$  is unbiased for  $\tau$  both conditional on  $Y_{ijt_0}$  and unconditionally. The formulas of  $\hat{\beta}_{1,ols}^{(2)}$  and its “unconditional” variance  $var(\hat{\beta}_{1,ols}^{(2)})$  are listed in Table 1. However, OLS assumes that the baseline weight  $Y_{ijt_0}$  is fixed. OLS targets the conditional variance of  $\hat{\beta}_{1,ols}^{(2)}$ , denoted by  $var(\hat{\beta}_{1,ols}^{(2)}|Y_{ijt_0})$ , instead of  $var(\hat{\beta}_{1,ols}^{(2)})$ . The formula of  $var(\hat{\beta}_{1,ols}^{(2)}|Y_{ijt_0})$  with a known common residual variance  $\sigma_{\epsilon^{(2)}}^2$  is presented in Table 1. Since  $\sigma_{\epsilon^{(2)}}^2$  is generally unknown, it is estimated by the following sample residual variance:

$$\hat{\sigma}_{\epsilon^{(2)}}^2 = \frac{\sum_{j=0}^1 \sum_{i=1}^{n_j} (y_{ijt_1} - \hat{y}_{ijt_1}^{(2)})^2}{(n_0 + n_1 - 3)}$$

, where  $\hat{y}_{ijt_1}^{(2)} = \hat{\beta}_{0,ols}^{(2)} + \hat{\beta}_{1,ols}^{(2)} G_{ij} + \hat{\beta}_{2,ols}^{(2)} Y_{ijt_0}$ , the predicted value from model (2). We let  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(2)}|Y_{ijt_0})$  denote the OLS model-based variance estimator with  $\hat{\sigma}_{\epsilon^{(2)}}^2$  substituted for  $\sigma_{\epsilon^{(2)}}^2$ . Note that  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(2)}|Y_{ijt_0})$  is reported by standard statistical softwares (e.g. “proc reg” in SAS). Its formula is presented in Table 1.

Since we want to generalize our conclusions to a general population and  $Y_{ijt_0}$  can take different values from those collected in the current sample, we may wonder whether significance tests based on the model-based conditional variance assuming  $Y_{ijt_0}$  is fixed (e.g.,  $t$

$= \frac{\hat{\beta}_{1,ols}^{(2)}}{\sqrt{\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(2)}|Y_{ijt_0})}}$ ) is comparable to unconditional in-

ference (e.g.,  $t = \frac{\hat{\beta}_{1,ols}^{(2)}}{\sqrt{var(\hat{\beta}_{1,ols}^{(2)})}}$ ), in which  $Y_{ijt_0}$  is treated as random variable, for testing  $H_0: \tau = 0$ . To establish this equivalence, we need to show: i)  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(2)}|Y_{ijt_0})$  is unbiased for  $var(\hat{\beta}_{1,ols}^{(2)}|Y_{ijt_0})$ ; ii)  $var(\hat{\beta}_{1,ols}^{(2)}|Y_{ijt_0})$  is unbiased for  $var(\hat{\beta}_{1,ols}^{(2)})$ . The first part is well established in a homoscedastic linear model. The second part holds because we can show that  $var(\hat{\beta}_{1,ols}^{(2)}) = E(var(\hat{\beta}_{1,ols}^{(2)}|Y_{ijt_0}))$  using the law of total variance formula and the fact

that  $\hat{\beta}_{1,ols}^{(2)}$  is unbiased for  $\tau$ . That is, the unconditional variance of  $\hat{\beta}_{1,ols}^{(2)}$  is the average of its conditional variance over the distribution of the baseline weight. Therefore, the usual model-based standard errors and associated  $p$ -values are valid for unconditional inference [3, 5, 17].

**Method 3: Repeated measures model (“RM”):** RM models the baseline and post-treatment weights ( $Y_{ijt_0}, Y_{ijt_1}$ ) jointly using the binary treatment indicator  $G_{ij}$ , the binary time factor  $T_{ij}$ , the time by treatment interaction  $G_{ij} \times T_{ij}$  as follows:

$$Y_{ijt} = \gamma_0^{(3)} + \gamma_1^{(3)} G_{ij} + \gamma_2^{(3)} T_{ij} + \gamma_3^{(3)} G_{ij} \times T_{ij} + e_{ijt}^{(3)}, i = 1, 2, \dots, n_j; j = 0, 1; t = t_0, t_1, \tag{3}$$

$$\begin{pmatrix} e_{ijt_0}^{(3)} \\ e_{ijt_1}^{(3)} \end{pmatrix} \sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma\right),$$

When  $t_0 = 0$  and  $t_1 = 1$ ,  $\gamma_0^{(3)} = \mu_{0t_0}$ ,  $\gamma_1^{(3)} = \mu_{1t_0} - \mu_{0t_0}$ ,  $\gamma_2^{(3)} = \mu_{0t_1} - \mu_{0t_0}$ , and  $\gamma_3^{(3)} = (\mu_{1t_1} - \mu_{1t_0}) - (\mu_{0t_1} - \mu_{0t_0})$ .  $\gamma_0^{(3)}$  represents the mean baseline weight of the control arm,  $\gamma_1^{(3)}$  represents the difference in the mean baseline weights of the treatment and control arms,  $\gamma_2^{(3)}$  represents the mean change from baseline in the control arm, and  $\gamma_3^{(3)}$  is generally interpreted as the difference in the mean change from baseline in a unit time interval between the treatment and control arms (“difference in difference”), also known as the difference in slopes. We have  $\mu_{1t_0} = \mu_{0t_0}$  from random allocation and it follows that  $\gamma_1^{(3)} = 0$  and  $\gamma_3^{(3)} = \mu_{1t_1} - \mu_{1t_0} = \tau$ . Thus, testing  $H_0: \gamma_3^{(3)} = 0$  is equivalent to testing  $H_0: \tau = 0$ .

The generalized least squares (GLS) model with correlated outcomes is routinely used to estimate the coefficients and standard errors of model (3). The GLS estimator of the treatment effect  $\hat{\gamma}_{3, gls}^{(3)}$  and its variance  $var(\hat{\gamma}_{3, gls}^{(3)})$  given known variance and covariance parameters are presented in Table 1.  $\hat{\gamma}_{3, gls}^{(3)}$  is estimated by the sample mean difference in body weight change between two arms and is unbiased for  $\tau$  in a large sample. The variance and covariance parameters are generally unknown and need to be estimated using the restricted maximum likelihood (REML). The conventional maximal likelihood estimation (MLE) should be avoided. The REML variance estimator  $\widehat{var}_{reml}(\hat{\gamma}_{3, gls}^{(3)})$  is derived by plugging the REML estimators of the variance and covariance parameters (i.e.,  $\sigma_0^2, \sigma_1^2, \rho\sigma_0\sigma_1$ ) into the formula of  $var(\hat{\gamma}_{3, gls}^{(3)})$ . We use Kenward and Roger method

[18] (“ddfm = kenwardroger” in SAS proc. mixed procedure) to adjust for the potential finite sample bias in  $\widehat{var}_{reml}(\hat{\gamma}_{3, gls}^{(3)})$  because of its failure to incorporate variabilities of the REML estimators of the variance and covariance parameters. This adjustment involves inflating the variance and covariance matrix and computing an adjusted approximation degrees of freedom.

**Method 4: Constrained Repeated measures Model (“cRM”):** By specifying  $\gamma_1^{(3)}$  in the model, **RM** model (3) assumes the mean baseline weight is different between two arms. Liang and Zeger [8] proposed the following **cRM** model by fixing  $\gamma_1^{(3)} = 0$  to force the treatment and control arms to have the same intercept. Intuitively, **cRM** is more efficient than **RM** because **cRM** estimates one less parameter. Formally, we model the baseline and post-treatment weights ( $Y_{ijt_0}, Y_{ijt_1}$ ) jointly using the binary factor  $T_{ij}$ , a time by treatment interaction  $G_{ij} \times T_{ij}$  in the following **cRM** model:

$$Y_{ijt} = \gamma_0^{(4)} + \gamma_2^{(4)} T_{ij} + \gamma_3^{(4)} G_{ij} \times T_{ij} + e_{ijt}^{(4)}, i = 1, 2, \dots, n_j; j = 0, 1; t = t_0, t_1 \quad (4)$$

$$\begin{pmatrix} e_{ijt_0}^{(4)} \\ e_{ijt_1}^{(4)} \end{pmatrix} \sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma\right),$$

where  $\gamma_0^{(4)} = \mu_{t_0}, \gamma_2^{(4)} = \mu_{0t_1} - \mu_{0t_0}$ , and  $\gamma_3^{(4)} = \tau$ . Interpretations of  $\gamma_0^{(4)}, \gamma_2^{(4)}$ , and  $\gamma_3^{(4)}$  are the same as their counterparts in **RM**. The formulas of the GLS point estimator  $\hat{\gamma}_{3, gls}^{(4)}$  and its variance  $var(\hat{\gamma}_{3, gls}^{(4)})$  are listed in Table 1.  $\hat{\gamma}_{3, gls}^{(4)}$  is unbiased for  $\tau$  asymptotically. The empirical or the model-based variance estimate for  $var(\hat{\gamma}_{3, gls}^{(4)})$  is derived using REML in the same way as a regular **RM** model.

**Method 5: ANOVA with change score (“ANOVA-Change”):** We model change score  $\Delta_{ij} = Y_{ijt_1} - Y_{ijt_0}$  using the binary treatment indicator  $G_{ij}$  as follows:

$$\Delta_{ij} = \beta_0^{(5)} + \beta_1^{(5)} G_{ij} + e_{ij}^{(5)}, i = 1, 2, \dots, n_j; j = 0, 1; \quad (5)$$

$$e_{ij}^{(5)} \sim N(0, \sigma_{e^{(5)}}^2) \text{ and } \sigma_{e^{(5)}}^2 = \sigma_1^2 + \sigma_0^2 - 2\rho\sigma_0\sigma_1,$$

where  $\beta_0^{(5)} = \mu_{0t_1} - \mu_{0t_0}, \beta_1^{(5)} = (\mu_{1t_1} - \mu_{1t_0}) - (\mu_{0t_1} - \mu_{0t_0})$ , and  $e_{ij}^{(5)}$  is i.i.d random error.  $\beta_0^{(5)}$  measures the mean difference score in the control arm.  $\beta_1^{(5)}$  measures the treatment effect  $\tilde{\tau}$ . Since  $\mu_{1t_0} = \mu_{0t_0}$  due to randomization at baseline,  $\beta_1^{(5)}$  is reduced to  $\tau$ . The closed-form expressions of  $\hat{\beta}_{1,ols}^{(5)}$  and  $var(\hat{\beta}_{1,ols}^{(5)})$  are listed in Table 1.  $\hat{\beta}_{1,ols}^{(5)}$  is derived as the sample mean difference in the change score between two arms (“difference in difference”) and

is unbiased for  $\tau$ . The OLS model-based variance of  $\hat{\beta}_{1,ols}^{(5)}$  assuming known  $\sigma_{e^{(5)}}^2$  is

$$var_{ols}(\hat{\beta}_{1,ols}^{(5)}) = \frac{\sigma_{e^{(5)}}^2}{\sum_{j=0}^1 \sum_{i=1}^{n_j} (G_{ij} - G_{..})^2},$$

where  $G_{..} = \frac{\sum_{j=0}^1 \sum_{i=1}^{n_j} G_{ij}}{n_0 + n_1} = \frac{n_1}{n_0 + n_1}$ .  $\sigma_{e^{(5)}}^2$  is estimated by

$$\hat{\sigma}_{e^{(5)}}^2 = \frac{\sum_{j=0}^1 \sum_{i=1}^{n_j} (\Delta_{ij} - \hat{\Delta}_{ij}^{(5)})^2}{(n_0 + n_1 - 2)},$$

where  $\hat{\Delta}_{ij}^{(5)}$  is the fitted value from model (5). We let  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(5)})$  denote the OLS model-based variance estimator with  $\hat{\sigma}_{e^{(5)}}^2$  substituted for  $\sigma_{e^{(5)}}^2$  Table 1, which is reported by standard statistical softwares. Since  $\sum_{j=0}^1 \sum_{i=1}^{n_j} (G_{ij} - G_{..})^2 = \frac{n_0 n_1}{n_0 + n_1}$ , it follows that  $var_{ols}(\hat{\beta}_{1,ols}^{(5)}) = var(\hat{\beta}_{1,ols}^{(5)})$ . It is well established that  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(5)})$  is unbiased for  $var_{ols}(\hat{\beta}_{1,ols}^{(5)})$ , and thus for  $var(\hat{\beta}_{1,ols}^{(5)})$ . The usual OLS model-based inference is valid for unconditional hypothesis testing.

**When the study population is heterogeneous**

**Method 6: ANCOVAII:** Different variance and covariance structures in the treatment and control arms suggest a baseline measurement by treatment interaction term in ANCOVA [2, 3, 9, 10]. To estimate  $\tau$  using an interaction model, we first compute the mean centered baseline weight  $\tilde{Y}_{ijt_0}$  by subtracting the overall mean baseline weight from individual baseline weights. i.e.,  $\tilde{Y}_{ijt_0} = Y_{ijt_0} - \mu_{t_0}$ . We then model the post-treatment body weight  $Y_{ijt_1}$  using the binary treatment indicator  $G_{ij}$ , the mean centered baseline weight  $\tilde{Y}_{ijt_0}$ , and the baseline weight by treatment interaction  $G_{ij} \times \tilde{Y}_{ijt_0}$  as follows:

$$Y_{ijt_1} = \beta_0^{(6)} + \beta_1^{(6)} G_{ij} + \beta_2^{(6)} \tilde{Y}_{ijt_0} + \beta_3^{(6)} G_{ij} \times \tilde{Y}_{ijt_0} + e_{ij}^{(6)}, i = 1, 2, \dots, n_j; j = 0, 1; \quad (6)$$

$$e_{i0}^{(6)} \sim N\left(0, \sigma_{e_0^{(6)}}^2\right) \text{ and } \sigma_{e_0^{(6)}}^2 = (1 - \rho_0^2) \sigma_{01}^2$$

$$e_{i1}^{(6)} \sim N\left(0, \sigma_{e_1^{(6)}}^2\right) \text{ and } \sigma_{e_1^{(6)}}^2 = (1 - \rho_1^2) \sigma_{11}^2$$

, where  $\beta_0^{(6)} = \mu_{0t_1}, \beta_1^{(6)} = \tau, \beta_2^{(6)} = \rho_0 \frac{\sigma_{0t_0}}{\sigma_0}$ , and  $\beta_3^{(6)} = \rho_1 \frac{\sigma_{1t_1}}{\sigma_0} - \rho_0 \frac{\sigma_{0t_0}}{\sigma_0}$ .  $e_{i0}^{(6)}$  and  $e_{i1}^{(6)}$  are i.i.d random errors in the control and treatment arms.  $\beta_1^{(6)}$  measures the treatment effect.  $\beta_2^{(6)}$  is the regression slope of the baseline body weight in the control arm.  $\beta_3^{(6)}$  measures the difference

in the regression slopes of the baseline weight between the treatment and control arms. Model (6) is heteroscedastic because the error terms in the treatment and control arms have different residual variances.

As presented in Table 2, the OLS estimator  $\hat{\beta}_{1,ols}^{(6)}$  is the adjusted mean difference in the post-treatment body weights controlling for a weighted mean difference of the baseline body weights between two arms with unequal weighting coefficients for treatment and control arms (i.e.,  $\hat{\beta}_{2,ols}^{(6)} + \hat{\beta}_{3,ols}^{(6)}$  for the treatment group, and  $\hat{\beta}_{2,ols}^{(6)}$  for the control group).  $\hat{\beta}_{1,ols}^{(6)}$  is unbiased for  $\tau$ . The conditional variance of  $\hat{\beta}_{1,ols}^{(6)}$ , denoted by  $var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$ , incorporates two different residual variances  $\sigma_{\epsilon_0}^2$  and  $\sigma_{\epsilon_1}^2$  (Table 2). Standard statistical softwares such as SAS do not output  $var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  because OLS incorrectly assumes a common residual variance  $\sigma_{\epsilon^{(6)}}^2$ , which is the following weighted average of  $\sigma_{\epsilon_0}^2$  and  $\sigma_{\epsilon_1}^2$ :

$$\sigma_{\epsilon^{(6)}}^2 = \frac{n_0}{n_0 + n_1} \sigma_{\epsilon_0}^2 + \frac{n_1}{n_0 + n_1} \sigma_{\epsilon_1}^2$$

We let  $var_{ols}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  denote the OLS model-based conditional variance of  $\hat{\beta}_{1,ols}^{(6)}$  incorporating  $\sigma_{\epsilon^{(6)}}^2$  (Table 2). Since  $\sigma_{\epsilon^{(6)}}^2$  is generally unknown,  $\hat{\sigma}_{\epsilon^{(6)}}^2$  is estimated by

$$\hat{\sigma}_{\epsilon^{(6)}}^2 = \frac{\sum_{j=0}^1 \sum_{i=1}^{n_j} (y_{ijt_1} - \hat{y}_{ijt_1})^2}{(n_0 + n_1 - 4)},$$

where  $\hat{y}_{ijt_1}$  is the predicted value of  $y_{ijt_1}$ . We let  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  denote the OLS model-based variance estimator of  $\hat{\beta}_{1,ols}^{(6)}$  with  $\hat{\sigma}_{\epsilon^{(6)}}^2$  substituted for  $\sigma_{\epsilon^{(6)}}^2$ . and known constant  $\hat{\mu}_{t_0}$  (Table 2).  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  is reported by standard statistical softwares (e.g., “proc reg” in SAS). To assess the validity of the model-based standard errors and  $p$ -values from a regular *ANCOVA* model for unconditional inference, we need to examine: i) whether  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  is unbiased for  $var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$ ; ii) whether  $var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  is unbiased for  $var(\hat{\beta}_{1,ols}^{(6)})$ .

First,  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  is unbiased for  $var_{ols}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$ . However, the unbiasedness of  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  as an estimator of  $var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  depends on the relation-

ship between  $var_{ols}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  and  $var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$ . Asymptotically, we have

$$\begin{aligned} \Delta_{\hat{\beta}_{1,ols}^{(6)}} &= var_{ols}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0}) - var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0}) \\ &= \left( \sigma_{\epsilon_0}^2 - \sigma_{\epsilon_1}^2 \right) \left( \frac{1}{n_1} - \frac{1}{n_0} \right) \end{aligned}$$

It can be shown in a balanced design ( $n_0 = n_1$ ),

$$var_{ols}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0}) \approx var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0}).$$

Thus,  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  is nearly unbiased for  $var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  [3]. When the design is unbalanced ( $n_0 \neq n_1$ ),

$$var_{ols}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0}) \neq var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0}).$$

Hence,  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  is biased for  $var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$ . Due to heteroscedasticity,  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  overestimates  $var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  if the group with a larger residual variance has larger sample size and the group with a smaller residual variance has smaller sample size, and otherwise may underestimate  $var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  [3, 4].

Second, the common mean baseline weight  $\mu_{t_0}$  is generally unknown. We need to estimate  $\mu_{t_0}$  in  $\tilde{Y}_{ijt_0}$  using the overall sample mean  $\hat{\mu}_{t_0} = \frac{\sum_{j=0}^1 \sum_{i=1}^{n_j} Y_{ijt_0}}{n_0 + n_1}$  but ANCOVA treats  $\hat{\mu}_{t_0}$  as fixed and fails to capture this additional variability in the conditional variances. As shown below, it turns out that  $var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  underestimates  $var(\hat{\beta}_{1,ols}^{(6)})$  by a factor of  $\beta_{3,ols}^{(6)2} var(\hat{\mu}_{t_0})$  [3]:

$$var(\hat{\beta}_{1,ols}^{(6)}) = E \left( var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0}) \right) + \beta_{3,ols}^{(6)2} var(\hat{\mu}_{t_0}).$$

Thus, the OLS model-based conditional inference is biased for unconditional hypothesis testing because of heteroscedasticity and neglecting of sampling variability in  $\hat{\mu}_{t_0}$ . To fix these two problems, we can use the following adjusted heteroscedasticity-consistent (HC) variance estimator to replace  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  for valid unconditional inference:

$$\widehat{var}_{aHC}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0}) = \widehat{var}_{HC}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0}) + \hat{\beta}_{3,ols}^{(6)2} \frac{\hat{\sigma}_0^2}{n_0 + n_1},$$

where  $\widehat{var}_{HC}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  is a HC variance estimator for  $var(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  [19] and can be output from standard softwares. HC variance estimators are consistent (i.e., unbiased in large sample). Among all available HC variance estimators, HC2 was shown to have the best

**Table 2** Estimators of treatment effect and variance estimators in a heterogeneous study population

Model	Estimator of treatment effect (t)	Type	True variance of treatment effect estimator	Variance estimator from OLS model
ANCOVA-Post II	$\hat{\beta}_{1,ob}^{(6)} = \bar{Y}_{1,1} - (\hat{\beta}_{2,ob}^{(6)} + \hat{\beta}_{3,ob}^{(6)})\bar{Y}_{1,10}$ $- (\bar{Y}_{0,0} - \hat{\beta}_{2,ob}^{(6)}\bar{Y}_{0,0})$	C	$\text{var}(\hat{\beta}_{1,ob}^{(6)}   \bar{Y}_{1,0}) = \left(\frac{1}{n_0} + \frac{\bar{Y}_{1,0}^2}{\sum_{i=1}^{n_0} (\bar{Y}_{1,0} - \bar{Y}_{1,0})^2}\right) \sigma_{\epsilon_1}^{(6)} + \left(\frac{1}{n_1} + \frac{\bar{Y}_{1,0}^2}{\sum_{i=1}^{n_1} (\bar{Y}_{1,0} - \bar{Y}_{1,0})^2}\right) \sigma_{\epsilon_2}^{(6)}$ $\sigma_{\epsilon_1}^{(6)} = (1 - \rho_0^2) \sigma_{\epsilon_1}^2, \quad \sigma_{\epsilon_2}^{(6)} = (1 - \rho_1^2) \sigma_{\epsilon_1}^2$	$\widehat{\text{var}}_{ob}(\hat{\beta}_{1,ob}^{(6)}   \bar{Y}_{1,0}) = \left(\frac{1}{n_0} + \frac{1}{n_1} + \frac{\bar{Y}_{1,0}^2}{\sum_{i=1}^{n_0} (\bar{Y}_{1,0} - \bar{Y}_{1,0})^2} + \frac{\bar{Y}_{1,0}^2}{\sum_{i=1}^{n_1} (\bar{Y}_{1,0} - \bar{Y}_{1,0})^2}\right) \hat{\sigma}_{\epsilon}^{(6)}$ $\hat{\sigma}_{\epsilon}^{(6)} = \frac{\sum_{i=1}^{n_0} \sum_{j=1}^{n_1} (Y_{1i} - \bar{Y}_{1,0})^2}{(n_0 + n_1 - 3)}$
ANCOVA-Post I	$\hat{\beta}_{1,ob}^{(7)} = \bar{Y}_{1,1} - \bar{Y}_{0,0} - \hat{\beta}_{2,ob}^{(7)}(\bar{Y}_{1,0} - \bar{Y}_{0,0})$	C	$\text{var}(\hat{\beta}_{1,ob}^{(6)}   \bar{Y}_{1,0}) = \frac{1}{n_0} (1 - \rho_0^2) \sigma_{\epsilon_1}^2 + \frac{1}{n_1} (1 - \rho_1^2) \sigma_{\epsilon_2}^2 + \left(\rho_1 \frac{\sigma_{11}}{\sigma_0} - \rho_0 \frac{\sigma_{11}}{\sigma_0}\right)^2 \frac{\sigma_{\epsilon_1}^2}{n_0 + n_1}$ $\text{var}(\hat{\beta}_{1,ob}^{(7)}   \bar{Y}_{1,0}) = \left(\frac{1}{n_0} + \frac{\sum_{i=1}^{n_0} (Y_{1i} - \bar{Y}_{0,0})^2 (\bar{Y}_{1,0} - \bar{Y}_{0,0})}{\sum_{j=0}^{n_0} \sum_{i=1}^{n_0} (Y_{1i} - \bar{Y}_{0,0})^2}\right) \sigma_{\epsilon_1}^{(7)} + \left(\frac{1}{n_1} + \frac{\sum_{i=1}^{n_1} (Y_{1i} - \bar{Y}_{1,0})^2 (\bar{Y}_{1,0} - \bar{Y}_{1,0})}{\sum_{j=0}^{n_1} \sum_{i=1}^{n_1} (Y_{1i} - \bar{Y}_{1,0})^2}\right) \sigma_{\epsilon_2}^{(7)}$ $\sigma_{\epsilon_1}^{(7)} = (1 - \rho_0^2) \sigma_{\epsilon_1}^2, \quad \sigma_{\epsilon_2}^{(7)} = (1 - \rho_1^2) \sigma_{\epsilon_1}^2$	$\widehat{\text{var}}_{ob}(\hat{\beta}_{1,ob}^{(7)}   \bar{Y}_{1,0}) = \left(\frac{1}{n_0} + \frac{1}{n_1} + \frac{\sum_{i=1}^{n_0} (Y_{1i} - \bar{Y}_{0,0})^2 (\bar{Y}_{1,0} - \bar{Y}_{0,0})}{\sum_{j=0}^{n_0} \sum_{i=1}^{n_0} (Y_{1i} - \bar{Y}_{0,0})^2} + \frac{\sum_{i=1}^{n_1} (Y_{1i} - \bar{Y}_{1,0})^2 (\bar{Y}_{1,0} - \bar{Y}_{1,0})}{\sum_{j=0}^{n_1} \sum_{i=1}^{n_1} (Y_{1i} - \bar{Y}_{1,0})^2}\right) \hat{\sigma}_{\epsilon}^{(7)}$ $\hat{\sigma}_{\epsilon}^{(7)} = \frac{\sum_{i=1}^{n_0} \sum_{j=1}^{n_1} (Y_{1i} - \bar{Y}_{0,0})^2}{(n_0 + n_1 - 4)}$
CRM	$\hat{\gamma}_{3, ob}^{(4)} = \bar{Y}_{1,1} - \bar{Y}_{0,0} - \frac{(\rho_0 \sigma_{01} - \rho_1 \sigma_{01})}{\sigma_0^2} (\bar{Y}_{1,0} - \bar{Y}_{0,0})$ $- \frac{\rho_0 \sigma_{01}}{\sigma_0^2} (\bar{Y}_{1,0} - \bar{Y}_{0,0})$	U	$\text{var}(\hat{\gamma}_{3, ob}^{(4)}) = \frac{1}{n_0} [(1 - \rho_0^2) \sigma_{01}^2 + ((\rho_1 \frac{\sigma_{11}}{\sigma_0} - \rho_0 \frac{\sigma_{11}}{\sigma_0}) \rho_1)^2 \sigma_0^2 + (1 - \rho_1^2) \sigma_{11}^2] + \frac{1}{n_1} [(1 - \rho_1^2) \sigma_{11}^2 + ((\rho_1 \frac{\sigma_{11}}{\sigma_0} - \rho_0 \frac{\sigma_{11}}{\sigma_0}) \rho_0)^2 \sigma_0^2]$ $\text{var}(\hat{\gamma}_{3, ob}^{(4)}) = \frac{1}{n_0} [(1 - \rho_0^2) \sigma_{01}^2 + ((\rho_1 \frac{\sigma_{11}}{\sigma_0} - \rho_0 \frac{\sigma_{11}}{\sigma_0}) \rho_1)^2 \sigma_0^2 + \frac{1}{n_1} [(1 - \rho_1^2) \sigma_{11}^2 + ((\rho_1 \frac{\sigma_{11}}{\sigma_0} - \rho_0 \frac{\sigma_{11}}{\sigma_0}) \rho_0)^2 \sigma_0^2]$	

performance in finite samples [3, 4] (e.g. “HCCMETHOD = 2” in proc. reg or “EMPIRICAL” in proc. mixed, SAS).  $\hat{\beta}_{3,ols}^{(6)}$  is the OLS estimator of  $\beta_3^{(6)}$ , and  $\hat{\sigma}_0^2$  is the overall sample variance of the baseline body weight. It follows directly that  $\widehat{var}_{aHC}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})$  is asymptotically unbiased for  $var(\hat{\beta}_{1,ols}^{(6)})$  and we can construct a valid test  $t = \frac{\hat{\beta}_{1,ols}^{(6)}}{\sqrt{\widehat{var}_{aHC}(\hat{\beta}_{1,ols}^{(6)}|\tilde{Y}_{ijt_0})}}$  for testing  $H_0: \tau = 0$  unconditionally.

**Method 7ANCOVAI:** We model the post-treatment weight  $Y_{ijt_1}$  using the binary treatment  $G$  and the baseline weight  $Y_{ijt_0}$ :

$$Y_{ijt_1} = \beta_0^{(7)} + \beta_1^{(7)}G_{ij} + \beta_2^{(7)}Y_{ijt_0} + e_{ij}^{(7)} \tag{7}$$

$$e_{i0}^{(7)} \sim N\left(0, \sigma_{e_0}^2\right) \text{ and } \sigma_{e_0}^2 = (1-\rho_0^2)\sigma_{01}^2 + \left(\beta_3^{(6)}p_1\right)^2\sigma_0^2$$

$$e_{i1}^{(7)} \sim N\left(0, \sigma_{e_1}^2\right) \text{ and } \sigma_{e_1}^2 = (1-\rho_1^2)\sigma_{11}^2 + \left(\beta_3^{(6)}p_0\right)^2\sigma_0^2$$

, where  $\beta_0^{(7)} = \beta_0^{(6)} - \beta_3^{(6)}p_0\mu_0$ , and  $\beta_1^{(7)} = \tau$ .  $e_{i0}^{(7)}$  and  $e_{i1}^{(7)}$  are random errors in the control and treatment arms. Since  $e_{i0}^{(7)}$  and  $e_{i1}^{(7)}$  have different variances in general, model (7) is heteroscedastic and the severity of heteroscedasticity is determined by the correlation coefficient, the variances of the post-treatment weights in two arms, and whether the design is balanced.

As shown in Table 2, the OLS estimator  $\hat{\beta}_{1,ols}^{(7)}$  is an adjusted mean difference in the post-treatment weights controlling for a weighted mean difference of the baseline weights between two arms with equal weighting coefficient for the treatment and control arms (i.e.,  $\hat{\beta}_{2,ols}^{(7)}$  for both arms).  $\hat{\beta}_{1,ols}^{(7)}$  is unbiased for  $\tau$ . The true conditional variance  $var(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  incorporates two different residual variances. Similar to ANCOVAII, the OLS model-based inference for ANCOVAI also mistakenly assumes a constant residual variance  $\sigma_{e^{(7)}}^2$ , which is a weighted average of  $\sigma_{e_0}^2$  and  $\sigma_{e_1}^2$ , as follows:

$$\sigma_{e^{(7)}}^2 = \frac{n_0}{n_0 + n_1} \sigma_{e_0}^2 + \frac{n_1}{n_0 + n_1} \sigma_{e_1}^2.$$

Since  $\sigma_{e^{(7)}}^2$  is unknown, it is estimated by

$$\hat{\sigma}_{e^{(7)}}^2 = \frac{\sum_{j=0}^1 \sum_{i=1}^{n_j} (y_{ijt_1} - \hat{y}_{ijt_1})^2}{n_0 + n_1 - 3},$$

where  $\hat{y}_{ijt_1}$  is the predicted value of  $y_{ijt_1}$  from model (7). The closed form expressions of the OLS model-based conditional variance  $var_{ols}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  incorporating  $\sigma_{e^{(7)}}^2$  and the OLS model-based variance estimator  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  with  $\hat{\sigma}_{e^{(7)}}^2$  substituted for  $\sigma_{e^{(7)}}^2$  are given in Table 2. Recall that standard statistical softwares report  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$ . To show the model-based standard errors and  $p$ -values are valid for unconditional inference, we need to examine: i) whether  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  is unbiased for  $var(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$ ; ii) whether  $var(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  is unbiased for  $var(\hat{\beta}_{1,ols}^{(7)})$ .

First,  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  is unbiased for  $var_{ols}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  but the unbiasedness of  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  as an estimator of  $var(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  depends on the relationship between  $var_{ols}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  and  $var(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$ . Asymptotically, we have

$$\begin{aligned} \Delta_{\hat{\beta}_{1,ols}^{(7)}} &= var_{ols}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0}) - var(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0}) \\ &= \left(\sigma_{e_0}^2 - \sigma_{e_1}^2\right) \left(\frac{1}{n_1} - \frac{1}{n_0}\right) \end{aligned}$$

When sample sizes are equal between two arms, we have

$$var_{ols}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0}) \approx var(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0}).$$

Thus,  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  is nearly unbiased for  $var(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  in a balanced design [3]. When sample sizes are not equal between two arms,

$$var_{ols}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0}) \neq var(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0}),$$

it follows directly that  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  is biased for  $var(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  due to heteroscedasticity.  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  may over-estimate  $var(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  when the group with a larger residual variance has larger sample size and the group with a smaller residual variance has smaller sample size, and otherwise may underestimate  $var(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  [3, 4]. ANCOVAI is robust against heteroscedasticity in a balanced design, but not in an unbalanced design.

Second, different from *ANCOVAII*,  $var(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  is unbiased for  $var(\hat{\beta}_{1,ols}^{(7)})$  because  $var(\hat{\beta}_{1,ols}^{(7)}) = E(var(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0}))$ .

Thus, the model-based standard errors and *p*-values are valid for unconditional inference in a balanced design but are biased in an unbalanced design only due to heteroscedasticity. This bias can be easily corrected by replacing  $\widehat{var}_{ols}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  with an HC variance estimator  $\widehat{var}_{HC}(\hat{\beta}_{1,ols}^{(7)}|Y_{ijt_0})$  [4, 19] and corrected *ANCOVAI* will provide valid unconditional inference.

**Constrained Repeated Measures heterogeneous variance model (“cRM”):** We model the baseline and post-treatment weights ( $Y_{ijt_0}, Y_{ijt_1}$ ) jointly using the binary time point  $T_{ij}$ , time by treatment interaction  $G_{ij} \times T_{ij}$ :

$$Y_{ijt} = \gamma_0^{(8)} + \gamma_1^{(8)} T_{ij} + \gamma_2^{(8)} G_{ij} \times T_{ij} + e_{ijt}^{(8)} \quad j = 0, 1; i = 1, 2, \dots, n_j. \tag{8}$$

$$\begin{pmatrix} e_{i0t_0}^{(8)} \\ e_{i0t_1}^{(8)} \end{pmatrix} \sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma_0\right) \text{ in the control arm,}$$

$$\begin{pmatrix} e_{i1t_0}^{(8)} \\ e_{i1t_1}^{(8)} \end{pmatrix} \sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma_1\right) \text{ in the treatment arm,}$$

where  $\gamma_0^{(8)} = \mu_{t_0}, \gamma_1^{(8)} = \mu_{0t_1} - \mu_{0t_0}$ , and  $\gamma_2^{(8)} = \tau$ . Noting that subjects in the treatment and control arms have different variance-covariance structures for the association between the pre- and post-treatment weights, we fit a *cRM* heterogeneous variance GLS model with group specific variance-covariance structure (“repeated/group=” in SAS proc. mixed procedure specifies distinct variance-covariance structure for each treatment arm). The formulas of  $\hat{\gamma}_{2,gl_s}^{(8)}$  and  $var(\hat{\gamma}_{2,gl_s}^{(8)})$  are listed in Table 2. The GLS estimator  $\hat{\gamma}_{2,gl_s}^{(8)}$  is asymptotically unbiased for  $\gamma_2^{(8)}$ . REML is used to derive the empirical or model-based variance estimator  $\widehat{var}_{reml}(\hat{\gamma}_{2,gl_s}^{(8)})$ .

**Results**

All treatment effect estimators, except the ANOVA estimator, are expressed as the mean difference in post-treatment measurements adjusting for the chance imbalance in baseline measurement between two arms in certain ways. Nonetheless, all estimators are unbiased for  $\tau$ . To compare these competing methods, we evaluate the efficiency of point estimators of treatment effect by comparing their “unconditional” variances. Since the hypothesis testing of no treatment effect is based on dividing the point estimator by its standard error (i.e., variance divided by sample size) and rejecting the null hypothesis

when this ratio exceeds a given threshold, the method that produces unbiased point estimate with the smallest unconditional variance is preferred because standard error in the dominator of statistical test determines the statistical power.

**When study population is homogeneous**

*ANCOVAI* is a more efficient alternative to *ANOVA* because  $var(\hat{\beta}_{1,ols}^{(2)}) \leq var(\hat{\beta}_{1,ols}^{(1)})$  (Table 1). This advantage of ANCOVA over ANOVA can also be observed from the fact that the residual error variance of *ANCOVAI* is less than the residual error variance of *ANOVA* (i.e.,  $(1-\rho^2)\sigma_1^2 \leq \sigma_1^2$ ). When the correlation coefficient  $\rho$  becomes larger, the *ANCOVAI* estimator has smaller variance. Since  $Y_{ijt_1}$  and  $Y_{ijt_0}$  are highly correlated in general, the inclusion of  $Y_{ijt_0}$  in *ANCOVAI* explains away some variability in  $Y_{ijt_1}$  and thus reduces the residual variance and yields a more efficient estimator of treatment effect than *ANOVA*.

*ANOVA-Change* and *RM* have exactly same point estimators of  $\tau$  and thus have the same variances (Table 1). To compare *ANOVA-Change* or *RM* with *ANOVA*, we can derive the difference between the unconditional variances of their treatment effect estimators as follows:

$$\Delta_1 = \sigma_0(1-2\rho\sigma_1).$$

When  $\rho < \frac{1}{2\sigma_1}$ ,  $\Delta_1 > 0$  and *ANOVA* outperforms *ANOVA-Change* and *RM* because the *ANOVA* estimator has smaller variance. When  $\rho > \frac{1}{2\sigma_1}$ ,  $\Delta_1 < 0$  and *ANOVA* underperforms the other two methods.

It can be shown that the difference between the unconditional variances of the *ANCOVAI* or *cRM* estimators and those of the *ANOVA-Change* or *RM* estimators are always nonnegative:

$$\begin{aligned} \Delta_2 &= (\sigma_1^2 + \sigma_0^2 - 2\rho\sigma_0\sigma_1) - (1-\rho^2)\sigma_1^2 \\ &= (\sigma_0 - \rho\sigma_1)^2 \geq 0 \end{aligned}$$

Thus, *ANOVA-Change* or *RM* is less efficient than either *ANCOVAI* or *cRM* because their estimators have larger variances. Intuitively *ANCOVAI* or *cRM* assumes that mean baseline weights in two arms are equal in a randomized study but *ANOVA-Change* or *RM* assumes that there is a baseline difference and needs to estimate an extra parameter.

As shown in Table 1, the *ANCOVAI* and *cRM* estimators of  $\tau$  are equivalent because  $\beta_{1,ols}^{(2)} = \frac{\rho\sigma_0\sigma_1}{\sigma_0^2}$ . However, *ANCOVAI* plugs in the OLS estimators  $\hat{\beta}_{1,ols}^{(2)}$ , whereas *cRM* plugs in the REML estimators of the variance and covariance parameters. The numerical difference between  $\hat{\beta}_{1,ols}^{(2)}$  and  $\hat{\gamma}_{3,gl_s}^{(4)}$  becomes negligible as sample size

increases. Because of this equivalence between  $\hat{\beta}_{1,ols}^{(2)}$  and  $\hat{\gamma}_{3,gl_s}^{(4)}$ ,  $var(\hat{\beta}_{1,ols}^{(2)})$  and  $var(\hat{\gamma}_{3,gl_s}^{(4)})$  are equal [3]. As discussed previously, *ANCOVAI* is a conditional model assuming fixed baseline covariates. Even though the model-based variance estimates are conditional, they are unbiased for the unconditional variance and thus the usual model-based conditional inference is still valid for unconditional hypothesis testing. *ANCOVAI* performs comparably to *cRM* [3, 17].

#### When study population is heterogeneous

A heterogeneous study population justifies the inclusion of a treatment by baseline weight interaction term. Thus, *ANCOVAII* is the correctly specified model, whereas *ANCOVAI* is a mis-specified model. In this case, the “conditional” treatment effect is not constant across different values of baseline weight. The “marginal” treatment effect  $\tau$  is simply the average of the conditional treatment effect over the distribution of the baseline weight and measures an overall treatment effect. As shown previously, both ANCOVA models can be used to estimate  $\tau$  even though *ANCOVAI* is mis-specified. Then, what is the advantage of using a more complex interaction model over a main effect model? It turns out the *ANCOVAII* estimator  $\hat{\beta}_{1,ols}^{(6)}$  is more efficient than the *ANCOVAI* estimator  $\hat{\beta}_{1,ols}^{(7)}$  because  $var(\hat{\beta}_{1,ols}^{(6)}) \leq var(\hat{\beta}_{1,ols}^{(7)})$  [5]. Only in a balanced design  $var(\hat{\beta}_{1,ols}^{(6)}) = var(\hat{\beta}_{1,ols}^{(7)})$  and the two ANCOVA models perform comparably. Note that the OLS model-based variance estimates for *ANCOVAI* and *II* are both biased for the corresponding unconditional variances, but the HC-variance estimators provide simple fixes.

The *ANCOVAII* and *cRM* estimators of  $\tau$  are equivalent because  $\beta_2^{(6)} + \beta_3^{(6)} = \frac{\rho_0\sigma_0\sigma_{01}}{\sigma_0^2}$  and  $\beta_2^{(6)} = \frac{\rho_1\sigma_0\sigma_{11}}{\sigma_0^2}$  (Table 2). Two methods only differ in the way two estimators are estimated. *ANCOVAII* plugs in the OLS estimators  $\hat{\beta}_{2,ols}^{(6)}$  and  $\hat{\beta}_{3,ols}^{(6)}$ , whereas *cRM* plugs in the REML estimators of the variance and covariance parameters. The numerical difference between the *ANCOVAII* and *cRM* estimators becomes smaller as sample size increases. As discussed previously, standard statistical softwares such as SAS does not output unconditional variance for *ANCOVAII* directly but the usual OLS model-based standard errors and  $p$ -values are biased for unconditional inference in heterogeneous scenario. The adjusted HC-variance estimator fixes this bias. Corrected *ANCOVAII* provides valid unconditional inference and performs comparably to *cRM*. Another alternative approach to estimate variances of the *ANCOVAI* and *II* estimators is to use bootstrap method [20].

#### Data example

No human data was used in this study. Instead we simulated three weight loss trial data sets based on a published study for three scenarios: homogeneous data, heterogeneous data with balanced and unbalanced designs as follows [21]:

- 1) The baseline weights for the control and treatment arms were generated from normal distribution with mean 88 kg and standard deviation 14 kg. Weights at 6 month after treatment for the control arm have mean 86 kg and standard deviation 15 kg. This gives a ~2.3% change from baseline. The mean and standard deviation of body weight at the sixth month in the treatment arm are 83 kg and 15 kg, respectively; This corresponds to a 5.7% change from baseline.
- 2) In the homogeneous data, the correlation coefficient between the pre- and post-treatment weights is 0.9. One hundred eighty subjects were assigned to the treatment and control arms equally. In the heterogeneous data, the correlation coefficient between the pre- and post-treatment weights in the control arm is 0.9 and 0.7 in the treatment arm. Sample sizes are ( $n_0 = 90, n_1 = 90$ ) for the balanced design and ( $n_0 = 60, n_1 = 120$ ) for the unbalanced design. We analyzed the data examples using the methods outlined in section *Methods*. The statistical results were reported in Table 3 (SAS programs are provided in the Additional file 1).

In the first data example, *ANOVA* produced the largest standard error and the largest  $p$ -value. *ANOVA-Change* and *RM* both outperformed *ANOVA* with much smaller standard errors and  $p$ -values. *ANCOVAI* and *cRM* outperformed *ANOVA-Change* and *RM* with smaller standard errors and  $p$ -values. Although *ANCOVAI* and *cRM* are equivalent when sample size is large, there are still minor numerical differences between the two in finite sample.

For the second data example with a balanced design, Fig. 2a shows that there is a strong baseline weight by treatment interaction. Both *ANCOVAI* and *II* have heteroscedastic errors by treatment arm (Fig. 2b and c). As shown in Table 2, the OLS model-based standard error of *ANCOVAI* is very similar to its HC and bootstrap standard errors. Thus, heteroscedasticity does not bias the model-based standard error of *ANCOVAI*. Although *ANCOVAII* is robust against heteroscedasticity in the balanced design, the OLS model-based standard error of *ANCOVAII* ( $s.e = 1.333$ ) is still not correct because OLS fails to consider the variability of estimating the overall mean baseline weight. The adjusted HC standard error for *ANCOVAII* is 1.402, which is closer to the model-

**Table 3** Statistical analysis of the three simulated data examples

Scenario	Method	Estimate	Standard error	<i>p</i> -value
Homogeneous	<i>ANOVA</i>	-3.089	2.106	0.144
	<i>ANCOVA I</i>	-2.422	0.955	0.0121
	<i>ANOVA-Change</i>	-2.354	0.971	0.0163
	<i>RM</i>	-2.354	0.971	0.0163
	<i>cRM</i>	-2.434	0.944	0.0108
Heterogeneous ( $n_0 = 90, n_1 = 90$ )	<i>ANCOVA I</i>	-3.203	1.403 <sup>a</sup>	0.0235
			1.397 <sup>b</sup>	0.0231
			1.400 <sup>d</sup>	n/a
	<i>ANCOVA II</i>	-3.165	1.333 <sup>a</sup>	0.0187
			1.402 <sup>c</sup>	0.0252
			1.397 <sup>d</sup>	n/a
<i>cRM</i>	-3.203	1.405	0.0241	
Heterogeneous ( $n_0 = 60, n_1 = 120$ )	<i>ANCOVA I</i>	-3.416	1.415 <sup>a</sup>	0.0167
			1.279 <sup>b</sup>	0.0083
			1.281 <sup>d</sup>	n/a
	<i>ANCOVA II</i>	-3.399	1.376 <sup>a</sup>	0.0145
			1.258 <sup>c</sup>	0.0076
			1.260 <sup>d</sup>	n/a
<i>cRM</i>	-3.396	1.262	0.0078	

<sup>a</sup>OLS regression model-based standard error

<sup>b</sup>HC standard error for ANCOVA I (main effect) model

<sup>c</sup>Modified HC standard error for ANCOVA II (interaction) model

<sup>d</sup>Bootstrapping standard error ( $n = 5000$ )

based and HC standard errors of *ANCOVA I*. The bootstrapping standard errors for *ANCOVA I* and *II* are close to their HC or adjusted HC standard errors, which suggests the HC and adjusted HC variances perform well in estimating the unconditional variances. The *cRM* estimate and its standard error are close to those from *ANCOVA I* and *II*.

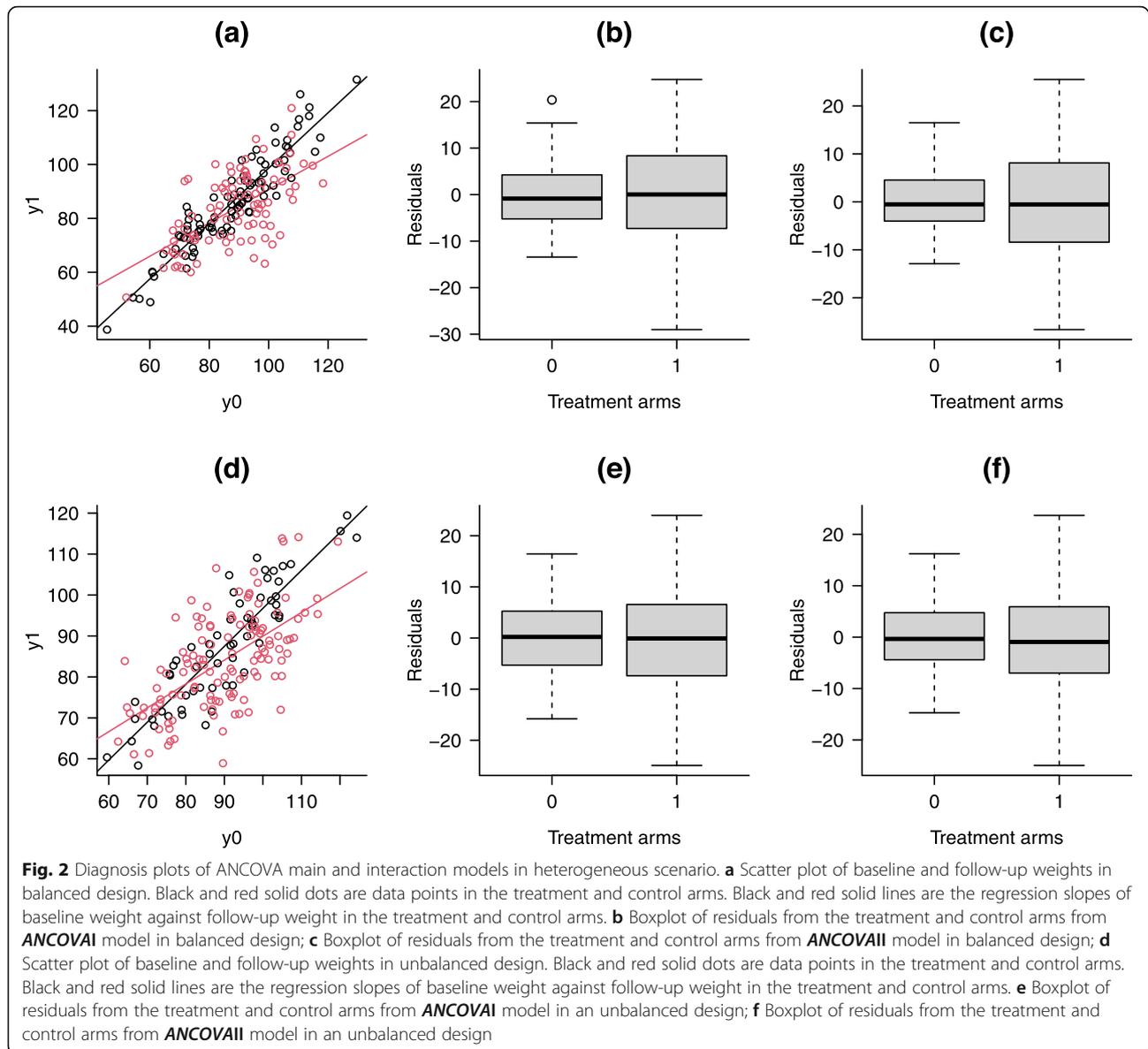
For the third example with an unbalanced design, Fig. 2d also reveals a baseline weight by treatment interaction. Both ANCOVA models have heteroscedastic errors by treatment arm (Fig. 2e and f). The model-based standard errors of *ANCOVA I* and *II* are not valid. The model-based standard errors were larger than the HC standard errors and thus overestimated the true conditional variances. Compared with *ANCOVA I*, *ANCOVA II* has a smaller HC standard error (also smaller *p*-value) and thus is slightly more efficient. The adjusted HC standard error for *ANCOVA II* is very close to the model-based standard error for *cRM*. The bootstrapping standard errors for *ANCOVA I* and *II* are very close to their HC or adjusted HC standard errors.

## Discussion

In this study we compare the efficiency of six unbiased methods analyzing pre-post designs. We found

*ANCOVA* and *cRM* are the equally most efficient methods compared with other alternatives in homogeneous and heterogeneous scenarios. In this study, we focus on the scenario in which randomization is properly performed and these competing methods all target the same causal quantity. In the scenarios where the treatment is not properly randomized or not randomized at all (e.g., in an observational study), the baseline score will not be balanced by design. In this case these competing methods may target different causal quantities. Debate over using change-score analysis (or *RM*) versus *ANCOVA* in the non-randomized setting, generally known as the lord's paradox, is a well-known example [22, 23].

The majority of previous studies has only examined homogeneous study population. In this setting, *ANOVA* is one of the least efficient approaches for analyzing pre-post designs because it does not utilize any baseline information. *ANOVA-Change* and *RM* incorporate the baseline score as part of outcome, whereas *ANCOVA I* includes the baseline score as a covariate. *ANCOVA I* outperforms *ANOVA-Change* and *RM* because *ANCOVA I* utilizes the assumption that the baseline scores are balanced between two arms in a randomized study. Thus, change score is a less efficient way to utilize the



baseline score than including the baseline score as a covariate. Since we seldom can control the values of the baseline score in randomized trials, the OLS assumption that the baseline score is fixed casts doubt on the validity of ANCOVA for hypothesis testing [6, 12]. Crager proved **ANCOVA I** is valid for unconditional inference in homogeneous scenario [6]. This conclusion can be simply attributed to that the conditional variance of the **ANCOVA I** estimator is an unbiased estimate for its unconditional variance [3].

A few studies investigated further a heterogeneous scenario [3, 4, 10, 12, 24]. Although the heterogeneity justifies the inclusion of the baseline measurement by treatment interaction term, **ANCOVA I** and **II** are both unbiased. Yang and Tsiatis showed that

**ANCOVA II** has a smaller unconditional variance estimator than that of **ANCOVA I** unless in a balanced design [9]. However, the OLS model-based variances of the **ANCOVA I** and **II** estimators, reported by standard statistical softwares, are conditional variances, not unconditional variances. The OLS model-based standard errors and associated  $p$ -values for **ANCOVA II** are generally questionable for unconditional inference, and the model-based inference for **ANCOVA I** is biased only when the design is unbalanced [3, 4, 10, 24]. With the corrected HC variance estimators, both models provide valid unconditional inference. Choosing between **ANCOVA I** and **II** then becomes an evaluation of a trade-off between simplicity and some gains in efficiency.

In homogenous setting, *cRM* was suggested as a superior choice to *ANCOVAI* because the unconditional variance of the *cRM* estimator is smaller than the conditional variance of the *ANCOVAI* estimator [25]. Kenward et al. pointed out that such direct comparison between the conditional and unconditional variances is not meaningful. Since both estimators are equivalent, it can be shown that *cRM* coupled with REML and Kenward-roger adjustment performs almost identically to *ANCOVAI* in finite samples [17]. In heterogeneous scenario, *cRM* is comparable to *ANCOVAII* [3]. In presence of missing data, applied researchers often prefer *cRM* over ANCOVA because it can utilize all observed data but ANCOVA uses only complete cases. However, imputation methods which utilize the strong pre-post correlation, such as weighting and regression imputation, can improve the statistical power for ANCOVA without biasing estimates, making it comparable to *cRM* [17].

Furthermore, ANCOVA has several advantages over *cRM*: first, outcome should only be the variable that can be influenced by treatment. Baseline measurement is certainly not an outcome by this definition. It is conceptually more appropriate to include the baseline score as covariate, not model it as outcome [5]; Second, it is very convenient to include other baseline variables in a regression model for more efficient estimates of treatment effect. Third, it is easy to adjust for other patterns of heteroscedastic errors in an OLS regression. For example, we may expect larger variability in the post-treatment weights associated with larger baseline weights. *cRM* cannot handle this more complex type of heteroscedasticity easily. HC-variance estimators for ANCOVA are simple fixes and readily implemented in statistical softwares.

## Conclusion

Comparing with other alternative methods, ANCOVA is a simple and the most efficient approach analyzing a pre-post randomized design. When there exists a baseline score by treatment interaction, we need to assess the heteroscedasticity of ANCOVA particularly when the design is not balanced. The HC-variances should be used for valid inference when heteroscedasticity is present. Adding an interaction term in ANCOVA can gain some efficiency but not including this term does not bias results.

## Abbreviations

**ANOVA**: Analysis of variance model; **ANCOVAI**: Analysis of covariance model adjusting for the baseline measurement; **ANCOVAII**: Analysis of covariance model adjusting for the baseline measurement by treatment interaction; **RM**: Constrained repeated measure model; **cRM**: Constrained repeated measure model; **HC**: Heteroscedasticity-consistent

## Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12874-021-01323-9>.

### Additional file 1.

## Acknowledgements

Not applicable

## Author's contributions

FW developed the idea for the paper, performed analysis, and drafted the manuscript. The author(s) read and approved the final manuscript.

## Funding

Not applicable

## Availability of data and materials

SAS code is provided as the Additional file 1. There is no real data used. All data generated or analyzed during this study are included in this published article [and its supplementary information files].

## Declarations

### Ethics approval and consent to participate

Not applicable

### Consent for publication

Not applicable

### Competing interests

The author declare that he has no competing risk

Received: 10 January 2021 Accepted: 19 May 2021

Published online: 24 July 2021

## References

- Vickers AJ. Analysis of variance is easily misapplied in the analysis of randomized trials: a critique and discussion of alternative statistical approaches. *Psychosom Med*. 2005;67(4):652–5. <https://doi.org/10.1097/01.psy.0000172624.52957.a8>.
- O'Connell NS, Dai L, Jiang Y, Speiser JL, Ward R, Wei W, et al. Methods for analysis of pre-post data in clinical research: a comparison of five common methods. *J Biom Biostat*. 2017;8(1):1–8. <https://doi.org/10.4172/2155-6180.1000334>.
- Wan F. Analyzing pre-post randomized studies with one post-randomization score using repeated measures and ANCOVA models. *Stat Methods Med Res*. 2019;28(10-11):2952–74. <https://doi.org/10.1177/0962280218789972>.
- Wan F. Analyzing pre-post designs using the analysis of covariance models with and without the interaction term in a heterogeneous study population. *Stat Methods Med Res*. 2020;29(1):189–204. <https://doi.org/10.1177/0962280219827971>.
- Senn S. Change from baseline and analysis of covariance revisited. *Stat Med*. 2006;25(24):4334–44. <https://doi.org/10.1002/sim.2682>.
- Crager MR. Analysis of covariance in parallel-group clinical trials with pretreatment baseline. *Biometrics*. 1987;43(4):895–901. <https://doi.org/10.2307/2531543>.
- Frison L, Pocock SJ. Repeated measures in clinical trials: analysis using mean summary statistics and its implications for design. *Stat Med*. 1992;11(13):1685–704. <https://doi.org/10.1002/sim.4780111304>.
- Brogan DR, Kutner MH. Comparative analyses of pretest-posttest research designs. *Am Stat*. 1980;34:229–32.
- Yang L, Tsiatis AA. Efficiency study of estimators for a treatment effect in a pretest-posttest trial. *Am Stat*. 2001;55(4):314–21. <https://doi.org/10.1198/000313001753272466>.
- Winkens B, van Breukelen GJ, Schouten HJ, Berger MP. Randomized clinical trials with a pre- and a post-treatment measurement: repeated measures versus ANCOVA models. *Contemp Clin Trials*. 2007;28(6):713–9. <https://doi.org/10.1016/j.cct.2007.04.002>.
- Dimitrov DM, Rumrill J, Phillip D. Pretest-posttest designs and measurement of change. *Work*. 2003;20:159–65.

12. Chen X. The adjustment of random baseline measurements in treatment effect estimation. *J Stat Plan Inference*. 2006;136(12):4161–75. <https://doi.org/10.1016/j.jspi.2005.08.046>.
13. Jennings E. Models for pretest-posttest data: repeated measures ANOVA revisited. *J Educ Behav Stat*. 1988;13(3):273–80. <https://doi.org/10.3102/10769986013003273>.
14. Liang K, Zeger S. Longitudinal data analysis of continuous and discrete responses for pre-post designs. *Sankhya*. 2000;62:134–48.
15. Donald AB, Gregory DA. Symmetrized percent change for treatment comparisons. *Am Stat*. 2006;60:27–31.
16. Cole TJ, Altman DG. Statistics notes: what is a percentage difference? *BMJ*. 2017;358:j3663. <https://doi.org/10.1136/bmj.j3663>.
17. Kenward MG, White IR, Carpenter JR. Re: should baseline be a covariate or dependent variable in analyses of change from baseline in clinical trials? (Liu GF et al., *Stat Med* 2009; 28: 2509–30). *Stat Med*. 2010;29(13):1455–6. <https://doi.org/10.1002/sim.3868>.
18. Kenward MG, Roger JH. Small sample inference for fixed effects from restricted maximum likelihood. *Biometrics*. 1997;53(3):983–97. <https://doi.org/10.2307/2533558>.
19. Long J, Ervin L. Using heteroscedasticity: consistent standard errors in the linear regression model. *Am Stat*. 2000;54:217–24.
20. Efron B, Tibshirani RJ. An introduction to the bootstrap. New York: Chapman & Hall; 1993. <https://doi.org/10.1007/978-1-4899-4541-9>.
21. Boozer CN, Daly PA, Homel P, et al. Herbal ephedra/caffeine for weight loss: a 6-month randomized safety and efficacy trial. *Int J Obes Relat Metab Disord*. 2002;6:593–604.
22. Lord FM. A paradox in the interpretation of group comparisons. *Psychol Bull*. 1967;68(5):304–5. <https://doi.org/10.1037/h0025105>.
23. Egbewale BE, Lewis M, Sim J. Bias, precision and statistical power of analysis of covariance in the analysis of randomized trials with baseline imbalance: a simulation study. *BMC Med Res Methodol*. 2014;14(1):49. <https://doi.org/10.1186/1471-2288-14-49>.
24. Senn S. Various varying variances: the challenge of nuisance parameters to the practicing biostatistician. *Stat Methods Med Res*. 2015;24(4):403–19. <https://doi.org/10.1177/0962280214520728>.
25. Liu GF, Lu K, Mogg R, Mallick M, Mehrotra DV. Should baseline be a covariate or dependent variable in analyses of change from baseline in clinical trials? *Stat Med*. 2009;28(20):2509–30. <https://doi.org/10.1002/sim.3639>.

## Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Ready to submit your research? Choose BMC and benefit from:**

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

**At BMC, research is always in progress.**

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

