

Covariate Adjustment: practical binary endpoints

Worked Examples using Resampled Data from MISTIE III Trial

Kelly Van Lancker (kelly.vanlancker@ugent.be)

2024-07-18 01:39

Contents

Prerequisites	1
Using This Tutorial	1
Installing and Loading R Packages	1
Background: MISTIE III Study Design	2
Creating Simulated Data:	2
Exploring the data	5
Covariate Adjustment	7
Design Parameters	7
Unadjusted Analysis	7
Covariate-Adjusted Analysis	8

Prerequisites

Using This Tutorial

This tutorial contains an example dataset as well as code to illustrate how to perform covariate adjustment in practice using R.

Installing and Loading R Packages

The following packages and their dependencies need to be installed:

- table1 - Creating simple tabulations in aggregate and by treatment arm
- tidyverse - An ecosystem of packages for working with data

```
required_packages <-  
  c("table1", "tidyverse")  
  
install.packages(required_packages)
```

Once the required packages are installed, they can be loaded using `library()`

```
library(table1)
library(tidyverse)
```

Background: MISTIE III Study Design

Data used in this example are simulated using data based on the **Minimally Invasive Surgery with Thrombolysis in Intracerebral haemorrhage Evacuation** trial (MISTIE III: NCT01827046). MISTIE III was an open-label, blinded endpoint, Phase III clinical trial of minimally invasive surgery with thrombolysis in intracerebral haemorrhage evacuation. The goal was to assess whether minimally invasive catheter evacuation followed by thrombolysis, with the aim of decreasing clot size to 15 mL or less, would improve functional outcome in patients with intracerebral haemorrhage (a severe form of stroke). To this end, participants were randomized 1:1 to standard-of-care medical management, or minimal invasive surgery with Alteplase for ICH removal. Outcomes were measured at 30, 180, and 365-days post-randomization using the Modified Rankin Scale (mRS). The primary outcome was defined as having a mRS score of 0-3 measured 365 days from enrollment (defined as a *success*). Survival was also assessed, with patients administratively censored on the date of their final MRS assessment. Though the trial used covariate adaptive randomization, we ignore that in our discussion below, for simplicity (since analogous computations taking this into account give similar results), and we will assume simple randomization.

Creating Simulated Data:

The data in this template are simulated data, generated from probability models fit to the original study data, and *not the actual data from the MISTIE III trial*. A new synthetic dataset was created by resampling baseline covariates from the original data with replacement. The outcome columns in the synthetic dataset were sequentially replaced using simulated values based on predictions from a sequence of regression models based on the actual study data.

Load MISTIE III Data

The data can be loaded directly from Github:

```
data_url <-
  "https://github.com/jbetz-jhu/CovariateAdjustmentTutorial/raw/main/Simulated_MISTIE_III_v1.2.csv"

sim_miii <-
  read.csv(file = url(data_url))

# Read in data: Recast categorical variables as factors
sim_miii <-
  sim_miii %>%
  dplyr::tibble() %>%
  dplyr::mutate(
    male =
      factor(
        x = male,
        levels = 0:1,
```

```

    labels = c("0. Female", "1. Male")
  ),
  across(
    .cols = all_of(
      x = c(
        "hx_cvd",
        "hx_hyperlipidemia",
        "on_anticoagulants",
        "on_antiplatelets"
      )
    ),
    .fns = function(x)
      factor(x, levels = 0:1, labels = c("0. No", "1. Yes"))
  ),
  across(.cols = starts_with("gcs") | starts_with("mrs"),
    .fns = factor),
  ich_location =
    factor(x = ich_location,
      levels = c("Deep", "Lobar")),
  arm =
    factor(x = arm,
      levels = c("medical", "surgical")),
  tx = 1 * (arm == "surgical")
) %>%
dplyr::rename(id = sim_participant_id)

```

The complete simulated trial data without any missing values are in a `data.frame` named `sim_miii`.

- Participant Identifier
 - `id`: Participant Identifier
- Baseline Covariates
 - `age`: Age at baseline in years
 - `male`: Participant sex (1 for Male or 0 for female)
 - `hx_cvd`: Cardiovascular disease history
 - `hx_hyperlipidemia`: Hyperlipidaemia medication compliant history
 - `on_anticoagulants`: On anticoagulants medication
 - `on_antiplatelets`: On antiplatelet medication
 - `ich_location`: Intracerebral haemorrhage clot location: (Lobar, Deep)
 - `ich_s_volume`: Intracerebral hemorrhage volume on stability scan
 - `ivh_s_volume`: Intraventricular hemorrhage volume on stability scan
 - `gcs_category`: Severity of impairment as measured by Glasgow Coma Score (GCS)
- Treatment:
 - `arm`: Treatment arm (surgical versus standard medical care)
 - `tx`: Treatment arm (binary; 1 for surgical and 0 for standard medical care)
 - `ich_eot_volume`: Intracerebral hemorrhage volume on end-of-treatment scan
- Outcomes:
 - `mrs_30d`: mRS at 30 days (0-3, 4, 5, 6)
 - `mrs_30d_complete`: mRS at 30 days if no data were missing
 - `mrs_180d`: mRS at 180 days (0-2, 3, 4, 5, 6)
 - `mrs_180d_complete`: mRS at 180 days if no data were missing

- `mrs_365d`: mRS at 365 days (0-1, 2, 3, 4, 5, 6)
- `mrs_365d_complete`: mRS at 365 days if no data were missing
- `days_on_study`: days until death or administrative censoring
- `died_on_study`: participant died (1) or is censored (0)

The outcomes `mrs_30d`, `mrs_180d`, and `mrs_365d` contain missing values: the actual values before the missingness mechanism is applied are also included with the `_complete` suffix.

Define Dichotomized Outcomes

The primary outcome was defined as having a modified Rankin Scale (mRS) score of 0-3 measured 365 days from enrollment (defined as a *success*). We therefore dichotomized `mrs_30d_complete`, `mrs_180d_complete` and `mrs_365d_complete`.

```
sim_miii$mrs_bin_30d_complete = ifelse(sim_miii$mrs_30d_complete == "0-3", 1, 0)
sim_miii$mrs_bin_180d_complete = ifelse(
  sim_miii$mrs_180d_complete == "0-2" ,
  1,
  ifelse(sim_miii$mrs_180d_complete ==
    "3", 1,
    0)
)
sim_miii$mrs_bin_365d_complete = ifelse(
  sim_miii$mrs_365d_complete == "0-1" ,
  1,
  ifelse(
    sim_miii$mrs_365d_complete == "2",
    1,
    ifelse(sim_miii$mrs_365d_complete ==
      "3", 1,
      0)
  )
)
```

- Dichotomized Outcomes:

- `mrs_bin_30d_complete`: Dichotomized mRS score at 30 days (1 if mRS equals 0–3, 0 otherwise)
- `mrs_bin_180d_complete`: Dichotomized mRS score at 180 days (1 if mRS equals 0–3, 0 otherwise)
- `mrs_bin_365d_complete`: Dichotomized mRS score at 365 days (1 if mRS equals 0–3, 0 otherwise)

Reference Level for Treatment

When the treatment is a `factor` variable, we can use the `levels()` function to see the reference level (i.e., the comparator/control group): it will appear as the first level.

```
# Check reference level
levels(sim_miii$arm)
```

```
## [1] "medical" "surgical"
```

Make sure that the reference level is appropriately chosen before running analyses. In this case study, ‘medical’ is the reference level.

Exploring the data

Baseline Demographics & Stratum

Below are summary statistics of participant characteristics at baseline:

```
table1(  
  ~ age + male + hx_cvd + hx_hyperlipidemia + on_anticoagulants + on_antiplatelets +  
    ich_location + ich_s_volume + ivh_s_volume + gcs_category + ich_eot_volume |  
  arm,  
  data = sim_miii  
)
```

	medical	surgical	Overall
	(N=500)	(N=500)	(N=1000)
age			
Mean (SD)	60.2 (12.9)	60.4 (12.3)	60.3 (12.6)
Median [Min, Max]	62.0 [28.0, 85.0]	61.0 [29.0, 90.0]	61.0 [28.0, 90.0]
male			
0. Female	194 (38.8%)	200 (40.0%)	394 (39.4%)
1. Male	306 (61.2%)	300 (60.0%)	606 (60.6%)
hx_cvd			
0. No	427 (85.4%)	425 (85.0%)	852 (85.2%)
1. Yes	73 (14.6%)	75 (15.0%)	148 (14.8%)
hx_hyperlipidemia			
0. No	303 (60.6%)	316 (63.2%)	619 (61.9%)
1. Yes	197 (39.4%)	184 (36.8%)	381 (38.1%)
on_anticoagulants			
0. No	468 (93.6%)	466 (93.2%)	934 (93.4%)
1. Yes	32 (6.4%)	34 (6.8%)	66 (6.6%)
on_antiplatelets			
0. No	365 (73.0%)	343 (68.6%)	708 (70.8%)
1. Yes	135 (27.0%)	157 (31.4%)	292 (29.2%)
ich_location			
Deep	309 (61.8%)	322 (64.4%)	631 (63.1%)
Lobar	191 (38.2%)	178 (35.6%)	369 (36.9%)
ich_s_volume			
Mean (SD)	48.4 (17.3)	48.1 (16.7)	48.3 (17.0)
Median [Min, Max]	45.8 [9.70, 112]	46.9 [12.9, 106]	46.2 [9.70, 112]
ivh_s_volume			
Mean (SD)	2.17 (3.66)	2.08 (3.77)	2.13 (3.72)
Median [Min, Max]	0 [0, 32.0]	0 [0, 43.0]	0 [0, 43.0]
gcs_category			
1. Severe (3-8)	136 (27.2%)	134 (26.8%)	270 (27.0%)
2. Moderate (9-12)	226 (45.2%)	201 (40.2%)	427 (42.7%)
3. Mild (13-15)	138 (27.6%)	165 (33.0%)	303 (30.3%)
ich_eot_volume			
Mean (SD)	47.1 (17.9)	14.2 (9.67)	30.7 (21.8)
Median [Min, Max]	43.4 [11.5, 114]	12.9 [0, 59.0]	26.9 [0, 114]

Modified Rankin Scale Outcomes

Here we summarize the outcomes of the study (without missing data):

```
table1(
  ~ mrs_30d_complete + mrs_180d_complete + mrs_365d_complete + mrs_bin_30d_complete +
    mrs_bin_180d_complete + mrs_bin_365d_complete | arm,
  data = sim_miii %>% dplyr::mutate(
    mrs_bin_30d_complete =
      factor(x = mrs_bin_30d_complete,
            levels = c("0", "1")),
    mrs_bin_180d_complete =
      factor(x = mrs_bin_180d_complete,
            levels = c("0", "1")),
    mrs_bin_365d_complete =
      factor(x = mrs_bin_365d_complete,
            levels = c("0", "1"))
  )
)
```

	medical	surgical	Overall
	(N=500)	(N=500)	(N=1000)
mrs_30d_complete			
0-3	45 (9.0%)	69 (13.8%)	114 (11.4%)
4	106 (21.2%)	143 (28.6%)	249 (24.9%)
5	262 (52.4%)	259 (51.8%)	521 (52.1%)
6	87 (17.4%)	29 (5.8%)	116 (11.6%)
mrs_180d_complete			
0-2	68 (13.6%)	91 (18.2%)	159 (15.9%)
3	119 (23.8%)	136 (27.2%)	255 (25.5%)
4	103 (20.6%)	112 (22.4%)	215 (21.5%)
5	78 (15.6%)	93 (18.6%)	171 (17.1%)
6	132 (26.4%)	68 (13.6%)	200 (20.0%)
mrs_365d_complete			
0-1	31 (6.2%)	42 (8.4%)	73 (7.3%)
2	54 (10.8%)	76 (15.2%)	130 (13.0%)
3	127 (25.4%)	129 (25.8%)	256 (25.6%)
4	100 (20.0%)	110 (22.0%)	210 (21.0%)
5	52 (10.4%)	63 (12.6%)	115 (11.5%)
6	136 (27.2%)	80 (16.0%)	216 (21.6%)
mrs_bin_30d_complete			
0	455 (91.0%)	431 (86.2%)	886 (88.6%)
1	45 (9.0%)	69 (13.8%)	114 (11.4%)
mrs_bin_180d_complete			
0	313 (62.6%)	273 (54.6%)	586 (58.6%)
1	187 (37.4%)	227 (45.4%)	414 (41.4%)
mrs_bin_365d_complete			
0	288 (57.6%)	253 (50.6%)	541 (54.1%)
1	212 (42.4%)	247 (49.4%)	459 (45.9%)

Covariate Adjustment

Design Parameters

The estimand in the trial was defined as the (absolute) risk difference, that is, the difference between the population proportion of successes under assignment to treatment versus control (where control was standard of care using medical management). The total sample size of approximately 498 patients was calculated based on the assumption that 25% of the patients would have an mRS score of 0–3 in the standard medical care group versus 38% of patients in the MISTIE group and provides a power of 88% to detect such a population risk difference of 13% at a 5% significance level.

Unadjusted Analysis

Question 1.

Calculate an unadjusted estimator for the (absolute) risk difference of interest, along with a confidence interval and p -value. Make use of the function `prop.test`.

```
mistieiii_prop_test <-  
  prop.test(  
    x = with(sim_miii, table(mrs_bin_365d_complete, arm))  
  )  
  
print(mistieiii_prop_test)  
  
##  
## 2-sample test for equality of proportions with continuity correction  
##  
## data: with(sim_miii, table(mrs_bin_365d_complete, arm))  
## X-squared = 4.6553, df = 1, p-value = 0.03096  
## alternative hypothesis: two.sided  
## 95 percent confidence interval:  
## 0.00642914 0.13451859  
## sample estimates:  
## prop 1 prop 2  
## 0.5323475 0.4618736  
  
# Estimate  
-diff(mistieiii_prop_test$estimate)  
  
## prop 2  
## 0.07047387  
  
# Confidence Interval  
mistieiii_prop_test$conf.int  
  
## [1] 0.00642914 0.13451859  
## attr(,"conf.level")  
## [1] 0.95
```

```
# P-Value
mistieiii_prop_test$p.value
```

```
## [1] 0.03095778
```

Covariate-Adjusted Analysis

We can use a logistic regression to adjust for baseline covariates with a binary outcome. However, the coefficient in a logistic regression is a *conditional association*: it is the associated change in the log odds of the outcome with a unit change in the predictor when holding all other variables constant. The average treatment effect is a *marginal, not conditional quantity*. In order to obtain the average treatment effect, we must marginalize by averaging over the covariates.

Question 2.

First, regress the final outcome `mrs_bin_365d_complete` on the treatment assignment indicator `arm` and the baseline covariates `ich_s_volume`, `age`, `ivh_s_volume`, `ich_location` and `gcs_category` using the function `glm` with `family = binomial(link = "logit")`.

```
# Fit the `glm` object
mistieiii_glm <-
  glm(
    formula =
      mrs_bin_365d_complete ~
      arm + ich_s_volume + age + ivh_s_volume +
      ich_location + gcs_category,
    data = sim_miii,
    family = binomial(link = "logit")
  )
```

```
# Print a summary of the `glm` object
summary(mistieiii_glm)
```

```
##
## Call:
## glm(formula = mrs_bin_365d_complete ~ arm + ich_s_volume + age +
##      ivh_s_volume + ich_location + gcs_category, family = binomial(link = "logit"),
##      data = sim_miii)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.663157   0.408438   4.072 4.66e-05 ***
## armsurgical     0.276494   0.137076   2.017  0.0437 *
## ich_s_volume   -0.006690   0.004228  -1.583  0.1135
## age            -0.042618   0.006019  -7.081 1.43e-12 ***
## ivh_s_volume   -0.024213   0.019494  -1.242  0.2142
## ich_locationLobar  0.075840   0.153170   0.495  0.6205
## gcs_category2. Moderate (9-12) 0.977206   0.177434   5.507 3.64e-08 ***
## gcs_category3. Mild (13-15)   1.679583   0.195808   8.578 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```



```
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 1379.6 on 999 degrees of freedom
## Residual deviance: 1231.4 on 992 degrees of freedom
## AIC: 1247.4
##
## Number of Fisher Scoring iterations: 4
```

Question 3.

We obtain an estimate of the average treatment effect by plugging in estimates of the probabilities from a fitted model: this involves obtaining a prediction for each individual as if they were assigned to treatment (using the function `predict`), and a prediction for each individual assigned to control, and marginalizing over the covariates by taking the sample average.

Calculate the adjusted estimator for the (absolute) risk difference of interest, along with the variance and a confidence interval. The variance can be calculated as $1/n$ times the sample variance of

$$\frac{Z}{\hat{P}(Z=1)}[Y - \hat{E}(Y|Z=1, X)] + \hat{E}(Y|Z=1, X) - \left\{ \frac{1-Z}{1 - \hat{P}(Z=1)}[Y - \hat{E}(Y|Z=0, X)] + \hat{E}(Y|Z=1, X) \right\}$$

```
# Predict Pr{Y = 1 | Z = 1, X}
pr_y1_z1 <-
  predict(
    object = mistieiii_glm,
    newdata =
      sim_miii %>%
      dplyr::mutate(
        arm = "surgical"
      ),
    type = "response"
  )

# Predict Pr{Y = 1 | Z = 0, X}
pr_y1_z0 <-
  predict(
    object = mistieiii_glm,
    newdata =
      sim_miii %>%
      dplyr::mutate(
        arm = "medical"
      ),
    type = "response"
  )

# Estimate
adj_mean = mean(pr_y1_z1) - mean(pr_y1_z0)
print(adj_mean)
```

```
## [1] 0.05927086
```

```

# Standard Error
p_arm = mean(sim_miii$arm=="surgical")
adj_se = sqrt(
  var(
    (sim_miii$arm=="surgical")/p_arm*
    (sim_miii$mrs_bin_365d_complete-pr_y1_z1) + pr_y1_z1 -
    ((sim_miii$arm=="medical")/(1-p_arm)*
    (sim_miii$mrs_bin_365d_complete-pr_y1_z0) + pr_y1_z0) )/
    nrow(sim_miii)
  )
print(adj_se)

```

```
## [1] 0.02918874
```

```

# Confidence Interval
c(adj_mean-qnorm(0.975)*adj_se, adj_mean+qnorm(0.975)*adj_se)

```

```
## [1] 0.002061971 0.116479744
```

Question 4.

Compare the confidence intervals of the unadjusted and covariate-adjusted estimators. What do you observe?