

Illustrative GEE Analysis of Stepped Wedge Cluster Randomized Trial Data

Fan Li

Trial Data with Continuous Responses

We first read in the simulated trial data set with continuous individual responses.

```
dir<-"D:/Research/CRT Methodology/SWDSampSize/Latex/Submission/R Code"

setwd(dir)

simdata_cont<-read.csv("simdata_cont.csv", header = TRUE)
```

The first 6 rows of the data set looks like the following:

```
head(simdata_cont)

##           y ind cluster period period1 period2 period3 period4 period5
## 1 -0.8017075  1      1      1      1      0      0      0      0
## 2  1.5416701  2      1      1      1      0      0      0      0
## 3 -0.4889857  3      1      1      1      0      0      0      0
## 4 -1.4624861  4      1      1      1      0      0      0      0
## 5 -1.7894035  5      1      1      1      0      0      0      0
## 6 -1.2769037  6      1      1      1      0      0      0      0
## treatment
## 1      0
## 2      0
## 3      0
## 4      0
## 5      0
## 6      0
```

From left to right, the columns of data are response (in long format), individual id, cluster id, period id, indicator for period 1, indicator for period 2, indicator for period 3, indicator for period 4, indicator for period 5 and indicator for treatment. We will extract the following key elements from the data set. Note that this is a simulated cohort stepped wedge cluster randomized trial with 20 clusters, 20 individuals per cluster (cohort size), and 5 periods (1 baseline period).

```
# responses
y<-as.numeric(simdata_cont$y)

# cluster identifier
id<-as.numeric(simdata_cont$cluster)

# period identifier
period<-as.numeric(simdata_cont$period)
# marginal mean design matrix
X<-simdata_cont[,c("period1", "period2", "period3", "period4", "period5", "treatment")]
X<-as.matrix(X)

# treatment indicator
trt<-X[, "treatment"]
```

```

# number of clusters
n<-length(unique(id))

# number of periods
t<-dim(X)[2]-1

# cluster size (across all periods)
clsize<-as.numeric(table(id))

# cohort size per cluster (balanced)
m<-clsize/t

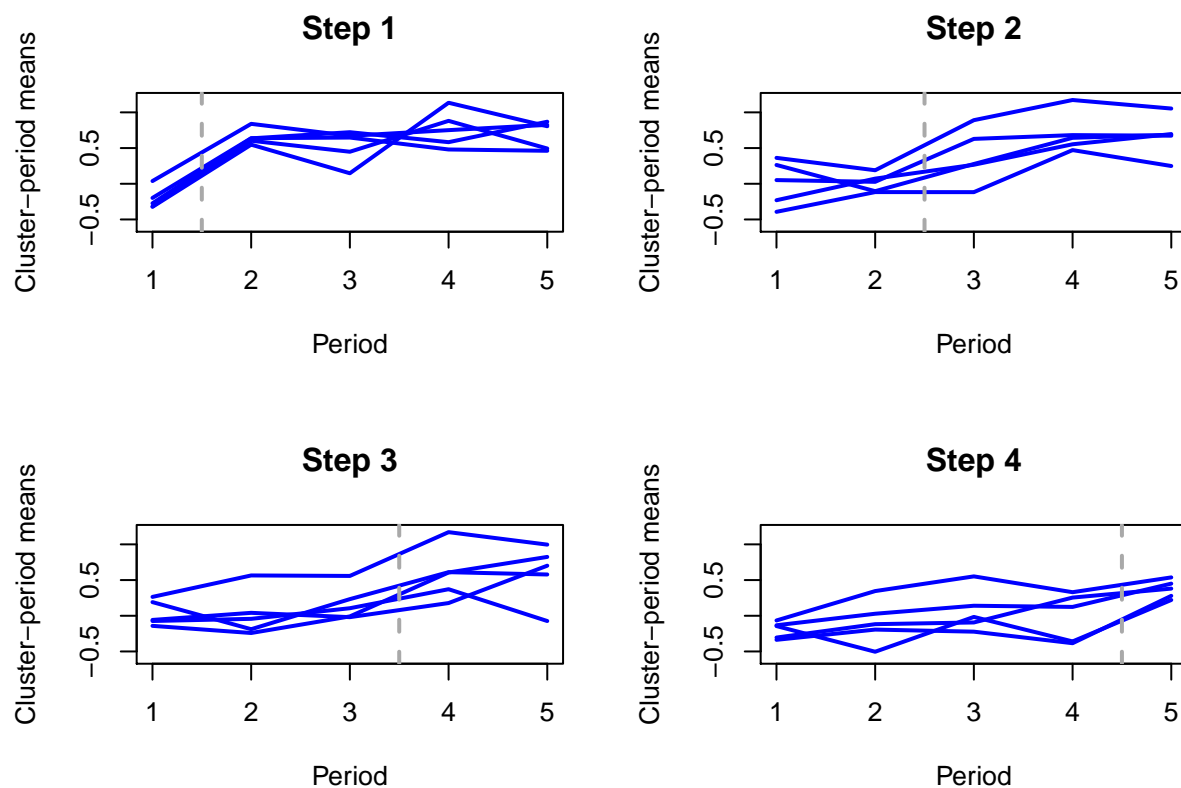
```

By summarizing the cluster-period means by the following plot, we observe that there is a gently increasing period effect over time, and the treatment appears to have a positive effect.

```

# Cluster-period means
clp_mu<-tapply(y,list(id,period),FUN=mean)

```



For the GEE and MAEE analysis of the trial data, we need to obtain the design matrix for the correlation parameters as follows.

```

# Create (large) design matrix for correlations
CREATEZ<-function(n,m,t){
  # correlation position indicators
  alpha0_pos<-1
  alpha1_pos<-2
  alpha2_pos<-3

```

```

zrow<-diag(3)
Z<-NULL

for(i in 1:n){
  mi<-m[i]
  bm1<-(1-alpha2_pos-alpha0_pos+alpha1_pos)*diag(t*mi)
  bm2<-(alpha2_pos-alpha1_pos)*kronecker(matrix(1,t,t),diag(mi))
  bm3<-(alpha0_pos-alpha1_pos)*kronecker(diag(t),matrix(1,mi,mi))
  bm4<-alpha1_pos*matrix(1,t*mi,t*mi)
  POS<-bm1+bm2+bm3+bm4

  for(j in 1:(t*mi-1)){
    for(k in (j+1):(t*mi)){
      Z<-rbind(Z,zrow[POS[j,k],])
    }
  }
  # print(i)
}
return(Z)
}

# large matrix (may take a minute to run)
Z<-CREATEZ(n,m,t)

```

We confirm the exploratory analysis of the trial data by fitting the GEE and MAEE using the following code. Detailed descriptions of the input arguments are available in the `contMAEE.R` program. Following the notations in Li, Turner and Preisser, we use the marginal mean model

$$\mu_{ijk} = \beta_j + X_{ij}\delta,$$

where the link is the identity function, β_j is the j th period effect, X_{ij} is the treatment indicator of cluster i in period j , δ is the marginal intervention effect. As indicated in Li, Turner and Preisser, the block exchangeable correlation structure is parameterized with three parameters ($\alpha_0, \alpha_1, \alpha_2$).

```

# Source the function
source("contMAEE.R")

# Implement the function
contMAEE(y=y,X=X,id=id,n=clsize,Z=Z,maxiter=25,epsilon=0.001,printrange="NO",
        shrink="ALPHA",makevone="NO")

```

```

## Loading required package: MASS

## GEE for correlated Gaussian data
## Number of Clusters: 20
## Maximum Cluster Size: 100
## Minimum Cluster Size: 100
## Number of Iterations: 6
## Results for marginal mean parameters
##      Beta      Estimate MB-stderr BC0-stderr BC1-stderr BC2-stderr
## [1,]    0 -0.09124742 0.06139865 0.04875328 0.05001981 0.05131924
## [2,]    1  0.03679558 0.06370980 0.05792417 0.06001125 0.06219307
## [3,]    2  0.06176615 0.07018814 0.06649575 0.06936958 0.07239740
## [4,]    3  0.16979517 0.07982543 0.09104943 0.09516587 0.09950651
## [5,]    4  0.12809893 0.09163032 0.08509015 0.08882935 0.09277746
## [6,]    5  0.45719138 0.06801707 0.06765181 0.07142571 0.07541453

```

```
##      BC3-stderr
## [1,] 0.04995324
## [2,] 0.05987887
## [3,] 0.06868531
## [4,] 0.09484761
## [5,] 0.08881937
## [6,] 0.07080831
##
## Results for correlation parameters
##      Alpha   Estimate BC0-stderr BC1-stderr BC2-stderr BC3-stderr
## [1,]      0 0.02864983 0.01656862 0.01699621 0.01743664 0.01700313
## [2,]      1 0.02909247 0.01728845 0.01773451 0.01819387 0.01772013
## [3,]      2 0.19823026 0.06311718 0.06527458 0.06751977 0.06499727
```

We observe a positive treatment effect 0.46, which is slightly larger but close to the value, 0.4, used in the data generation. A gently increasing period effect is reflected in the parameter estimates. The estimates for within-period correlation α_0 and within-individual correlation α_2 are close to the true values, 0.03 and 0.2, although the inter-period correlation α_1 is overestimated. For the class of bias-corrected sandwich variances, we observe that $BC0 < BC1 < BC2$ and $BC1 \approx BC3$. Note that this is an illustrative analysis of only one data set with a limited number of clusters.

Trial Data with Binary Responses

We read in the simulated trial data set with binary individual responses.

```
simdata_bin<-read.csv("simdata_bin.csv", header = TRUE)
```

The first 6 rows of the data set looks like the following:

```
head(simdata_bin)

##   y ind cluster period period1 period2 period3 period4 period5 treatment
## 1 0   1     1     1       1       0       0       0       0         0
## 2 0   2     1     1       1       0       0       0       0         0
## 3 0   3     1     1       1       0       0       0       0         0
## 4 0   4     1     1       1       0       0       0       0         0
## 5 0   5     1     1       1       0       0       0       0         0
## 6 0   6     1     1       1       0       0       0       0         0
```

We will extract the following key elements from the trial data set as before.

```
# responses
y<-as.numeric(simdata_bin$y)

# cluster identifier
id<-as.numeric(simdata_bin$cluster)

# period identifier
period<-as.numeric(simdata_bin$period)

# marginal mean design matrix
X<-simdata_bin[,c("period1", "period2", "period3", "period4", "period5", "treatment")]
X<-as.matrix(X)

# treatment indicator
trt<-X[, "treatment"]
```

```

# number of clusters
n<-length(unique(id))

# number of periods
t<-dim(X)[2]-1

# cluster size (across all periods)
clsize<-as.numeric(table(id))

# cohort size per cluster (balanced)
m<-clsize/t

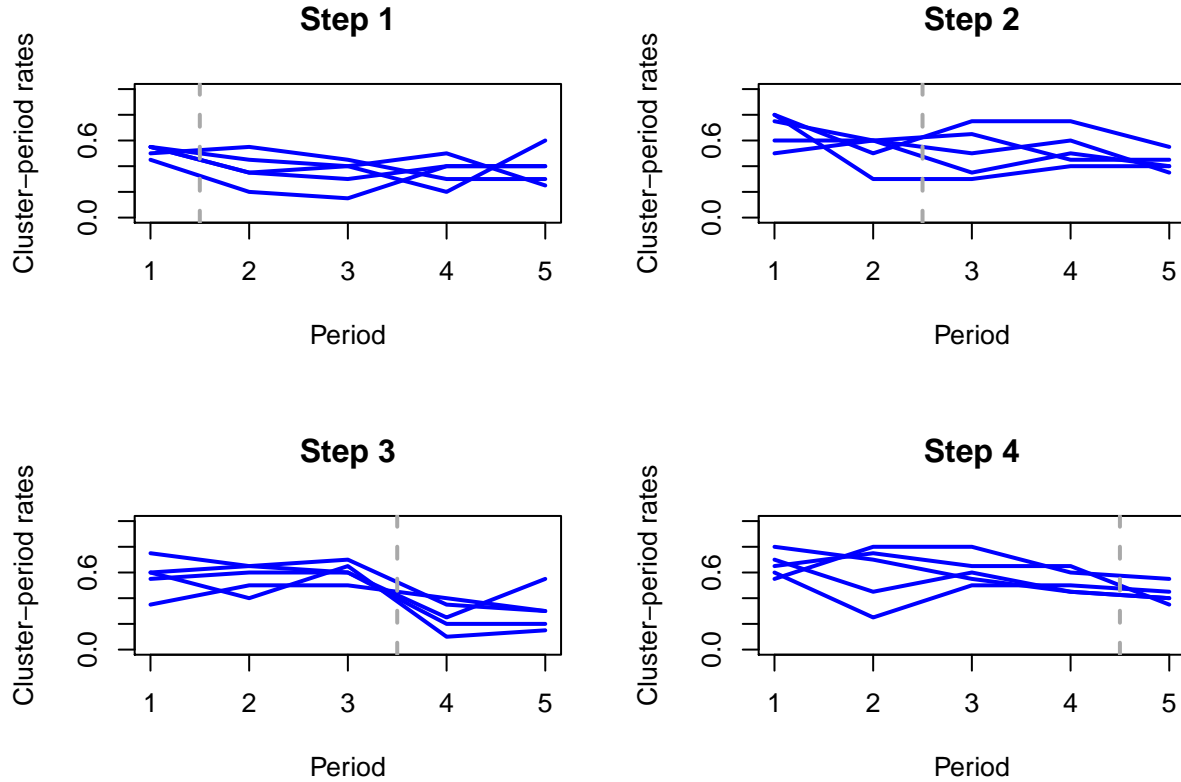
```

By summarizing the cluster-period rates by the following plot, we confirm that there is a gently decreasing period effect over time, and the treatment appears to be associated with decreased rates.

```

# Cluster-period means
clp_mu<-tapply(y,list(id,period),FUN=mean)

```



Since this trial data possess the same structure as the previous one (20 clusters, 20 individuals per cluster and 5 periods), we could use the same design matrix for estimating correlation parameters. We confirm the exploratory analysis of the trial data by fitting the GEE and MAEE in the following analysis. We use the marginal mean model

$$\text{logit}(\mu_{ijk}) = \beta_j + X_{ij}\delta,$$

where the link is the logistic function, β_j is the j th period effect, X_{ij} is the treatment indicator of cluster i in period j , δ is the marginal intervention effect on the log odds ratio scale. Again, the block exchangeable correlation structure is parameterized with three parameters $(\alpha_0, \alpha_1, \alpha_2)$.

```

# Source the function
source("binMAEE.R")

# Implement the function
binMAEE(y=y,X=X,id=id,n=clsize,Z=Z,maxiter=25,epsilon=0.001,prinrange="NO",
        shrink="ALPHA",makevone="NO")

## GEE for correlated binary data
## Number of Clusters: 20
## Maximum Cluster Size: 100
## Minimum Cluster Size: 100
## Number of Iterations: 5
## Results for marginal mean parameters
##      Beta   Estimate MB-stderr BC0-stderr BC1-stderr BC2-stderr BC3-stderr
## [1,]    0  0.4468718 0.1272800  0.1138047  0.1167623  0.1197967  0.1165604
## [2,]    1  0.2093945 0.1313458  0.1408497  0.1459532  0.1512875  0.1449121
## [3,]    2  0.4044053 0.1497445  0.1217652  0.1276724  0.1339235  0.1252425
## [4,]    3  0.1668318 0.1739910  0.1501134  0.1577612  0.1659134  0.1553603
## [5,]    4  0.1874213 0.2049487  0.1864827  0.1959001  0.2058687  0.1953348
## [6,]    5 -0.6456953 0.1606787  0.1767356  0.1864173  0.1966525  0.1848739
##
## Results for correlation parameters
##      Alpha   Estimate BC0-stderr BC1-stderr BC2-stderr BC3-stderr
## [1,]      0 0.028514568 0.012991123 0.013329305 0.013676297 0.013329762
## [2,]      1 0.007437099 0.008609916 0.008833773 0.009063449 0.008833968
## [3,]      2 0.200789999 0.022404304 0.022985365 0.023581671 0.022987060

```

We observe the marginal treatment effect in the odds ratio scale to be $\exp(\delta) = 0.52$, which is slightly larger than the value, 0.45, used in the data generation. A gently decreasing period effect is reflected in the parameter estimates. The estimates for within-period correlation α_0 and within-individual correlation α_2 are close to the true values, 0.03 and 0.2, although the inter-period correlation α_1 is underestimated. We note that this is an illustrative analysis of only one data set with a limited number of clusters.