

Average Effects Based on Regressions with a Logarithmic Link Function: A New Approach
with Stochastic Covariates

Supplementary Material A

.

Average Effects Based on Regressions with a Logarithmic Link Function: A New Approach
with Stochastic Covariates

Simulation details

In this section we provide additional information on the simulation study. First, a detailed description of the data generation process is given. Second, we state the explicit moment-based estimators for each unconditional distribution of Z used in the simulation. Third, we provide the listings for the data generation procedure as well as for the estimators.

Data simulation

As the population model, we choose a regression with log link according to Equation (9) and conditional NB distributions for Y . The model includes a dichotomous treatment variable X , a covariate Z and their interaction XZ .

In order to generate scenarios that are comparable across different distributions of Z , we computed the regression coefficients $\beta_{00}, \beta_{01}, \beta_{11}$ from their standardized versions $\beta_{00}^*, \beta_{01}^*, \beta_{11}^*$ (i.e., with regard to a z -transformed covariate $Z^* = (Z - E(Z))/\sqrt{\text{Var}(Z)}$).

These are derived as follows:

$$\begin{aligned}
 \log[E(Y|X, Z^*)] &= \beta_{00}^* + \beta_{10}^*X + \beta_{01}Z^* + \beta_{11}^*XZ^* \\
 &= \beta_{00}^* + \beta_{10}^*X + \beta_{01} \left(\frac{Z - E(Z)}{\sqrt{\text{Var}(Z)}} \right) + \beta_{11}^*X \left(\frac{Z - E(Z)}{\sqrt{\text{Var}(Z)}} \right) \\
 &= \left(\beta_{00}^* - \frac{\beta_{01}^*}{\sqrt{\text{Var}(Z)}} E(Z) \right) + \left(\beta_{10}^* - \frac{\beta_{11}^*}{\sqrt{\text{Var}(Z)}} E(Z) \right) X \\
 &\quad + \left(\frac{\beta_{01}^*}{\sqrt{\text{Var}(Z)}} \right) Z + \left(\frac{\beta_{11}^*}{\sqrt{\text{Var}(Z)}} \right) XZ
 \end{aligned}$$

Starting from the regression on X and the standardized Z^* , we derived the

corresponding regression coefficients for the unstandardized Z , which are

$$\beta_{00} = \beta_{00}^* - \frac{\beta_{01}^*}{\sqrt{\text{Var}(Z)}} E(Z), \quad (1)$$

$$\beta_{01} = \frac{\beta_{01}^*}{\sqrt{\text{Var}(Z)}}, \quad (2)$$

$$\beta_{11} = \frac{\beta_{11}^*}{\sqrt{\text{Var}(Z)}} \quad (3)$$

The standardized regression coefficients are divided by the standard deviation of Z , because the covariates are on a different scale depending on their distribution in the simulation. We manipulated the effect size of the ATE, and thus the remaining regression coefficient β_{10} was not a free parameter but determined by

$$\beta_{10} = \log \left(\frac{\Delta \cdot \sqrt{\text{Var}(Y|X=0)} + M_Z(\beta_{01}) \exp(\beta_{00})}{\exp(\beta_{00}) M_Z(\beta_{01} + \beta_{11})} \right)$$

The standardized regression coefficients were chosen $\beta_{00}^* = 0$ and $\beta_{01}^* = 0.5$ to simulate a moderate covariate effect in the control group $X = 0$. The standardized interaction parameter was varied from no interaction ($\beta_{11}^* = 0$) to strong interaction ($\beta_{11}^* = 1$) with a step width of 0.25. As the derivations above suggest, the difference between the fixed-covariate and stochastic-covariate-based standard errors should increase with the strength of the interaction parameter. The population distributions of Z were chosen to resemble distinct distributional shapes: a standard normal distribution $\mathcal{N}(0, 1)$ for a symmetric shape with zero skewness and kurtosis, a Poisson distribution $\mathcal{P}(1)$ for a positively skewed shape, and a uniform distribution $U(-\sqrt{3}, \sqrt{3})$ for a flat shape with negative kurtosis. The parameters of $\mathcal{N}(0, 1)$ and $U(-\sqrt{3}, \sqrt{3})$ represent standardized variables, that is, $E(Z) = 0$, $\text{Var}(Z) = 1$. For the Poisson distribution, however, the parameter $\lambda = 1$ was chosen to create a moderate skew. We further varied the effect size Δ (Glass, 1976) of the ATE and the sample size N . In total, we evaluated 225 experimental

scenarios (see Table 3 for more details on the design parameters).

The simulation consisted of three steps, which were carried out using the statistical software R (R Core Team, 2018). In the first step, N covariate values z_i were generated by drawing from the populational distribution of Z . In the second step, the values of the dichotomous treatment x_i were drawn from a Binomial distribution $B(1, p)$ where

$$p = P(X|Z^*) = 1 - \frac{1}{1 + \exp(Z^*)} \quad (4)$$

which was generated as a logistic function with standardized Z^* (i.e., the higher the value of Z , the higher the probability of getting into the treatment group). This means that we simulated a quasi-experimental setting, where treatment assignment is not random but depending on Z . As mentioned above, ATE and UTE are not equivalent in this case, and consideration of the covariate Z is necessary to estimate the average (causal) effect. In the third step, the conditional expectations $E(Y|X = x_i, Z = z_i)$ were computed and the observed y_i were drawn from a negative binomial distribution. The overdispersion parameter was $\alpha = 1$. Each experimental condition was replicated $R = 5,000$ times. For more details and the corresponding R code, see the supplementary material.

Moment-based estimators by distribution

For normal distributed $Z \sim \mathcal{N}(\hat{\mu}, \hat{\sigma})$ with $\hat{\mu} = \bar{z}$ and $\hat{\sigma} = \hat{\sigma}_z$

$$\widehat{\text{ATE}}_{\text{MOM}}(\hat{\beta}, \hat{\theta}) = \exp(\hat{\beta}_{00} + \hat{\beta}_{10}) \exp\left(\hat{\mu}(\hat{\beta}_{01} + \hat{\beta}_{11}) + \frac{\hat{\sigma}^2(\hat{\beta}_{01} + \hat{\beta}_{11})^2}{2}\right) - \exp(\hat{\beta}_{00}) \exp\left(\hat{\mu}\hat{\beta}_{01} + \frac{\hat{\sigma}^2\hat{\beta}_{01}^2}{2}\right) \quad (5)$$

For Poisson distributed $Z \sim \mathcal{P}(\hat{\lambda})$ where $\hat{\lambda} = \bar{z}$

$$\widehat{\text{ATE}}_{\text{MOM}}(\hat{\beta}, \hat{\theta}) = \exp(\hat{\beta}_{00} + \hat{\beta}_{10}) \exp(\hat{\lambda}(\exp(\hat{\beta}_{01} + \hat{\beta}_{11}) - 1)) - \exp(\hat{\beta}_{00}) \exp(\hat{\lambda}(\exp(\hat{\beta}_{01}) - 1)) \quad (6)$$

For Uniform distributed $Z \sim U(\hat{a}, \hat{b})$ where $\hat{a} = \bar{z} - \sqrt{3}\hat{\sigma}_z$ and $\hat{b} = \bar{z} + \sqrt{3}\hat{\sigma}_z$

$$\widehat{\text{ATE}}_{\text{MOM}}(\hat{\beta}, \hat{\theta}) = \exp(\hat{\beta}_{00} + \hat{\beta}_{10}) \frac{e^{(\hat{\beta}_{01} + \hat{\beta}_{11})\hat{b}} - e^{(\hat{\beta}_{01} + \hat{\beta}_{11})\hat{a}}}{(\hat{\beta}_{01} + \hat{\beta}_{11})(\hat{b} - \hat{a})} - \exp(\hat{\beta}_{00}) \frac{e^{\hat{\beta}_{01}\hat{b}} - e^{\hat{\beta}_{01}\hat{a}}}{\hat{\beta}_{01}(\hat{b} - \hat{a})} \quad (7)$$

Estimation of sample moments for moment-based approach

In our simulation, we used a two-step procedure for obtaining sample moment estimates and their covariance. The first step was based on a decomposition of the unconditional expectation and variance of Z with respect to the treatment groups:

$$E(Z) = E(Z|X = 1) \cdot P(X = 1) + E(Z|X = 0) \cdot P(X = 0) \quad (8)$$

$$\text{Var}(Z) = \text{Var}(Z|X = 1) \cdot P(X = 1) + \text{Var}(Z|X = 0) \cdot P(X = 0) \quad (9)$$

$$+ (E(Z|X = 1) - E(Z))^2 \cdot P(X = 1) + (E(Z|X = 0) - E(Z))^2 \cdot P(X = 0) \quad (10)$$

We used a multigroup structural equation model with stochastic group sizes to estimate the $X = x$ -conditional expectations and variances, and the probabilities $P(X = 0)$ and $P(X = 1)$.

$$\bar{z}_0 = \hat{E}(Z|X = 0) \quad (11)$$

$$\bar{z}_1 = \hat{E}(Z|X = 1) \quad (12)$$

$$s_0^2 = \widehat{\text{Var}}(Z|X = 0) \quad (13)$$

$$s_1^2 = \widehat{\text{Var}}(Z|X = 1) \quad (14)$$

$$p_0 = \hat{P}(X = 0) \quad (15)$$

$$p_1 = \hat{P}(X = 1) \quad (16)$$

The sample mean \bar{z} and variance s^2 were then obtained by

$$\bar{z} = \bar{z}_0 \cdot p_0 + \bar{z}_1 \cdot p_1 \quad (17)$$

$$s^2 = s_0^2 \cdot p_0 + s_1^2 \cdot p_1 + (\bar{z}_0 - \bar{z})^2 \cdot p_0 + (\bar{z}_1 - \bar{z})^2 \cdot p_1 \quad (18)$$

This approach allows us to model the group sizes as random variables. Thus, variability in sample moment estimates due to different group weights are accounted for. For more information on the underlying model and assumptions, see Mayer, Dietzfelbinger, Rosseel, and Steyer (2016, p. 4).

In the second step, we used the estimated sample moments \bar{z} and s^2 to compute the estimated parameters of the distribution. Depending on the distribution this was

$$\text{Normal } \mathcal{N}(\mu, \sigma): \quad \hat{\mu} = \bar{z}, \quad \hat{\sigma}^2 = s^2$$

$$\text{Poisson } \mathcal{P}(\lambda): \quad \hat{\lambda} = \bar{z}$$

$$\text{Uniform } U(a, b) \quad \hat{a} = \bar{z} - \sqrt{3s^2}, \quad \hat{b} = \bar{z} + \sqrt{3s^2}$$

For the three distributions in our simulation study, this procedure provided unbiased and efficient estimates for the parameters of the distribution. We especially compared the two-step procedure to conventional maximum likelihood estimates. For the uniformly distributed Z , the MLE provided biased estimates in small samples, while the two step-procedure was unbiased. For normally or Poisson distributed Z , point estimates from MLE and the two-step procedure were identical, however, the MLE was a little less efficient.

Data generation function

```

1 data.sim <- function(N, cohensd, b101, distribution){
2   # Standardized Regression Coefficients in Control Group
3   b000 <- 0    # Intercept
4   b001 <- 0.5  # Slope
5
6   # Remaining parameters are set depending on the covariate's distribution
7   if (distribution == "normal"){
8     muZ <- 0    # Mean of Z
9     sigmaZ <- 1 # Variance of Z
10    a000 <- b000 - b001*muZ/sqrt(sigmaZ) # Unstandardized Intercept in Control Group
11    a001 <- b001/sqrt(sigmaZ) # Unstandardized Slope in Control Group
12    a101 <- b101/sqrt(sigmaZ) # Unstandardized Interaction Coefficient
13    varY0 <- 1.089 # Var(Y|X=0)
14    a100 <- log(sqrt(varY0)*cohensd + exp(.5*b001^2)) - .5*(b001 + b101)^2 # Unstandardized
      Coefficient of Treatment Variable
15    z <- rnorm(N, muZ, sigmaZ) # Z from a Standard Normal Distribution
16  } else if (distribution == "poisson"){
17    lambdaZ = 1 # Mean and Variance of Z
18    varY0 <- 1.06 # Var(Y|X=0)
19    a000 <- b000 - b001*lambdaZ/sqrt(lambdaZ) # Unstandardized Intercept in Control Group
20    a001 <- b001/sqrt(lambdaZ) # Unstandardized Slope in Control Group
21    a101 <- b101/sqrt(lambdaZ) # Unstandardized Interaction Coefficient
22    a100 <- log(sqrt(varY0)*cohensd + exp(a000)*exp(lambdaZ*(exp(a001) - 1))) -
      log(exp(a000)*exp(lambdaZ*(exp(a001+a101)-1))) # Unstandardized Coefficient of
      Treatment Variable
23    z <- rpois(N, lambdaZ) # Z from a Poisson distribution
24  } else if (distribution == "uniform"){
25    aZ <- -sqrt(3) # lower limit of Z
26    bZ <- sqrt(3) # upper limit of Z
27    muZ <- (sqrt(3)-sqrt(3))/2 # Mean of Z
28    sigmaZ <- sqrt((sqrt(3)+sqrt(3))^2/12) # Variance of Z
29    a000 <- b000 - b001*muZ/sqrt(sigmaZ) # Unstandardized Intercept in Control Group
30    a001 <- b001/sqrt(sigmaZ) # Unstandardized Slope in Control Group
31    a101 <- b101/sqrt(sigmaZ) # Unstandardized Interaction Coefficient
32    varY0 <- 1.09 # Var(Y|X=0)
33    a100 <- log(sqrt(varY0)*cohensd + exp(b000)*(exp(b001*bZ) -
      exp(b001*aZ))/(b001*(bZ-aZ))) -
34    log(exp(b000)*(exp((b001+b101)*bZ)-exp((b001+b101)*aZ))/((b001+b101)*(bZ-aZ))) #
      Unstandardized Coefficient of Treatment Variable
35    z <- runif(N, aZ, bZ) # Z from a Uniform distribution
36  }
37  zS <- scale(z) # z-transformation of Z for generation of X
38  x <- rbinom(N, 1, plogis(zS)) # Treatment with treatment probability as logistic function of Z
39  muY <- exp(a000 + a001*z + a100*x + a101*x*z) # Conditional Expectation of Y depending on Z
      and X
40  y <- rnbinom(N, mu = muY, size = 1) # Y from negative binomial distribution
41
42  d <- data.frame(y,x,z)
43  return(d)
44 }

```

Estimator functions

```

1  library(lavaan)
2  library(car)
3
4  ##### function estimating TEMZ #####
5  computeTemz <- function(m1, d){
6    # Get regression coefficients
7    coefs <- coef(m1)
8    b00 <- coefs[1]
9    b10 <- coefs[2]
10   b01 <- coefs[3]
11   b11 <- coefs[4]
12
13   # Get mean of covariate z
14   zM <- mean(d$z)
15
16   # Compute and return TEMZ
17   temz <- exp(b00+b10+b01*zM+b11*zM) - exp(b00 + b01*zM)
18   return(temz)
19 }
20
21
22 ##### function estimating ATE (empirical approach, fixed covariate) #####
23 computeAveEmp <- function(m1, d){
24
25   # Compute all conditional effects
26   newdata0 <- newdata1 <- d
27   newdata0$x <- 0
28   newdata1$x <- 1
29
30   tau0 <- predict(m1, newdata=newdata0, type="response") ## tau0 is almost constant
31   tau1 <- predict(m1, newdata=newdata1, type="response") ## more variation in tau1
32   delta10 <- tau1 - tau0
33
34   # Compute VCOV of all conditional effects
35   vcovs <- vcov(m1)
36   mat <- cbind(delta10, tau1, d$z*delta10, d$z*tau1)
37   vcov_delta10 <- mat %*% vcovs %*% t(mat)
38
39   # Compute standard error for ATE
40   mat <- rbind(1/nrow(d),
41     as.numeric(d$x==0)/sum(d$x==0),
42     as.numeric(d$x==1)/sum(d$x))
43   vcov_eff <- mat %*% vcov_delta10 %*% t(mat)
44   se_Ave <- sqrt(vcov_eff[1,1])
45
46   # Return results
47   res <- list(Ave=mean(delta10),
48     se_Ave=se_Ave)
49
50   return(res)
51 }
52

```



```

53 ##### function estimating ATE (moment-based approach) #####
54 computeAveMom <- function(m1, d, distribution, se = "stochastic"){
55   # Get regression coefficients and their covariance
56   coefs <- coef(m1)
57   vcovs <- vcov(m1)
58
59   pnames <- c("g000", "g100", "g001", "g101")
60   names(coefs) <- row.names(vcovs) <- colnames(vcovs) <- pnames
61
62   # lavaan model for joint estimation of sample moments
63   modelz <- '
64     z ~ c(meanz001, meanz101)*1
65     z ~~ c(varz001, varz101)*z
66     group % c(groupw0, groupw1)*w
67     gw0 := groupw0
68     gw1 := groupw1
69     mz001 := meanz001
70     mz101 := meanz101
71     vz001 := varz001
72     vz101 := varz101
73     N := exp(gw0) + exp(gw1)
74     relfreq0 := exp(gw0)/N
75     relfreq1 := exp(gw1)/N
76     Ez1 := mz001*relfreq0 + mz101*relfreq1
77     Vz1 := vz001*relfreq0 + vz101*relfreq1 + relfreq0*(mz001-Ez1)^2 + relfreq1*(mz101-Ez1)^2
78     Px0 := relfreq0
79     Px1 := relfreq1
80     Pk0gx0 := relfreq0/Px0
81     Pk0gx1 := relfreq1/Px1
82     Ez1gx0 := mz001*Pk0gx0
83     Ez1gx1 := mz101*Pk0gx1
84     Vz1gx0 := vz001*Pk0gx0
85     Vz1gx1 := vz101*Pk0gx1
86   '
87   mz <- sem(modelz, data=d, group="x", group.label=c(0,1), group.w.free=TRUE)
88
89   ## augment coefs and vcovs
90   acoefs <- c(coefs, coef(mz, type="user")[-c(1:6)])
91
92   if (se == "stochastic"){
93     # for stochastic covariate
94     avcovs <- lav_matrix_bdiag(vcovs, lavInspect(mz, "vcov.def", add.class = FALSE))
95   } else if (se == "fixed"){
96     # for fixed covariate: covariance of distribution parameters fixed to zero
97     avcovs <- lav_matrix_bdiag(vcovs, lavInspect(mz, "vcov.def", add.class = FALSE)*0)
98   }
99
100  row.names(avcovs) <- colnames(avcovs) <- names(acoefs)
101
102  if (distribution == "normal"){
103    Eg1 <- deltaMethod(acoefs,
104      "exp(g000+g100)*exp((g001+g101)*Ez1+((g001+g101)^2*Vz1/2))-exp(g000)*exp(g001*Ez1+(g001^2*Vz1/2))",
105      avcovs, func="Eg1")
106  } else if (distribution == "uniform"){

```

```
105     Eg1 <- deltaMethod(acoefs,
106       "exp(g000+g100)*(exp((g001+g101)*(Ez1+sqrt(3*Vz1))))-exp((g001+g101)*(Ez1-sqrt(3*Vz1))))/((g001+g101)*(2*s
107     avcovs, func="Eg1")
108   } else if (distribution == "poisson"){
109     Eg1 <- deltaMethod(acoefs,
110       "exp(g000+g100)*exp(Ez1*(exp(g001+g101)-1))-exp(g000)*exp(Ez1*(exp(g001)-1))", avcovs,
111       func="Eg1")
112   } else if (distribution == "chisquare"){
113     Eg1 <- deltaMethod(acoefs,
114       "(exp(g000+g100)/(1-2*(g001+g101))^(Ez1/2))-exp(g000)/(1-2*g001)^(Ez1/2)", avcovs,
115       func="Eg1")
116   }
117 }
118
119 res <- list(Ave=Eg1$Estimate,
120 se__Ave=Eg1$SE)
121
122 return(res)
123 }
```

Simulation procedure

```

1 library(doParallel)
2 library(foreach)
3 library(MASS)
4 source("dataSimulation.R")
5 source("estimators.R")
6
7 # Set design parameters
8 sampleN <- c(50, 100, 250, 500, 1000)
9 interact <- c(0, 0.25, 0.5, 0.75, 1)
10 cohensd <- c(-.5, 0, .5)
11 distribution <- c("normal", "poisson", "uniform")
12 reps <- 5000
13 conditions <- expand.grid(sampleN, interact, cohensd, distribution)
14 colnames(conditions) <- c("sampleN", "interact", "cohensd", "distribution")
15
16 ### Simulation procedure
17 #setup parallel backend to use 4 processors
18 cluster <- makeCluster(4, outfile = "debug.txt")
19 registerDoParallel(cluster)
20 results <- list()
21 for (i in 1:nrow(conditions)){
22   # Parallel replications within each scenario
23   res <- foreach::foreach(icount(reps), .combine = rbind, .packages = c("car", "lavaan", "MASS"),
24     .errorhandling = "remove") %dopar% {
25     # Generate data set
26     d <- data.sim(conditions$sampleN[i],
27       conditions$cohensd[i],
28       conditions$interact[i],
29       conditions$distribution[i])
30     # Estimate NB-GLM
31     m1 <- glm.nb(y ~ x*z, data=d)
32     # Get results from ATE estimators
33     tmp <- c(unlist(computeAveEmp(m1, d, se = "fixed")),
34       unlist(computeAveAnaStoch(m1, d, distribution = conditions$distribution[i])),
35       unlist(computeAveAnaFixed(m1, d, distribution = conditions$distribution[i])),
36       unlist(computeAveAnaStoch(m1, d, distribution = "normal")),
37       computeTemz(m1, d),
38       m1$converged)
39     tmp
40   }
41   # Combine and save results
42   results <- data.frame(res,
43     conditions$sampleN[i],
44     conditions$cohensd[i],
45     conditions$interact[i],
46     conditions$distribution[i])
47   names(results) <- c("aveEmpFixed", "seEmpFixed",
48     "aveAnaStochastic", "seAnaStochastic",
49     "aveAnaFixed", "seAnaFixed",
50     "aveStochNorm", "seStochNorm",
51     "TEMZ",
52     "conv", "N", "cohensd", "interact", "distribution")

```

```
52   saveRDS(results, file = paste0("simulation/res_",
53     conditions$sampleN[i], "_",
54     conditions$cohensd[i], "_",
55     conditions$interact[i], "_",
56     conditions$distribution[i],
57     ".rds"))
58 }
59 stopCluster(cluster)
```

CountEffects package

As a supplement to this article, we provide an R package (R Core Team, 2018) called CountEffects, which comprises the functions related to the moment-based approach presented in this article.

It is available from GitHub and can be installed as follows:

```
1 library(devtools)
2 install_github('chkiefer/CountEffects')
```

References

- Glass, G. V. (1976). Primary, secondary, and meta-analysis of research. *Educational Researcher*, 5(10), 3-8. doi: 10.2307/1174772
- Mayer, A., Dietzfelbinger, L., Rosseel, Y., & Steyer, R. (2016). The EffectLiteR approach for analyzing average and conditional effects. *Multivariate Behavioral Research*, 51, 374–391. doi: <http://dx.doi.org/10.1080/00273171.2016.1151334>
- R Core Team. (2018). R: A language and environment for statistical computing [Computer software manual]. Vienna, Austria. Retrieved from <http://www.R-project.org/>