Symposium on Data Science and Statistics (SDSS) 2019, Bellevue, WA

## An R Package for Linear Mediation Analysis with Complex Survey Data



Yujiao Mai, Jiahui Xu, Deo Kumar Srivastava, Hui Zhang Department of Biostatistics

5/31/2019



## Outline

- 1. Classic Mediation Analysis
- 2. Complex Surveys
- 3. Problem & Solution
- 4. Software & Application

## Section 1 Classic Mediation Analysis



### 



• Estimator:  $F_{ML}(\theta; \mathbf{S}_n) = log |\mathbf{\Sigma}(\theta)| + trace [\mathbf{S}_n \mathbf{\Sigma}(\theta)^{-1}] - log |\mathbf{S}_n| - p$ 

 $S_n$ : sample covariance matrix

 $\Sigma(\theta)$ : estimated covariance matrix,  $\theta = (\mu_M, \mu_Y, \alpha, \gamma, \beta, \sigma_{\varepsilon_M}^2, \sigma_{\varepsilon_Y}^2)$ , p: number of parameters VAR $(M; \theta) = VAR(\mu_M + \alpha X) + \sigma_{\varepsilon_M}^2$ 

## Model (Continue)

• Significance Test:  $H_0: \alpha\beta = 0; H_1: \alpha\beta \neq 0$ 

1) Sobel's Test (Sobel, 1982)

 $T_{Sobel} = \hat{\alpha}\hat{\beta} / \sqrt{\hat{\beta}^2 \hat{Var}(\hat{\alpha}) + \hat{\alpha}^2 \hat{Var}(\hat{\beta})} \quad \stackrel{H_0}{\sim} N(0,1) \text{ as } n \to \infty$ 

2) Resampling/bootstrap

### **Assumptions**

True Model Independent Residuals Normality Independent-identically-distributed sample data

## Section 2 Complex Surveys

## 

7

## Multi-stage Sampling (Wolter, 2007)

Suppose the population (40,000 students) locate in two states.

Each state has 20 districts.

Each district has 10 schools.

Each school has 100 students.

Sample weights  $(\omega_0) = the$ number of students represented

State	District	School	Student	Y	$\omega_0$
1	1	1	1	9.8	2500
1	1	1	2	7.5	2500
1	1	2	3	8.3	2500
1	1	2	4	4.5	2500
1	2	3	5	4.5	2500
1	2	3	6	5.1	2500
1	2	4	7	2.3	2500
1	2	4	8	6.5	2500
 2	3	5	9	8.1	2500
2	3	5	10	3.2	2500
2	3	6	11	4.5	2500
2	3	6	12	5.8	2500
2	4	7	13	6.6	2500
2	4	7	14	8.1	2500
2	4	8	15	1.2	2500
2	4	8	16	6.5	2500

$$\omega_{0} = \frac{1}{\pi_{\text{PSU}} \times \pi_{\text{sch}|\text{PSU}} \times \pi_{\text{stud}|\text{sch}}} = \frac{1}{\frac{2}{20} \times \frac{2}{10} \times \frac{2}{100}} = 2,500$$

## **Problem: Estimate the Mean of** *Y*

- Disaggregated estimates vs. aggregated estimates (Group-specific effects) (Generalized effects)
- For aggregated estimates:
- Point estimate is consistent when including the sample weights
- Standard errors is underestimated even when taking into account the sample weights

### Adjustments:

Taylor series linearization (TSL)

Bootstrap

Jackknife repeated replications (JRR)

Balanced repeated replications (BRR)

## **Balanced repeated replications**

(Wolter, 2007) (Fay & Train, 1995)

State	District	School	Student	Y	$\omega_0$	$oldsymbol{\omega}_1'$	$\omega_2'$	$\boldsymbol{\omega}_3'$	$oldsymbol{\omega}_4'$
1	1	1	1	9.8	2500	0	2×2500	0	2×2500
1	1	1	2	7.5	2500	0	2×2500	0	2×2500
1	1	2	3	8.3	2500	0	2×2500	0	2×2500
1	1	2	4	4.5	2500	0	2×2500	0	2×2500
1	2	3	5	4.5	2500	2×2500	0	2×2500	0
1	2	3	6	5.1	2500	2×2500	0	2×2500	0
1	2	4	7	2.3	2500	2×2500	0	2×2500	0
1	2	4	8	6.5	2500	2×2500	0	2×2500	0
2	3	5	9	8.1	2500	0	0	2×2500	2×2500
2	3	5	10	3.2	2500	0	0	2×2500	2×2500
2	3	6	11	4.5	2500	0	0	2×2500	2×2500
2	3	6	12	5.8	2500	0	0	2×2500	2×2500
2	4	7	13	6.6	2500	2×2500	2×2500	0	0
2	4	7	14	8.1	2500	2×2500	2×2500	0	0
2	4	8	15	1.2	2500	2×2500	2×2500	0	0
2	4	8	16	6.5	2500	2×2500	2×2500	0	0

Note. The Worken The sampling samitpling built in the sound is District.

Replicate sampling weights  $\boldsymbol{\omega}'_r$ , r = 1, 2, ..., R

*R* is the number of replications, R = 4 in the case.

$$SE_{BRR}(\hat{\mu}) = \sqrt{\frac{\sum_{r=1}^{R} (\hat{\mu}_r - \hat{\mu})^2}{R}}$$

 $\hat{\mu}_r$  is the estimate using replicate weight  $w'_r$ .

 $\hat{\mu}$  is the estimate using original(main) sample weights

# Section 3 Problem & Solution



## **Problem to Solve**

### How can mediation analysis work with complex surveys?

<b>Classic Mediation Analysis</b>	Complex Survey	/S
i-i-d sample: Independent	Within cluster: Dependent	Disaggregated
i-i-d sample: Identically distributed	Between cluster: Heterogenous	Estimates
i-i-d sample: Equal possibility	Unequal-sized strata: Unequal possibility	Aggregated Estimates

## **Potential Solutions**

### Aggregated Estimates

Design-based method:

- Taylor series linearization (TSL)
- Bootstrap
- Jackknife repeated replications (JRR)
- Balanced repeated replications (BRR)

## Most national/international surveys do not provide the cluster indicator.

## Apply Balanced Repeated Replications to Analysis of $\alpha\beta$ (Mai et al., 2019)

### Estimator

Maximum Likelihood  $F_{ML}(\theta; \mathbf{S}_n) = log |\boldsymbol{\Sigma}(\theta)| + trace [\mathbf{S}_n \boldsymbol{\Sigma}(\theta)^{-1}] - log |\mathbf{S}_n| - p$ 

Weighted covariance matrix

Replace  $S_n$  with weighted sample covariance matrix  $S_{wn}$ 

### Test Statistic

Asymptotical Normality (Bishop, 1975; Rao, 1973)  $\hat{\alpha}\hat{\beta} \sim N(\mu, \sigma^2)$  as  $n \to \infty$ 

**Use BRR Standard Errors** 

$$T_{\rm BRR} = \frac{\hat{\alpha}\hat{\beta}}{SE_{\rm BRR}(\hat{\alpha}\hat{\beta})} \xrightarrow{H_0} N(0,1) \text{ as } n \to \infty$$

$$SE_{BRR}(\hat{\alpha}\hat{\beta}) = \sqrt{\frac{1}{R(1-f)^2}} \sum_{r=1}^{R} (\hat{\alpha}_r \hat{\beta}_r - \hat{\alpha}_0 \hat{\beta}_0)^2$$

## Section 4 **Software & Application**



## **Software Packages**

• R package 'MedSurvey' (Feb. 2019)

Flexible and complex models

https://CRAN.R-project.org/package=MedSurvey



### Application

#### Data:

2014-15 CPS Tobacco Use Supplement (TUS; U.S. Department of Commerce and U.S. Census Bureau 2016), employed adult daily smokers (Non-Hispanic White males only).

#### Survey Design:

Balanced Repeated Replications R = 160. f = 0.5 is suggested.



```
## Package and options
library("MedSurvey")
## Data and related information
MedData
R < -160
wgtnames <- paste("repwgt", seq(0,R,by=1), sep="")</pre>
mwgtname=wgtnames[1]
repwgtnames=wgtnames[2:(R+1)]
## Sepcify the model 1
model1 < -
  numcg ~ u0*1 + gamma0*workban + b1*sp_adltban + b2*sp_kidsban
  sp_adltban ~ u1*1 + a1*workban
  sp_kidsban ~ u2*1 + a2*workban
  sp_adltban ~~ sp_kidsban
 a1b1 := a1*b1
 a2b2 := a2*b2
 total := gamma0 + (a1*b1) + (a2*b2)
```

## View the summary results of the mediation analysis
med.summary(fit=fit.BRR, med.eff=c('a1b1', 'a2b2'))

## Results

Table 1. Tests for Mediation Effects between Smoking Ban at Work and Number of CigarettesSmoked per Day among Male NH White Daily Smokers

Mediator		Not We	ighted		BRR Weighted					
	$\hat{lpha}\hat{eta}$	ŜD	<i>p</i> -value	Adjusted <i>p</i> -value	$\hat{lpha}\hat{eta}$	$\widehat{SD}_{\mathrm{BRR}}$	<i>p</i> -value	Adjusted <i>p</i> -value		
<i>M</i> <sub>1</sub>	-0.0170	0.0066	0.010	0.020	-0.0156	0.0073	0.032	0.063		
<i>M</i> <sub>2</sub>	0.0018	0.0033	0.588	0.612	0.0004	0.0045	0.937	1.000		

*Note. M*1 = Supporting Smoking Ban in Adult-exclusive Areas.

*M*2 = Supporting Smoking Ban in Children's Areas.

n = 2,260 (Estimated population size is 3,294,568).

Adjusted p values adjusted for the multiple tests using Holm's method (Holm, 1979).

## **Software Packages**

• R package 'MedSurvey' (Feb. 2019)

Flexible and complex models

https://CRAN.R-project.org/package=MedSurvey

### • R shiny (March 2019)

User-friendly interface

https://sjbiostat.shinyapps.io/MedSurvey/



### **R** shiny

MedSurvey X	+							_		
+ → C	t.shinyapps.io/MedSurvey/								<b>e</b> :	
Model Fitting E	xample Manual User Map Contact									
A data file with header		Summary	Downloa	ad	-4					
Browse MedData.csv		Multimediation	with Comp	liex Survey Da	ata:					
Uplo	bad complete	ſ	Effect	Estimate	BRR SE.	P-Value	Adjusted P-Value			
Y (Outcome)		sp_adltban	a1b1	-0.015456	0.004574	0.000728	0.001456			
numcg		sp_kidsban	a2b2	-0.005366	0.003852	0.1636	0.1636			
workban         M (Mediators)         sp_adltban sp_kidsban         Zm (Covariates for mediator)	• • • • • • • • • • • • • • • • • • •	p Value adju: Standard err	stment me ors type is	ethod is holm.					.ook,	this's our produc

## **Software Packages**

• R package 'MedSurvey' (Feb. 2019) Flexible and complex models

https://CRAN.R-project.org/package=MedSurvey

- R shiny (March 2019)
   User-friendly interface
   <a href="https://sjbiostat.shinyapps.io/MedSurvey/">https://sjbiostat.shinyapps.io/MedSurvey/</a>
- SAS macros 'MedBRR' (Dec. 2018)

Super large-scale datasets

https://github.com/YujiaoMai/MedSurvey



## References

- Bishop, Y. M., Fienberg, S. E., & Holland, P. W. (1975). *Discrete multivariate analysis: theory and practice*. Cambridge, Massachusetts: Cambridge, Mass., MIT Press.
- Fay, R. E., & Train, G. F. (1995). Aspects of survey and model-based postcensal estimation of income and poverty characteristics for states and counties. In *Proceedings of the Section on Government Statistics, American Statistical Association, Alexandria, VA* (pp. 154-159).
- Judkins, D. R. (1990). Fay's method for variance estimation. Journal of Official Statistics, 6(3), 223-239.
- Mai, Y., Ha, T., & Soulakova, J. N. (2019). Multimediation Method With Balanced Repeated Replications For Analysis Of Complex Surveys. *Structural Equation Modeling: A Multidisciplinary Journal*.
- Rao, C. R. (1973). Linear statistical inference and its applications (2nd ed., Vol. 2). New York, NY: John Wiley & Sons, Inc.
- Sobel, M. E. (1982). Asymptotic confidence intervals for indirect effects in structural equation models. Sociological Methodology, 13, 290–312.
- U.S. Department of Commerce, & U.S. Census Bureau. (2016). National Cancer Institute and Food and Drug Administration co-sponsored Tobacco Use Supplement to the Current Population Survey. 2014-15.
- Wolter, K. (2007). Introduction to variance estimation. New York, NY: Springer.

## Acknowledgements

 The research is sponsored by American Lebanese Syrian Associated Charities (ALSAC).





Questions?

Thank you !