Keio University



The Problem of Treating Imputed Data as Observed Data When We Estimate the Effect of Exposure to **Particulate Matter**

> **Tomoshige Nakamura*** Mihoko Minami

BOSTON UNIVERSITY/KEIO UNIVERSITY WORKSHOP 2016 Probability and Statistics

Boston University — August 15-19, 2016

our health, when the number of monitoring stations is limited

Particulate Matter

- Particulate Matter is Complex mixture of extremely small particles and liquid droplets, and are widely studied and concerned the relationship with various diseases.
- Estimating the effect of Particulate Matter exposure to our health is one of the active research topic of environmental epidemiology.

•





 When we estimate the effect of exposure, we often use the community health survey data, and need the information of exposure at survey areas.

Left figure shows the survey conducted at some area of Japan, we can see the number of monitoring stations in survey areas are limited.



BU Workshop 2016 @ Boston



· In many of researches for particulate matter exposure, missing data problem are not paid enough attention, and simple regression model are used to fill the missing of exposure • Then, the effect of exposure is estimated as if they were observed

Notation!

• The another reason why regression imputation is used is that "the analyst for exposure assessment and the analyst for estimating the effect exposure are not same one".

 \cdot So, the unified bayesian approach is not the realistic choice to analysis.







Procedure using for estimating the effect of exposure in common research.

The procedure of estimating the effect of exposure using regression imputation can be decompose to 3 steps.



BU Workshop 2016 @ Boston



Step1 (constructiong the model for exposure) : fitting linear regression model to the data observed at monitoring stations, and construct the prediction model for exposure.



BU Workshop 2016 @ Boston





Procedure using for estimating the effect of exposure in common research.

Step2 (filling the missing) : using the model constructed at first step and covariates, we predict the exposure of sub-survey area which cannot be observed exposure. • Then fill the missing by predicted values of model.



BU Workshop 2016 @ Boston





Procedure using for estimating the effect of exposure in common research.

imputed data, and estimate the effect of exposure



BU Workshop 2016 @ Boston

•

Step3 (Estimating the effect) Then poisson regression model is fitted to regression





The aim of this research

From the missing data analysis context

may be violated, the variance of estimator will be underestimated !!

Purpose

- when we use regression imputation in practical analysis.
- problem caused by the limited number of monitoring stations.

 \cdot If we treat the imputed values as if they are observed, the consistency of estimator

A. We try to <u>clarify the problem of using regression imputation</u> from analytical and practical points of view, and organize the important points of matter

B. To develop the method to estimate the effect of exposure, avoiding the





Today I'm going to talk...

Focus on…

when we use regression imputation in practical analysis.

BU Workshop 2016 @ Boston

A. We try to <u>clarify the problem of using regression imputation</u> from analytical and practical points of view, and organize the important points of matter

> At first, let us consider the problem of regression imputation from analytical point of view under the simple setting.





Analytical point of view - Setting

Setting we consider the case of using linear regression model for imputation and poisson regression model to estimate the effect …

· This setting, can be formulated as the problem of estimating the regression coefficients β by fitting model-(2) to data generated by Model-(1).

- (1) Generating scheme of data $Y_i | X_i \sim \text{Possion}(\lambda_i)$ $\log(\lambda_i) = \beta_0^* + X_i \beta_1^*$ $X_i | \mu_i \sim \text{Normal}(\mu_i, \sigma^2)$
 - μ_i : mean of exposure



under this setting, consider the properties of estimators



Analytical point of view - Consistency of regression coefficient estimator.

The estimator for $\beta = (\beta_0, \beta_1)^T$ has following property





- $\hat{\beta}_0$ is not a consistent estimator $\hat{\beta}_1$ is consistent estimator

We are interested in only, $\hat{\beta}_1$ so the inconsistency of intercept is not a significant problem

Boston-Keio Workshop 2016 @ Boston Univ.







Analytical point of view - Result

- If we know the true mean of exposure, the estimator for effect of exposure will be consistent.
- When we don't know the true mean of exposure, the estimator of exposure will be inconsistent
- In both cases, asymptotic variance will be smaller than when we treat imputed values properly.



What is the practitioner want to know is

- Now we showed, the inference based on regression imputation method is invalid.
- However, what is the practitioner want to know is...

How much the bias will be occur (A) when we use the regression imputation ?

(B) <u>When the large bias will be occur ?</u>



BU Workshop 2016 @ Boston



we perform the simulation of regression imputation under previous settings



Practical point of view - simulation setting

- variance.
- Procedure of Simulation
- 1.Generate 160 size data containing 70% missing exposure. 3.Impute the missing values of exposure by linear predictor. 4. Fit poisson regression model to regression imputed data, and estimate the coefficients. 5. iterate 1-3 procedure, <u>1000</u> times



· We visualize the inconsistency of estimator and underestimation of asymptotic

2. Construct the Linear regression model for imputation using data obtained at monitoring station.





Practical point of view - simulation setting (procedure)



Boston-Keio Workshop 2016 @ Boston Univ.

Result - Fitting Poisson Regression Model to Regression Imputed Data.

- Compute 95% Confidence interval for coefficient of poisson regression model, 1000 times. (sort by coef)
- Black Solid Line : Estimated coefficients
- <u>Red Band : 95%Cls based on Fisher Information</u>
 (85,2% contain the true value)
- <u>Green Band : 95%Cls based on Sandwich Estimator</u>
 <u>(88.6% contain the true value)</u>
- Mean of Estimated value does not consistent to true value. 95% Cls based on fisher information and Sandwich Estimator is shorter but not too much.



Result - analytical & practical point of view (simple setting).

By Analysis

Asymptotic variance is underestimated

By Simulation

 If we can <u>specify the mean model</u> for exposure <u>properly</u> Average of estimates is approximately equal to the true parameter value. • **Underestimation** of asymp. variance matters little

Result Using the regression imputation is **not appropriate in analytically**, but is not unacceptable in practice

BU Workshop 2016 @ Boston

• The estimator of the effect of exposure is **inconsistent**



Result - analytical & practical point of view.



Using the regression imputation is **not appropriate in analytically**, but is not unacceptable in practice

BU Workshop 2016 @ Boston



• The estimator of the effect of exposure is **inconsistent**

If we can specify the mean model for exposure properly

Average of estimates is approximately equal to the true parameter value.

• **Underestimation** of asymp. variance matters little, when R > 0.6





Simulation (practical settings) - Procedure

Let me consider the data generated from following model.

Model for Outcome

$$Y_i | X_i \sim \text{Poisson}(\lambda_i)$$

 $\log \lambda_i = \beta_0 + x_i \beta_1 + \sum_{k=2}^4 z_{ik} \beta_k$

Model for Exposure

$$X_i | \mu_i \sim \text{Normal}(\mu_i, \sigma^2)$$
$$\mu_i = \alpha_0 + \exp\left(\sum_{j=1}^4 w_{ij}\alpha_j + \gamma_{i0} + \sum_{j=1}^4 w_{ij}\gamma_{ij}\right)$$

where, α_j are fixed effect, and γ_{ij} are the random effect at each site, and ϕ_i is spatial component.

BU Workshop 2016 @ Boston





Longitude







Simulation (practical settings) - Procedure



Heatmap of Amount of Particulate Matter Exposure

BU Workshop 2016 @ Boston

- · In general, monitoring stations are not established with low particulate matter concentration.
- · At areas surrounded by black solid line, we assume that we cannot observe the exposure information.

 Then, we use regression imputation and estimate the effect of exposure β (same as previous simulation).

Iterate these procedure 1000 times.







Simulation (practical settings) - Result

Compute 95% Confidence interval for coefficient of poisson regression model, 1000 times. (sort by coef)

- <u>Red Band : 95%Cls based on Fisher Information (25,9%</u> <u>contain the true value)</u>
- Green Band : 95%Cls based on Sandwich Estimator <u>(29.7% contain the true value)</u>

Black Line : Estimated coefficients

 This figure shows, the consistency of estimator is violated, highly biased, and when there is a spatial correlation, the estimated value of exposure will be smaller than true value.

BU Workshop 2016 @ Boston







<u>Simple Settings</u>

<u>Realistic Settings</u>

In real, for harmful effect, the negative impact of misjudgment



BU Workshop 2016 @ Boston

- using regression imputation to fill the missing does not cause the serious problem
- if we ignore the effect of spatial structure and random effect of each sites, the effect of exposure will be underestimated!!



Conclusion

problem, so it is unacceptable in practice !!!!

Ongoing work

We now show the underestimation of the effect only by simulation, so next we show them by analytically

we develop the method alternative to regression imputation that can make robust inference under the mild model misspecification for exposure.

Using regression imputation to estimate the effect of exposure causes very serious





References - 1

- 124(1), 23-29.
- Acute Myocardial Infarction." Environ Health Perspect 121:192-196;
- Medicine, Vol. 192, No. 3, pp. 337-341.
- mortality using different air pollution exposure models: impacts in rural and urban California", International Journal of Environmental Health Research, Vol. 26, Iss. 2
- · U.S. EPA.(2009) ``2009 Final Report: Integrated Science Assessment for Particulate Matter." U.S. Environmental Protection Agency, Washington, DC, EPA/600/R-08/139F.
- Rubin, D. B. (1976). ``Inference and missing data''. Biometrika, 63, 581–592.
- Pollution and Mortality in a Cohort of More than a Million Adults in Rome." Environ Health Perspect 121:324-331; BU Workshop 2016 @ Boston

· Kioumourtzoglou, M.A., Schwartz, J. D., Weisskopf, M. G., Melly, S. J., Wang, Y., Dominici, F., Zanobetti, A. (2016). ``Long-term PM2.5 Exposure and Neurological Hospital Admissions in the Northeastern United States." Environmental Health Perspectives,

• Madrigano J, Kloog I, Goldberg R, Coull BA, Mittleman MA, Schwartz J. (2013). ``Long-term Exposure to PM2.5 and Incidence of

· Yongping Hao, Lina Balluz, Heather Strosnider, Xiao Jun Wen, Chaoyang Li, and Judith R. Qualters (2015). ``Ozone, Fine Particulate Matter, and Chronic Lower Respiratory Disease Mortality in the United States", American Journal of Respiratory and Critical Care

· Cynthia A. Garcia, Poh-Sin Yap, Hye-Youn Park, Barbara L. Weller (2016)., ``Association of long-term PM2.5 exposure with

· Cesaroni G, Badaloni C, Gariazzo C, Stafoggia M, Sozzi R, Davoli M, Forastiere F. (2013). ``Long-Term Exposure to Urban Air

references - 2

- Correlation". In JSM Proceedings, Statistical Computing Section. Alexandria, VA: American Statistical Association. 2509-2514.
- Fifth Berkeley Symp. on Math. Statist. and Prob., Vol. 1 (Univ. of Calif. Press, 1967), 221-233
- Rubin, D. B. (1987). Multiple imputation for nonresponse in surveys. New York: Wiley.
- Tsiatis, A. A. (2006). Semiparametric theory and missing data. New York: Springer.
- Vaart, A. W. (1998). Asymptotic statistic. Cambridge: Cambridge University Press.

• Nakamura, T. Minami, M.(2015). Prediction of PM10 and PM2.5 Concentration Using Land Use Data and Spatial

· Peter J. Huber. (1967), The behavior of maximum likelihood estimates under nonstandard conditions," in Proc.

· Wood, S. N. (2006), Generalized additive models: An introduction with R. Boca Raton, FL: Chapman & Hall/CRC.





