

MR4571186 62G05 62D05 62G20

Che, Menglu (1-YALE-DBH); Han, Peisong (1-MI-BPH);  
Lawless, Jerald F. (3-WTRL-S)

**Improving estimation efficiency for two-phase, outcome-dependent sampling studies. (English. English summary)**

*Electron. J. Stat.* **17** (2023), no. 1, 1043–1073.

When some covariates are expensive or difficult to measure, two-phase outcome dependent sampling (ODS) is often employed. To make the argument more concrete, denote the outcome, inexpensive-to-measure covariates, and expensive-to-measure covariates by  $Y$ ,  $\mathbf{X}$ , and  $\mathbf{Z}$ , respectively. In two-phase ODS,  $(Y, \mathbf{X})$  are taken from all individuals in Phase 1, and then  $\mathbf{Z}$  is taken from selected individuals in Phase 2.

In this article, the authors study the problem of estimating the parameter  $\beta$  in a parametric model of the conditional density  $f(Y|\mathbf{X}, \mathbf{Z}; \beta)$  under two-phase ODS. The conditional maximum likelihood (CML) method is popularly chosen, because it does not need to model  $f(\mathbf{Z}|\mathbf{X})$  unlike the full maximum likelihood method and it can handle zero selection probabilities in Phase 2. However, CML leads to efficient loss due to discarding the information on those who do not enter Phase 2.

To pursue efficiency gain over CML, the authors summarize the information on Phase 1 through a parametric model of the conditional density  $f(Y|\mathbf{X}; \theta)$ . More specifically, they investigate the maximum empirical likelihood (MEL) estimation of  $\beta$ , where the information is explicitly incorporated as some zero expectation restriction implied by this model, as well as the zero score restriction from CML.

Their analysis starts from the case with strictly positive selection probabilities in Phase 2 for all individuals. Under the oracle scenario where the selection probabilities are known, the MEL estimator for  $\beta$  is shown to attain efficiency gain over CML in the sense that the asymptotic variance of the former is smaller than the one for the latter. The analysis is further extended to the cases where positive selection probabilities are parametrically modeled (and estimated) and zero selection probabilities are allowed. Efficiency gain of MEL over CML in finite samples is confirmed via Monte Carlo simulations, and a real data analysis from the National Health and Nutrition Examination Survey is also conducted.

*Masayuki Hirukawa*

#### [References]

1. BAKSI, A. J., TREIBEL, T. A., DAVIES, J. E., HADJILOIZOU, N., FOALE, R. A., PARKER, K. H., FRANCIS, D. P., MAYET, J. and HUGHES, A. D. (2009). A meta-analysis of the mechanism of blood pressure change with aging. *Journal of the American College of Cardiology* **54** 2087–2092.
2. BARNETT, I. J., LEE, S. and LIN, X. (2013). Detecting rare variant effects using extreme phenotype sampling in sequencing association studies. *Genetic Epidemiology* **37** 142–151.
3. BJØRNLAND, T., BYE, A., RYENG, E., WISLØFF, U. and LANGAAS, M. (2018). Powerful extreme phenotype sampling designs and score tests for genetic association studies. *Statistics in Medicine* **37** 4234–4251. MR3879425 MR3879425
4. BRESLOW, N. and CAIN, K. (1988). Logistic regression for two-stage case-control data. *Biometrika* **75** 11–20. MR0932812 MR0932812
5. BRESLOW, N. E. and HOLUBKOV, R. (1997). Maximum likelihood estimation of logistic regression parameters under two-phase, outcome-dependent sampling. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **59** 447–461. MR1440590 MR1440590
6. BRESLOW, N. E. and HOLUBKOV, R. (1997). Weighted likelihood, pseudolikelihood

- and maximum likelihood methods for logistic regression analysis of two-stage data. *Statistics in Medicine* **16** 103–116.
7. CHATTERJEE, N., CHEN, Y.-H. and BRESLOW, N. E. (2003). A pseudoscore estimator for regression problems with two-phase sampling. *Journal of the American Statistical Association* **98** 158–168. MR1965682 MR1965682
  8. CHATTERJEE, N., CHEN, Y.-H., MAAS, P. and CARROLL, R. J. (2016). Constrained maximum likelihood estimation for model calibration using summary-level information from external big data sources. *Journal of the American Statistical Association* **111** 107–117. MR3494641 MR3494641
  9. CHE, M., LAWLESS, J. F. and HAN, P. (2020). Empirical and Conditional Likelihoods for Two-Phase Studies. *Canadian Journal of Statistics* [doi.org/10.1002/cjs.11566](https://doi.org/10.1002/cjs.11566). MR4267924 MR4267924
  10. DERKACH, A., LAWLESS, J. F. and SUN, L. (2015). Score tests for association under response-dependent sampling designs for expensive covariates. *Biometrika* **102** 988–994. MR3431569 MR3431569
  11. ESPIN-GARCIA, O., CRAIU, R. V. and BULL, S. B. (2018). Two-phase designs for joint quantitative-trait-dependent and genotype-dependent sampling in post-GWAS regional sequencing. *Genetic epidemiology* **42** 104–116.
  12. HAN, P. and LAWLESS, J. F. (2016). Comment on “Constrained Maximum Likelihood Estimation for Model Calibration Using Summary-level Information from External Big Data Source”. *Journal of the American Statistical Association* **111** 118–121. MR3494642 MR3494642
  13. HAN, P. and LAWLESS, J. F. (2019). Empirical likelihood estimation using auxiliary summary information with different covariate distributions. *Statistica Sinica* **29** 1321–1342. MR3932520 MR3932520
  14. HAN, P., TAYLOR, J. M. G. and MUKHERJEE, B. (2022). Integrating risk prediction models with no model details. *Canadian Journal of Statistics* **00** To appear. MR4595233
  15. HERMANSEN, K. (2000). Diet, blood pressure and hypertension. *British Journal of Nutrition* **83** S113–S119.
  16. HORVITZ, D. G. and THOMPSON, D. J. (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association* **47** 663–685. MR0053460 MR0053460
  17. HUANG, C.-Y., QIN, J. and TSAI, H.-T. (2016). Efficient estimation of the Cox model with auxiliary subgroup survival information. *Journal of the American Statistical Association* **111** 787–799. MR3538705 MR3538705
  18. IMBENS, G. W. (2002). Generalized method of moments and empirical likelihood. *Journal of Business & Economic Statistics* **20** 493–506. MR1973800 MR1973800
  19. KEOGH, R. H. and COX, D. R. (2014). *Case-control studies* **4**. Cambridge University Press, Cambridge, UK. MR3443808 MR3443808
  20. KIM, S., CAI, J. and LU, W. (2013). More efficient estimators for case-cohort studies. *Biometrika* **100** 695–708. MR3094446 MR3094446
  21. KULICH, M. and LIN, D. (2004). Improving the efficiency of relative-risk estimation in case-cohort studies. *Journal of the American Statistical Association* **99** 832–844. MR2090916 MR2090916
  22. LAWLESS, J., KALBFLEISCH, J. and WILD, C. (1999). Semiparametric methods for response-selective and missing data problems in regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **61** 413–438. MR1680310 MR1680310
  23. LIN, D.-Y., ZENG, D. and TANG, Z.-Z. (2013). Quantitative trait analysis in sequencing studies under trait-dependent sampling. *Proceedings of the National*

- Academy of Sciences of the United States of America* **110** 12247–12252. MR3105371 MR3105371
24. LITTLE, R. J. and RUBIN, D. B. (2019). *Statistical analysis with missing data* **793**. John Wiley & Sons, Hoboken, New Jersey, USA. MR1925014 MR1925014
  25. OWEN, A. B. (2001). *Empirical likelihood*. Chapman and Hall/CRC, Boca Raton, FL, USA.
  26. PEPE, M. S. and FLEMING, T. R. (1991). A nonparametric method for dealing with mismeasured covariate data. *Journal of the American Statistical Association* **86** 108–113. MR1137103 MR1137103
  27. PIEGORSCH, W. W., WEINBERG, C. R. and TAYLOR, J. A. (1994). Non-hierarchical logistic models and case-only designs for assessing susceptibility in population-based case-control studies. *Statistics in Medicine* **13** 153–162.
  28. QIN, J. (2000). Combining parametric and empirical likelihoods. *Biometrika* **87** 484–490. MR1782493 MR1782493
  29. QIN, J. and LAWLESS, J. (1994). Empirical likelihood and general estimating equations. *The Annals of Statistics* 300–325. MR1272085 MR1272085
  30. QIN, J., ZHANG, B. and LEUNG, D. H. (2009). Empirical likelihood in missing data problems. *Journal of the American Statistical Association* **104** 1492–1503. MR2750574 MR2750574
  31. REILLY, M. and PEPE, M. S. (1995). A mean score method for missing and auxiliary covariate data in regression models. *Biometrika* **82** 299–314. MR1354230 MR1354230
  32. RIVERA-RODRIGUEZ, C., HANEUSE, S., WANG, M. and SPIEGELMAN, D. (2020). Augmented pseudo-likelihood estimation for two-phase studies. *Statistical Methods in Medical Research* **29** 344–358. MR4064512 MR4064512
  33. ROBINS, J. M., ROTNITZKY, A. and ZHAO, L. P. (1994). Estimation of regression coefficients when some regressors are not always observed. *Journal of the American Statistical Association* **89** 846–866. MR1294730 MR1294730
  34. SCHAID, D. J., JENKINS, G. D., INGLE, J. N. and WEINSHILBOUM, R. M. (2013). Two-phase designs to follow-up genome-wide association signals with DNA resequencing studies. *Genetic epidemiology* **37** 229–238.
  35. SCOTT, A. J. and WILD, C. J. (1997). Fitting regression models to case-control data by maximum likelihood. *Biometrika* **84** 57–71. MR1450191 MR1450191
  36. SCOTT, A. J. and WILD, C. J. (2011). Fitting regression models with response-biased samples. *Canadian Journal of Statistics* **39** 519–536. MR2842429 MR2842429
  37. TAO, R., ZENG, D. and LIN, D.-Y. (2017). Efficient Semiparametric Inference Under Two-Phase Sampling, With Applications to Genetic Association Studies. *Journal of the American Statistical Association* **112** 1468–1476. MR3750869 MR3750869
  38. TSIATIS, A. (2007). *Semiparametric theory and missing data*. Springer Science & Business Media, New York, NY, USA. MR2233926 MR2233926
  39. WEAVER, M. A. and ZHOU, H. (2005). An estimated likelihood method for continuous outcome regression models with outcome-dependent sampling. *Journal of the American Statistical Association* **100** 459–469. MR2160550 MR2160550
  40. ZHANG, Z. and ROCKETTE, H. E. (2006). Semiparametric maximum likelihood for missing covariates in parametric regression. *Annals of the Institute of Statistical Mathematics* **58** 687–706. MR2345163 MR2345163
  41. ZHAO, Y., LAWLESS, J. F. and MCLEISH, D. L. (2009). Likelihood methods for regression models with expensive variables missing by design. *Biometrical Journal: Journal of Mathematical Methods in Biosciences* **51** 123–136. MR2667516 MR2667516
  42. ZHOU, H., WEAVER, M. A., QIN, J., LONGNECKER, M. and WANG, M. C. (2002).

A semiparametric empirical likelihood method for data from an outcome-dependent sampling scheme with a continuous outcome. *Biometrics* **58** 413–421. MR1908182  
MR1908182

*Note: This list reflects references listed in the original paper as accurately as possible with no attempt to correct errors.*