

The final version of this paper was published in Epidemiology (September 2011 - Volume 22 - Issue 5 - pp 680-685). The online link is [here](#).

BRIEF REPORT

Does More Accurate Exposure Prediction Improve Health Effect Estimates?

ADAM A. SZPIRO^{*,†}

Department of Biostatistics, University of Washington

CHRISTOPHER J. PACIOREK[‡]

*Department of Biostatistics, Harvard School of Public Health
Department of Statistics, University of California, Berkeley*

LIANNE SHEPPARD[†]

Department of Biostatistics, University of Washington

SUMMARY

A unique challenge in air pollution cohort studies and similar applications in environmental epidemiology is that exposure is not measured directly at subject locations. Instead, pollution data from monitoring stations at different locations than the subjects are used to predict exposures, and these predicted exposures are used to estimate the health effect parameter of interest. It has been widely assumed that it is desirable to minimize the error in predicting the true exposure in order to improve health effect estimation. We show in a simulation study that this is not always the case. We interpret our results in light of recently developed statistical theory for measurement error, and we discuss implications for the design and analysis of future epidemiological research.

^{*}To whom correspondence should be addressed: Department of Biostatistics, University of Washington, Seattle WA 98195, tel: 206-616-6846, fax: 206-543-3286, email: aszpiro@u.washington.edu

[†]Supported by the United States Environmental Protection Agency through R831697 and Assistance Agreement CR-83407101 and by the National Institute of Environmental Health Sciences through R01-ES009411 and 5P50ES015915.

[‡]Supported by the National Institute of Environmental Health Sciences through R01 ES017017.

1. ACCURATE EXPOSURE PREDICTION MAY NOT IMPROVE HEALTH EFFECT ESTIMATION

There has been a significant emphasis in air pollution epidemiology research on developing statistical models to predict exposures at subject locations where measurements are not available (Yanosky et al. 2009; Szpiro et al. 2010; Brauer 2010; Fanshawe et al. 2008; Su et al. 2009; Jerrett et al. 2005a; Hoek et al. 2008). These efforts are predicated on the assumption that exposure predictions with less measurement error relative to the unknown true values will improve health effect estimation (Jerrett et al. 2005b; Kunzli et al. 2005; Puett et al. 2009). We demonstrate in a simulation study that this is not always the case, and we interpret our results using recently developed statistical theory for measurement error resulting from spatially misaligned data (Szpiro et al. 2011).

2. MATHEMATICAL FRAMEWORK AND SIMULATION STUDY

Most modern statistical models for predicting long-term average air pollution concentrations are based on “land-use” regression (LUR). In LUR modeling, a linear regression model with geographic (land-use) covariates such as population density, proximity to traffic, and proximity to commercial areas is fit to monitoring data and is then used to predict concentrations at subject locations. Elaborations on this framework account for spatial and spatio-temporal correlation and various approaches to model selection, but LUR remains a central component. We focus on a pure LUR model in this paper.

2.1 Stochastic data-generating model

Consider an association study with the $N \times 1$ vector of observed health outcomes Y , $N \times 1$ vector of exposures X , and $N \times m$ matrix of covariates Z . Assume a linear regression model

$$Y = \beta_0 + X\beta_X + Z\beta_Z + \varepsilon, \quad (2.1)$$

with coefficient of interest β_X and ε an $N \times 1$ random vector with independent elements distributed as Gaussian random variables with mean 0 and variance σ_ε^2 (i.e., $N(0, \sigma_\varepsilon^2)$).

We are interested in the situation where Y and Z are observed, but instead of X we observe the $N^* \times 1$ vector X^* of exposures at different locations. N^* is the number of exposure monitors. Assume that X , the subject exposures, and X^* , the monitor concentrations, are jointly distributed as

$$\begin{pmatrix} X \\ X^* \end{pmatrix} = \begin{pmatrix} S \\ S^* \end{pmatrix} \alpha + \begin{pmatrix} \eta \\ \eta^* \end{pmatrix}. \quad (2.2)$$

In this expression, S and S^* are random $N \times k$ and $N^* \times k$ dimensional matrices of the k geographic covariates used in the LUR model observed without error, α is an unknown $k \times 1$ vector of coefficients, and η and η^* are independent vectors with elements distributed as $N(0, \sigma_\eta^2)$. The stochasticity in S and S^* derives from random selection of subject and monitor locations. If the exposure model is known, it is standard practice to estimate α based on X^* and then use $W = S\hat{\alpha}$ in place of X in equation (2.1) to estimate β_X . That is predictions from the LUR model are used as estimated exposures in place of the unknown true values, a form of regression calibration (Gryparis et al. 2009).

We quantify the accuracy in approximating X by W by

$$R_W^2 = 1 - \sum_{i=1}^N (W_i - X_i)^2 / \sum_{i=1}^N \left(X_i - \frac{1}{N} \sum_{i=1}^N X_i \right)^2,$$

where larger R_W^2 values correspond to less measurement error. This defines an out-of-sample measure of prediction accuracy since it is based on prediction error at subject locations, and it is not subject to bias from overfitting the exposure model to the monitoring data (Hastie et al. 2001, Ch. 7). R_W^2 is a random quantity that varies for each realization of the data-generating model, and we denote its expectation \bar{R}_W^2 .

There are a number of criteria for evaluating the validity and reliability of health effect estimates. Following Kim et al. (2009), we consider bias, standard deviation, root mean squared error (RMSE), and

coverage probability (the proportion of 95% confidence intervals that include the true β_X).

2.2 Misspecified exposure model

We generally do not know the exact form of the exposure model and may use a misspecified model for prediction. One form of model misspecification is to omit a geographic covariate from the LUR model. This corresponds to only observing the $N \times (k - 1)$ and $N^* \times (k - 1)$ matrices S' and $S^{*'}$ obtained by deleting the k th columns of S and S^* . We then estimate the corresponding $(k - 1) \times 1$ vector of coefficients α' and replace X in equation (2.1) by $W' = S'\hat{\alpha}'$ to obtain $\hat{\beta}'_X$. We denote measures of exposure prediction accuracy $R_W^{2'}$ and $\bar{R}_W^{2'}$ as in the case of the correctly specified exposure model.

We generally expect R_W^2 to be larger than $R_W^{2'}$, which from the perspective of exposure modeling implies that the correctly specified exposure model gives better predictions than the misspecified one. It is reasonable to expect that this will also lead to improved health effect estimation. However, in the next subsection we will demonstrate a class of examples in which R_W^2 is consistently larger than $R_W^{2'}$, but $\hat{\beta}_X$ has more error than $\hat{\beta}'_X$ as measured in terms of bias, variance, RMSE, and coverage probability. We emphasize that R_W^2 is not inflated by overfitting since it is based on the correctly specified exposure model and quantifies out-of-sample prediction accuracy at subject locations.

2.3 Simulation study

We set $k = 4$ (three geographic covariates and an intercept) and consider scenarios with N between 100 and 10,000 subjects and N^* equal to 100 monitors. We assume the three geographic covariates are independent of each other at all locations and are independent between subjects. In particular, for each subject i we assume the j th geographic covariate S_{ij} is independently distributed as $N(0, 1)$. Similarly

we assume the S_{ij}^* are distributed as $N(0, 1)$ for $j = 1, 2$, but the third geographic covariate for the monitoring sites is distributed as $N(0, \sigma^2)$ for $\sigma^2 = 0.1, 1.0$, or 4.0 . Finally, we set $\alpha_0 = 0$, $\alpha_j = 4$ for $j = 1, 2, 3$, $\beta_0 = 1$, $\beta_X = 2$, $\sigma_\varepsilon = 25$, and $\sigma_\eta = 4$, and we assume there are no additional covariates Z . Example simulation code in R (R Development Core Team 2010) can be found in the Appendix.

The choice of σ^2 controls the level of variability in the third geographic covariate at the monitoring locations. By comparing the misspecified model (i.e., the model that does not contain the third geographic covariate) to the correctly specified full model, we are able to assess the added value of including the third geographic covariate in predictions, depending on its variability. The situation with $\sigma^2 = 0.1$ is of particular interest, as it represents a geographic covariate that has limited variability in the monitoring data compared to the other geographic covariates but is equally variable in the subject data where it will be used to predict exposures. This is realistic, for example, if the covariate measures near-road traffic exposure. Regulatory monitors are often sited away from roadways in order to measure background pollution levels, so they may not span the full range of covariate values relevant for predicting exposures at subject home locations, a significant fraction of which are near major roads.

In Table 1 and Figure 1, we show the results from 80,000 Monte Carlo simulations with $N = 10,000$ subjects, $N^* = 100$ monitoring sites, and $\sigma^2 = 1.0$. The coefficient for the third geographic covariate α_3 is estimated well in the full model and is statistically significant in all simulations. The corresponding exposure prediction accuracy R_W^2 is consistently near 0.75, compared to $R_W^{2'}$ near 0.50 with the misspecified model. Health effect estimation efficiency is improved by using the correctly specified exposure model, which gives a standard deviation for $\hat{\beta}_X$ of 0.12 compared to 0.21 for $\hat{\beta}_X'$ with the misspecified model. The coverage probabilities for both models are poor since the standard error estimates fail to account for exposure measurement error. The correctly specified exposure model results in a modest improvement in

coverage probability, although it also introduces slightly more bias than the misspecified model.

Analogous results are shown in Table 1 and Figure 2 for $\sigma^2 = 0.1$, representing a situation where one of the geographic covariates is less variable in the distribution of monitoring locations than are the other geographic covariates. The smaller value of σ^2 results in more variability in estimating α_3 , but this parameter is still estimated well and is statistically significant in 83% of Monte Carlo simulations. There is clear improvement in the exposure predictions from using the full model with R_W^2 at least 0.67 in 95% of simulations, as compared to the misspecified model with $R_W^{2'}$ consistently near 0.50. But in this situation, the health effect estimation is more precise when we use the misspecified exposure model, with the standard deviation of $\hat{\beta}_X'$ equal to 0.16, compared to 0.23 for $\hat{\beta}_X$ using the fully specified model. The misspecified model also results in less bias and a modest improvement in coverage probability.

We vary the number of subjects as well as σ^2 and summarize the results in Figure 3 by plotting the difference between the standard deviation of $\hat{\beta}_X'$ based on the misspecified exposure model and $\hat{\beta}_X$ based on the correct exposure model on the vertical axis against N on the horizontal axis; a positive difference indicates the correctly specified model is more efficient. We restrict to 5,000 Monte Carlo simulations since this is sufficient to estimate the standard deviations (the biases are smaller and require more Monte Carlo simulations). The difference is positive for $\sigma^2 = 1.0$ and 4.0, consistent with the prior expectation that more accurate exposure predictions result in more efficient health effect estimation. But it is negative for $\sigma^2 = 0.1$ except for the case where there are only $N = 100$ subjects, demonstrating that in larger health studies the misspecified exposure model results in more efficient health effect estimation even though it gives less accurate exposure predictions. For all simulations we considered, the average out-of-sample exposure model prediction accuracies are \bar{R}_W^2 between 0.73 and 0.75 for the correctly specified model and $\bar{R}_W^{2'}$ between 0.49 and 0.50 for the misspecified model that omits the third geographic covariate.

3. THEORETICAL INTERPRETATION IN A MEASUREMENT ERROR FRAMEWORK

The results of our simulation study seem paradoxical in that we have shown a class of examples where more accurate exposure predictions do not lead to improved health effect estimation. Table 1 shows that for $\sigma^2 = 0.1$ the correctly specified model consistently gives more variable exposure predictions and more accurate out-of-sample prediction than the misspecified exposure model. However, a small part of the additional exposure variability is induced by error in estimating α_3 , which leads to less efficient estimation of β_X . These findings can be understood in a theoretical context by referring to the statistical measurement error framework developed for this setting by Szpiro et al. (2011); see also Gryparis et al. (2009).

Briefly, for a fairly general class of exposure models there are two components to the measurement error. The Berkson-like component of error results from smoothing the exposure surface using a model that may not account for all sources of variation and can be thought of as the part of the true exposure that is not predictable from the model. It is similar to standard Berkson error (Carroll et al. 2006) in that it inflates the health effect estimate standard deviation and introduces little or no bias, but it is different from Berkson error in that it is correlated in space and is not completely independent of the predicted exposures. The classical-like component results from uncertainty in estimating the exposure model parameters. It is similar to classical measurement error since it is a source of variability in the predicted exposures and can introduce bias in health effect estimates as well as change their standard errors. The classical-like component is different from classical measurement error since the additional variability from exposure model parameter estimation is shared across all prediction locations rather than being independent. Additional details on this decomposition can be found in Szpiro et al. (2011).

For the simple LUR exposure model considered here, the Berkson-like component is pure Berkson

error because there is no spatial dependence structure in η and η^* and the S_{ij} and S_{ij}^* are independent. When we use the correctly specified exposure model, the Berkson error is just η , but misspecifying the model by omitting the third geographic covariate increases the Berkson error substantially, resulting in a degradation of prediction accuracy. However, Berkson error plays the same role mathematically as the random ε in the disease model, so its impact on the health effect estimation error diminishes for large N . On the other hand, each coefficient that needs to be estimated in the exposure model contributes to the classical-like error, and this part of the error remains important regardless of the number of subjects. In some situations, this could result in a bias-variance tradeoff since classical-like error induces bias while Berkson-like error does not.

It turns out that for $\sigma^2 = 0.1$ in the monitoring data, we get relatively variable estimates of α_3 when using the full exposure model, while still improving out-of-sample prediction accuracy at subject locations. This results in substantial classical-like measurement error that (for sufficiently large N) is more important than the additional Berkson error that is introduced by omitting the corresponding geographic covariate. There is very little bias in any of our simulations, so the dominant classical-like error primarily results in more variable estimates of β_X .

4. IMPLICATIONS FOR FUTURE RESEARCH

We have shown a class of examples where more accurate exposure prediction does not lead to improved health effect estimation. It bears emphasis that this does not result from overfitting the exposure model, at least not as overfitting is traditionally understood for prediction models (Hastie et al. 2001, page 194). In all cases we considered, using the correctly specified model that includes all three geographic covariates results in improved prediction accuracy, as measured by out-of-sample R_{WV}^2 evaluated at the subject loca-

tions. In addition, the estimated coefficient $\hat{\alpha}_3$ for this covariate is nearly always statistically significant.

Our findings have important implications for the design and analysis of future environmental epidemiological studies. Most importantly, we believe that a paradigm shift is needed in environmental epidemiology. Development of models for exposure prediction and health effect estimation should be considered simultaneously, as opposed to the current practice of treating them distinctly by first selecting an exposure model to optimize prediction accuracy and then using the resulting predictions for health effect estimation. Recent papers that address measurement error in air pollution cohort studies represent progress in this direction (Kim et al. 2009; Szpiro et al. 2011; Gryparis et al. 2009; Madsen et al. 2008). Our results do not necessarily suggest employing a joint statistical estimation model for the exposure and health parameters in which the health data would influence estimation of the exposure model parameters. The issue we have highlighted relates more directly to model selection than parameter estimation.

There is extensive literature on penalization and other methods for optimizing accuracy of prediction models (Hastie et al. 2001), but these techniques are not directly applicable because better prediction accuracy may induce less precise health effect estimation. New statistical methodology is needed to select exposure models to optimize efficiency of health effects inference, perhaps involving alternative forms of penalization that account for the structure in both the monitoring and health outcome data. It is also worth exploring asymptotic methods to estimate the bias and variance of $\hat{\beta}_X$ in order to select optimal geographic covariates, particularly when there is a relatively large number of monitoring locations compared to the geographic covariates.

Since the relative benefits of different air pollution exposure models depend on the variability of geographic covariates in the subject population and monitor locations and on the size of the cohort, it is evident that study design can be improved by accounting for statistical issues at the intersection of expo-

sure prediction and health effect estimation. All else being equal, it is preferable to design an exposure monitoring campaign to maximize the variability of pertinent geographic covariates across monitor locations. An asset allocation algorithm (Kanaroglou et al. 2005) may be useful for optimizing the monitoring design to predict exposures in an epidemiology study with known subject locations.

We have only considered the relatively simple setting of a linear disease model with an exposure model that is LUR with independent geographic covariates. Even in this case we have shown that more accurate exposure prediction does not necessarily lead to improved health effect estimation. We expect that similar phenomena can occur in other settings, but further research is needed to identify general conditions and assess the implications of more complex situations.

REFERENCES

- M. Brauer. How much, how long, what, and where: Air pollution exposure assessment for epidemiologic studies of respiratory disease. *Proceedings of the American Thoracic Society*, 7:111–115, 2010.
- R. J. Carroll, D. Ruppert, L. A. Stefanski, and C. M Crainiceanu. *Measurement Error in Nonlinear Models: A Modern Perspective, Second Edition*. Chapman and Hall/CRC, 2006.
- T. R. Fanshawe, P. J. Diggle, S. Rushton, R. Sanderson, P. W. W. Lurz, S. V. Glinianaia, M. S. Pearce, L. Parker, M. Charlton, and T. Pless-Mulloli. Modelling spatio-temporal variation in exposure to particulate matter: a two-stage approach. *Environmetrics*, 19(6):549–566, 2008.
- A. Gryparis, C. J. Paciorek, A. Zeka, J. Schwartz, and B. A. Coull. Measurement error caused by spatial misalignment in environmental epidemiology. *Biostatistics*, 10(2):258–274, 2009.
- T. Hastie, R. Tibshirani, and J. Friedman. *Elements of Statistical Learning*. Springer, New York, 2001.
- G. Hoek, R. Beelen, K. de Hoogh, D. Vienneau, J. Gulliver, P. Fischer, and D. Briggs. A review of land-use regression models to assess spatial variation in outdoor air pollution. *Atmospheric Environment*, 42:7561–7578, 2008.
- M. Jerrett, A. Arain, P. Kanaroglou, B. Beckerman, D. Potoglou, T. Sahsuvaroglu, J. Morrison, and C. Giovis. A review and evaluation of intraurban air pollution exposure models. *Journal of Exposure Analysis and Environmental Epidemiology*, 15:185–204, 2005a.
- M. Jerrett, R. T. Burnett, R. Ma, C. A. Pope, D. Krewski, K. B. Newbold, G. Thurston, Y. Shi, N. Finkelstein, E. E. Calle, and M. J. Thun. Spatial analysis of air pollution mortality in Los Angeles. *Epidemiology*, 16(6):727–736, 2005b.
- P. S. Kanaroglou, M. Jerrett, J. Morrison, B. B. Beckerman, M. A. Arain, N. L. Gilbert, and J. R. Brook. Establishing an air pollution monitoring network for intraurban population exposure assessment: A location-allocation approach. *Atmospheric Environment*, 39:2399–2409, 2005.
- S. Y. Kim, L. Sheppard, and H. Kim. Health effects of long-term air pollution: influence of exposure prediction methods. *Epidemiology*, 20(3):442–450, 2009.
- N. Kunzli, M. Jerrett, W. J. Mack, B. Beckerman, L. LaBree, F. Gilliland, D. Thomas, J. Peters, and H. N. Hodis.

- Ambient air pollution and atherosclerosis in Los Angeles. *Environmental Health Perspectives*, 113(2):201–206, 2005.
- L. Madsen, D. Ruppert, and N. S. Altman. Regression with spatially misaligned data. *Environmetrics*, 19:453–467, 2008.
- R. C. Puett, J. E. Hart, J. D. Yanosky, C. Paciorek, J. Schwartz, H. Suh, F. E. Speizer, and F. Laden. Chronic fine and coarse particulate exposure, mortality, and coronary heart disease in the Nurses’ Health Study. *Environmental Health Perspectives*, 117(11):1697–1701, 2009.
- R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2010. URL <http://www.R-project.org>. ISBN 3-900051-07-0.
- J. G. Su, M. Jerrett, B. Beckerman, M. Wilhelm, J. K. Ghosh, and B. Ritz. Predicting traffic-related air pollution in Los Angeles using a distance decay regression selection strategy. *Environmental Research*, 109(6):657–670, 2009.
- A. Szpiro, P. D. Sampson, L. Sheppard, T. Lumley, S. D. Adar, and J. D. Kaufman. Predicting intra-urban variation in air pollution concentrations with complex spatio-temporal dependencies. *Environmetrics*, 21(4):606–631, 2010.
- A. Szpiro, L. Sheppard, and T. Lumley. Efficient measurement error correction for spatially misaligned data. *Bio-statistics*, In Press, 2011.
- J.D. Yanosky, C.J. Paciorek, and H. Suh. Predicting chronic fine and coarse particulate exposure using spatio-temporal models for the northeastern and midwestern US. *Environmental Health Perspectives*, 117:522–529, 2009.

	$\sigma^2 = 1$		$\sigma^2 = 0.1$	
	Correct model	Misspecified model	Correct model	Misspecified model
Exposure predictions				
\bar{R}_W^2	0.74	0.49	0.73	0.50
$\bar{\text{Var}}(W)$	48.5	32.7	50.2	32.3
Exposure model parameter estimate $\hat{\alpha}_3$				
Standard deviation	0.41	—	1.37	—
Statistically significant ($p < 0.05$)	100%	—	83%	—
Health effect parameter estimate $\hat{\beta}_X$				
Bias	−0.007	−0.001	−0.035	0.001
Standard deviation	0.12	0.21	0.23	0.16
RMSE	0.12	0.21	0.23	0.16
E(SE)	0.038	0.049	0.038	0.049
95% CI coverage	45%	35%	26%	46%

Table 1. Results from 80,000 Monte Carlo simulations with $N = 10,000$ and $N^* = 100$. The \bar{R}_W^2 and $\bar{\text{Var}}(W)$ are the out-of-sample prediction R_W^2 and variance of predicted exposures, respectively, averaged over 80,000 Monte Carlo simulations. The bias, standard deviation, and fraction of Monte Carlo runs statistically significant are given for estimates of $\hat{\alpha}_3$ in the correctly specified exposure model only, since this parameter is not included in the misspecified model. The bias, standard deviation, root mean squared error (RMSE), and 95% confidence interval coverage are given for estimates of the health effect parameter $\hat{\beta}_X$. We also report the average estimated standard error (SE) for $\hat{\beta}_X$. The Monte Carlo standard error in estimating the bias of $\hat{\beta}_X$ in all models is less than 0.001.

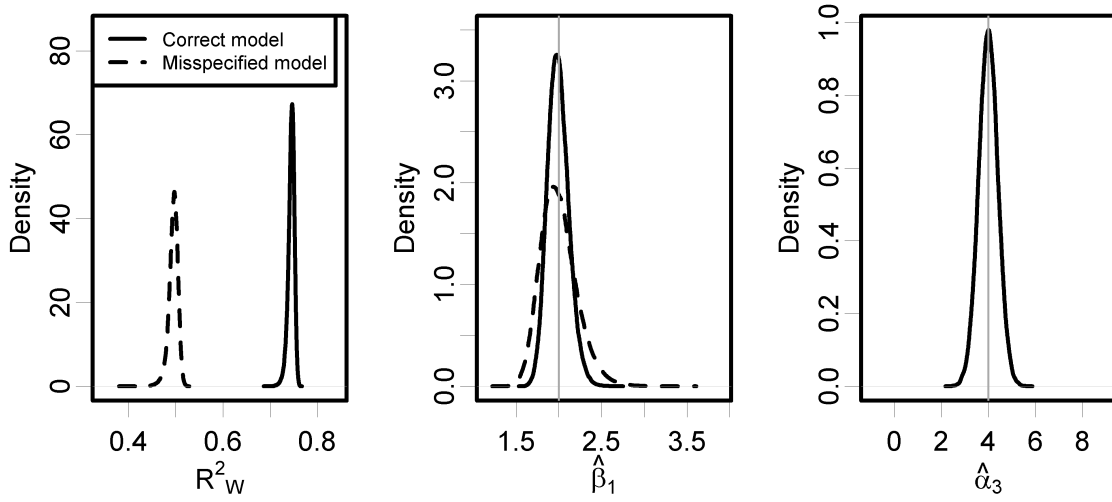


Fig. 1. Results from 80,000 Monte Carlo simulations with $N = 10,000$, $N^* = 100$, and $\sigma^2 = 1.0$. For the correctly specified exposure model the average out-of-sample prediction accuracy is $\bar{R}_W^2 = 0.74$ and the health effect estimation standard deviation is 0.12 with a bias of -0.007 (95% CI: -0.008 to -0.006). Corresponding statistics for the misspecified exposure model are $\bar{R}_W^{2'} = 0.49$ and health effect estimation standard deviation 0.21 with a bias of -0.001 (95% CI: -0.002 to 0.0006). The standard error of $\hat{\alpha}_3$ for the correctly specified model is 0.41, and $\hat{\alpha}_3$ is statistically significant in all simulations.

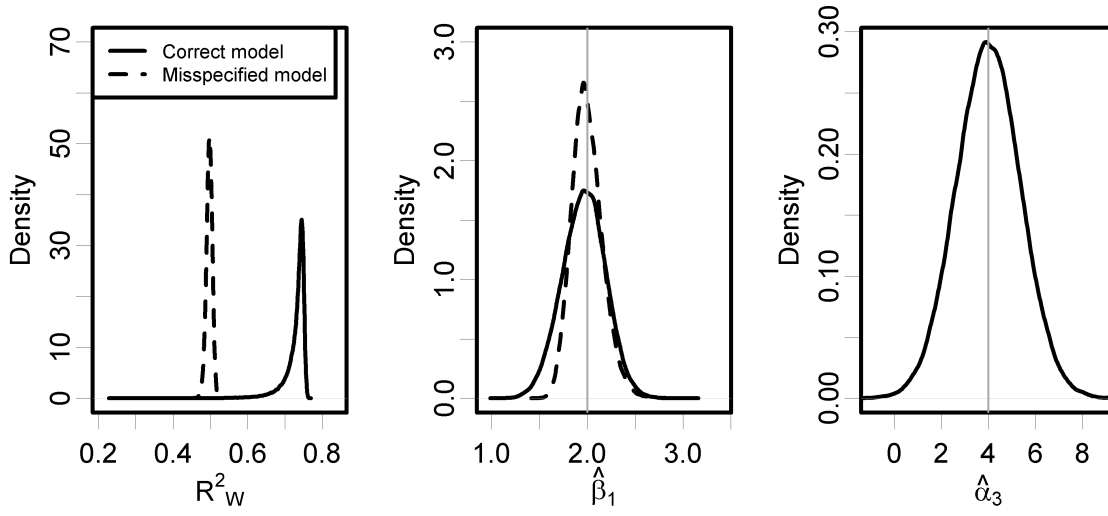


Fig. 2. Results from 80,000 Monte Carlo simulations with $N = 10,000$, $N^* = 100$, and $\sigma^2 = 0.1$. For the correctly specified exposure model the average out-of-sample prediction accuracy is $\bar{R}_W^2 = 0.73$ and the health effect estimation standard deviation is 0.23 with a bias of -0.035 (95% CI: -0.037 to -0.034). Corresponding statistics for the misspecified exposure model are $\bar{R}_W'^2 = 0.50$ and health effect estimation standard deviation 0.16 with a bias of 0.001 (95% CI: -0.0003 to 0.002). The density plot for R_W^2 shows some small outliers for the full model, but the prediction accuracy is better than for the misspecified model in all but 144 of the 80,000 simulations. The standard deviation of $\hat{\alpha}_3$ for the correctly specified model is 1.37, and $\hat{\alpha}_3$ is statistically significant in 83% of simulations.

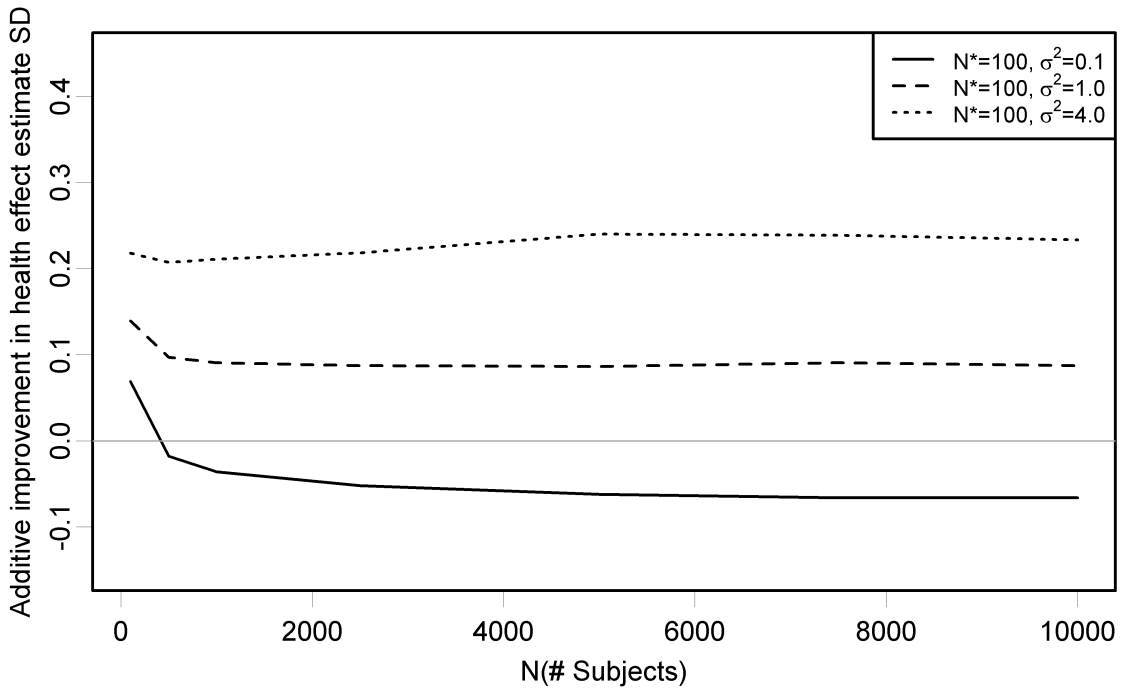


Fig. 3. Results from 5,000 Monte Carlo simulations with $N^* = 100$, $\sigma^2 = 0.1, 1.0, 4.0$, and N ranging from 100 to 10,000. The vertical axis shows the difference between standard deviation of $\hat{\beta}_X$ from the misspecified and correct exposure models. A positive difference indicates that the correctly specified model is more efficient. For all values of σ^2 , the average exposure model prediction accuracies are \bar{R}_W^2 between 0.73 and 0.75 and $\bar{R}_W^{2'}$ between 0.49 and 0.50.

APPENDIX

A. EXAMPLE R CODE TO PRODUCE THE RESULTS IN TABLE 1 .

```

set.seed(10000)
N <- 80000
n.subjs.list <- c(10000)
n.samps.list <- c(100,100)
z3.sd.list <- c(0.3,1.0)
for (i.scen in c(1:2)){
  n.samps <- n.samps.list[i.scen]
  z3.sd <- z3.sd.list[i.scen]
  fname <- paste("save.new2_",n.samps,"_",round(z3.sd,1),sep="")

  sd.eps <- 25
  sd.eta <- 4

  beta.est.list <- list()
  se.est.list <- list()
  r2.list <- list()
  exp.var.list <- list()
  alpha3.est.list <- list()
  alpha3.se.list <- list()

  alpha1 <- 4
  alpha2 <- 4
  alpha3 <- 4
  beta0 <- 1
  beta1 <- 2

  for (i.subjs.list in c(1:length(n.subjs.list))){

    n.subjs <- n.subjs.list[i.subjs.list]

    beta.est <- data.frame(matrix(rep(NA,N*3),ncol=3))
    se.est <- data.frame(matrix(rep(NA,N*3),ncol=3))
    r2 <- data.frame(matrix(rep(NA,N*3),ncol=3))
    exp.var <- data.frame(matrix(rep(NA,N*3),ncol=3))
    colnames(beta.est) <- c("true","model1","model2")
    colnames(se.est) <- c("true","model1","model2")
    colnames(r2) <- c("true","model1","model2")
    colnames(exp.var) <- c("true","model1","model2")
    alpha3.est <- rep(NA,N)
    alpha3.se <- rep(NA,N)

    for (i in 1:N){

      if (floor(i/1000)==i/1000) print(i)

      z1.subjs <- rnorm(n.subjs)
      z2.subjs <- rnorm(n.subjs)
      z3.subjs <- rnorm(n.subjs)

```

```

z1.samps <- rnorm(n.samps)
z2.samps <- rnorm(n.samps)
z3.samps <- rnorm(n.samps,0,z3.sd)

exp.subjs <- alpha1*z1.subjs + alpha2*z2.subjs + alpha3*z3.subjs +
  rnorm(n.subjs,0,sd.eta)
exp.samps <- alpha1*z1.samps + alpha2*z2.samps + alpha3*z3.samps +
  rnorm(n.samps,0,sd.eta)
y.subjs <- beta0 + beta1*exp.subjs+rnorm(n.subjs,0,sd.eps)

# true exposure
exp.subjs.est.true <- exp.subjs
lm.fit <- lm(y.subjs~exp.subjs.est.true)
beta.est[i,"true"] <- summary(lm.fit)$coef[2,1]
se.est[i,"true"] <- summary(lm.fit)$coef[2,2]
r2[i,"true"] <- summary(lm(exp.subjs.est.true~exp.subjs))$r.sq
exp.var[i,"true"] <- var(exp.subjs.est.true)

# model 1 (Misspecified Model)
exp.lm.fit <- lm(exp.samps~z1.samps+z2.samps)
exp.subjs.est.m1 <- exp.lm.fit$coef[1] + exp.lm.fit$coef[2]*z1.subjs +
  exp.lm.fit$coef[3]*z2.subjs
lm.fit <- lm(y.subjs~exp.subjs.est.m1)
beta.est[i,"model1"] <- summary(lm.fit)$coef[2,1]
se.est[i,"model1"] <- summary(lm.fit)$coef[2,2]
r2[i,"model1"] <- summary(lm(exp.subjs.est.m1~exp.subjs))$r.sq
exp.var[i,"model1"] <- var(exp.subjs.est.m1)

# model 2 (Correctly Specified Model)
exp.lm.fit <- lm(exp.samps~z1.samps+z2.samps+z3.samps)
exp.subjs.est.m2 <- exp.lm.fit$coef[1] + exp.lm.fit$coef[2]*z1.subjs +
  exp.lm.fit$coef[3]*z2.subjs + exp.lm.fit$coef[4]*z3.subjs
lm.fit <- lm(y.subjs~exp.subjs.est.m2)
beta.est[i,"model2"] <- summary(lm.fit)$coef[2,1]
se.est[i,"model2"] <- summary(lm.fit)$coef[2,2]
r2[i,"model2"] <- summary(lm(exp.subjs.est.m2~exp.subjs))$r.sq
exp.var[i,"model2"] <- var(exp.subjs.est.m2)
alpha3.est[i] <- exp.lm.fit$coef[4]
alpha3.se[i] <- summary(exp.lm.fit)$coef[4,2]
}

print(n.subjs)
beta.est.list[[i.subjs.list]] <- beta.est
se.est.list[[i.subjs.list]] <- se.est
r2.list[[i.subjs.list]] <- r2
exp.var.list[[i.subjs.list]] <- exp.var
alpha3.est.list[[i.subjs.list]] <- alpha3.est
alpha3.se.list[[i.subjs.list]] <- alpha3.se
}
save.image(file=paste(fname, ".Rdata", sep=""))
}

```