

## University of Groningen

### Educational note

Sudharsanan, Nikkil; Bijlsma, Maarten J

*Published in:*  
International Journal of Epidemiology

*DOI:*  
[10.1093/ije/dyab090](https://doi.org/10.1093/ije/dyab090)

**IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.**

*Document Version*  
Publisher's PDF, also known as Version of record

*Publication date:*  
2021

[Link to publication in University of Groningen/UMCG research database](#)

*Citation for published version (APA):*  
Sudharsanan, N., & Bijlsma, M. J. (2021). Educational note: Causal decomposition of population health differences using Monte Carlo integration and the g-formula. *International Journal of Epidemiology*, 50(6), 2098–2107. <https://doi.org/10.1093/ije/dyab090>

#### Copyright

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).

The publication may also be distributed here under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license. More information can be found on the University of Groningen website: <https://www.rug.nl/library/open-access/self-archiving-pure/taverne-amendment>.

#### Take-down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.


*Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.*



---

Education Corner

# Educational note: causal decomposition of population health differences using Monte Carlo integration and the g-formula

Nikkil Sudharsanan <sup>1\*</sup> and Maarten J Bijlsma<sup>2,3</sup>

<sup>1</sup>Heidelberg Institute of Global Health, Heidelberg University, Germany, <sup>2</sup>Laboratory of Population Health, Max Planck Institute for Demographic Research, Germany and <sup>3</sup>Groningen Research Institute of Pharmacy, Unit Pharmacotherapy, -Epidemiology & -Economics (PTEE), University of Groningen, the Netherlands

\*Corresponding author. Heidelberg Institute of Global Health, Im Neuenheimer Feld 130.3, 69120 Heidelberg, Germany.  
E-mail: nikkil.sudharsanan@uni-heidelberg.de

Editorial decision 25 March 2021; Accepted 19 April 2021

## Abstract

One key objective of the population health sciences is to understand why one social group has different levels of health and well-being compared with another. Whereas several methods have been developed in economics, sociology, demography, and epidemiology to answer these types of questions, a recent method introduced by Jackson and VanderWeele (2018) provided an update to decompositions by anchoring them within causal inference theory. In this paper, we demonstrate how to implement the causal decomposition using Monte Carlo integration and the parametric g-formula. Causal decomposition can help to identify the sources of differences across populations and provide researchers with a way to move beyond estimating inequalities to explaining them and determining what can be done to reduce health disparities. Our implementation approach can easily and flexibly be applied for different types of outcome and explanatory variables without having to derive decomposition equations. We describe the concepts of the approach and the practical steps and considerations needed to implement it. We then walk through a worked example in which we investigate the contribution of smoking to sex differences in mortality in South Korea. For this example, we provide both pseudocode and R code using our package, *cfdecomp*. Ultimately, we outline how to implement a very general decomposition algorithm that is grounded in counterfactual theory but still easy to apply to a wide range of situations.

**Key words:** Decomposition, causal inference, Monte Carlo, parametric g-formula, population models, health disparities

---

**Key Messages**

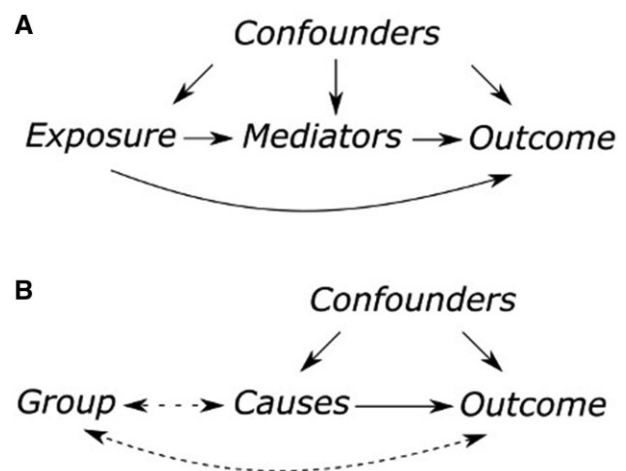
- Causal or counterfactual-based decomposition methods are of growing importance in epidemiology and the population health sciences.
- We develop and demonstrate a highly flexible implementation of the causal decomposition that is grounded in counterfactual theory but still easy to apply to a wide range of questions without having to derive specialized decomposition equations.
- We demonstrate how to use our decomposition algorithm to estimate the contribution of smoking to sex differences in the age-adjusted 1-year risk of mortality in South Korea, finding that smoking explains 27% of the male mortality disadvantage at ages  $\geq 50$  years.

**Introduction**

A central aim of the population health sciences is to understand why one social group has different levels of health and well-being compared with another. Recent examples of this question include understanding why Hispanics have worse congenital heart disease outcomes compared with non-Hispanics,<sup>1</sup> why adult mortality is higher in urban compared with rural Indonesia,<sup>2</sup> and why poorer individuals in Finland have higher mortality compared with more affluent individuals.<sup>3</sup> By identifying the sources of differences across populations, these studies provide an important first step in determining what can be done to reduce health disparities.

Decomposition analyses are one of the key tools for understanding the sources of differences in an outcome between groups and can help to move researchers from estimating to explaining health inequalities. At their core, decomposition analyses seek to determine how much of an observed difference in an outcome between two groups is due to the differing distribution of specific causes of that outcome between the groups. For example, in the example above on Finland, researchers may ask: ‘How much of the mortality difference between rich and poor individuals is due to the higher prevalence of smoking among poorer compared with richer individuals?’

Although such questions may sound like mediation analysis,<sup>4–8</sup> there is a key difference between mediation and decomposition. In a causal mediation analysis, we would first estimate the causal effect of poverty on mortality and then identify how much of this effect is driven through the causal effect of poverty on smoking. In a decomposition analysis, we are interested in how much smoking contributes to observed differences in mortality between poor and non-poor, and are agnostic to how much of the difference in smoking between poor and non-poor is due to the causal effect of smoking and how much is due to confounding causes. This crucial difference (depicted graphically using directed acyclic graphs in



**Figure 1** Directed acyclic graphs showing conceptual differences between mediation (A) and decomposition (B). Solid lines represent causal effects, whereas two-way dotted lines represent associations.

Figure 1) has consequences for the analytical approach to be taken and requires fewer confounding variables to be accounted for. Importantly, in a decomposition analysis, since we are not attempting to estimate the causal effect of the group variable (the exposure in a mediation analysis), we do not have to contend with the open issue of whether causal effects can be estimated for non-manipulable characteristics such as race.<sup>9</sup>

Various methods have been developed across disciplines for conducting decomposition analyses. Regression decompositions, such as the Oaxaca-Blinder decomposition<sup>10,11</sup> and its non-linear extensions,<sup>12,13</sup> use individual-level data and are employed frequently in economics and sociology,<sup>14</sup> whereas approaches using aggregate-level data are common in demography.<sup>15–18</sup> Recent advances in epidemiology provide a new perspective to decompositions, situating them in causal inference and counterfactual theory.<sup>2,9,19,20</sup> Among these, Jackson and VanderWeele’s (2018) provide an important advance by framing decomposition analyses around interventions to reduce disparities, where the

importance of specific characteristics to differences between populations is evaluated through hypothetical intervention scenarios to equalize these characteristics between groups.<sup>19</sup>

In this paper, we demonstrate a simple way to implement the counterfactual decomposition using parametric models and Monte Carlo integration. We focus on a worked example that asks ‘How much of the observed sex difference in mortality in South Korea is due to the higher prevalence of smoking among men compared with women?’ and demonstrates how to decompose sex differences in the age-adjusted 1-year mortality risk ratio between men and women. Our approach is based on a straightforward algorithm for estimating counterfactual decompositions for different outcome distributions without having to derive decomposition equations and can be easily applied within common statistical packages or implemented with our R package, *cfdecomp*.<sup>21</sup>

## A counterfactual approach to decomposition

### Concepts

We motivate and develop our approach through the question: ‘What is the contribution of smoking to sex differences in mortality in South Korea?’ We adopt a counterfactual perspective and define ‘contribution’ by asking ‘How large would the difference in mortality be if men and women counterfactually had an equal smoking prevalence?’

Our first main step is to specify exactly what level of smoking prevalence we are equalizing men and women to. When the relationship between an outcome (such as mortality) and a mediator (such as smoking) is non-linear, this choice can affect the contribution estimate.<sup>18</sup> Therefore, the choice of the counterfactual mediator distribution should be informed by substantive concerns (e.g. what makes sense from a policy perspective?) and inferential concerns (e.g. certain values may be outside the range observed in the data and should therefore be avoided). We choose to set men to have the smoking prevalence of women, since this maps to a clear intervention that public health policymakers may seek to achieve.

The second main step is to specify a summary population measure. This is the measure that we will use to compare the mortality of men and women in South Korea. For our example, we consider the age-adjusted 1-year risk of death. In theory, our approach can be extended to decompose more complicated summary measures, such as disability-adjusted life years lost or period life expectancy. However, decomposing such measures requires additional,

often stronger, assumptions. For this reason, we do not cover the application of our method to those summary measures here and choose rather to focus on common summary measures with clear assumptions.

Third, we need to specify contrasts of these summary measures between men and women (i.e. how are we going to compare the summary measure?). We consider the risk ratio for men relative to women (adjusted for age). Our method also allows us to decompose other contrasts, such as the risk difference—a point that we will return to when describing the decomposition algorithm below.

Based on Steps 1–3, we can construct our estimate of the ‘contribution’ of smoking by seeing how much the difference in the summary measure between men and women reduces when we set men to have the same smoking prevalence as women. For example, we would compare the mortality risk ratio between men and women in the observed data to the mortality risk ratio between men and women in a counterfactual world in which we set men to have the same smoking levels as women. We could then estimate the contribution of smoking as the percentage reduction in the male–female mortality disparity. Note that this contribution is not bounded between 0 and 1, and could result in negative contributions or contributions of >100%. This is not an issue, however; this situation occurs in both mediation and decomposition analyses when the indirect effect (the association via the mediators) and the direct effect (the association not via the mediators) are of opposite signs and hence partially cancel each other out in the total effect. Indeed, many recent papers using mediation and decomposition analyses have found contribution estimates of >100 or <0.<sup>2,19,22</sup> Contribution estimates of <0 or >1 could also occur due to imprecision in the underlying estimates. For this reason, it is important to present and interpret such estimates with their accompanying standard error. In [Supplementary Appendix 3](#), available as [Supplementary data](#) at *IJE* online, we provide a more general formal exposition of the causal decomposition.

### Parametric modelling and Monte Carlo-based estimation

The core estimand in our decomposition is the counterfactual summary measure of mortality for men if they were set to have the same smoking distribution as women. Estimating this counterfactual requires (i) a way to match the smoking distribution between men and women, and (ii) a way to re-estimate mortality as a function of the new smoking distribution. Importantly, since we are interested in the effect of changing the level of smoking on mortality,

our approach to re-estimating mortality needs to adjust for the confounders of the smoking–mortality relationship. Our solution to these two issues is to use the parametric g-formula and Monte Carlo integration.<sup>7,23–25</sup> This entire approach can be estimated by following a straightforward algorithm.

### Decomposition algorithm

*Step 0: Specify starting decisions.*

- i. Decide on a summary measure.
- ii. Decide on a contrast.
- iii. Decide on the counterfactual mediator distribution.

*Step 1: Estimate relationships in the data.*

- i. Fit regression model(s) for the mediator(s) of interest with confounders of the mediator–outcome relationship as covariates.
- ii. Fit regression model(s) for the outcome with the mediator(s) of interest and the same confounders as the mediator model.

*Step 2: Form the natural-course pseudo-population.*

- i. Use the coefficients from the mediator model(s) with the observed confounder values to simulate mediator values for each individual in the data.
- ii. Use the coefficients from the outcome model(s) together with the observed confounder values and the new simulated mediator values to simulate the outcome for each individual in the data. This is the natural-course pseudo-population.
- iii. Within this natural-course pseudo-population, estimate the summary measure for both groups and then form the contrast of interest across groups.

*Step 3: Form the counterfactual pseudo-population.*

- i. Use the coefficients from the mediator model(s) with the observed confounder values to simulate mediator values that follow the counterfactual mediator distribution.
- ii. Use the coefficients from the outcome model(s) together with the observed confounder values and simulated mediator values to simulate the outcome for each individual in the data. This is the counterfactual pseudo-population.
- iii. Within this counterfactual pseudo-population, estimate the summary measure for both groups and then form the contrast of interest across groups.

*Step 4: Compare the contrast of interest in the natural-course and counterfactual pseudo-populations.* To

estimate standard errors and to produce stable estimates of the contribution, we have to address two types of variability. First, since we are drawing values of the mediators and outcomes from probability distributions, the exact values assigned to individuals can change across multiple draws. This results in the estimate of the contribution also changing across draws (known as Monte Carlo error). To reduce this error, we conduct Steps 2 and 3 multiple times, each time drawing a new set of mediator and outcome values. We then construct the contrasts for each draw and then average across all these draws to produce stable natural-course and counterfactual estimates, before calculating the contribution in Step 4.

Second, because our results are based on a sample, we need to account for sampling variability. This is especially important for the construction of confidence intervals around the estimates. We use a bootstrap procedure to capture this uncertainty, drawing with replacement a fresh sample of size equal to the original data before Step 1, conducting the entire analysis  $k$  times, and then estimating the standard error of our decomposition estimates as the standard deviation of the estimates from the  $k$  bootstrap samples.

Our algorithm above treats the variables involved as time-fixed, which may not always be appropriate.<sup>5,8</sup> The algorithm can be easily expanded, however, to allow for time-varying variables; we present a time-varying version of the decomposition algorithm above in [Supplementary Appendix 2](#), available as [Supplementary data](#) at *IJE* online, based on Westreich *et al.* (2012).<sup>26</sup> A second important note is that the natural course is often used in g-formula analyses to validate the estimation models rather than as part of the estimand. In our algorithm, however, the natural course also forms part of the contribution estimate. We chose to use the natural-course estimate instead of the observed data in our estimand so that both the counterfactual and ‘as-is’ scenarios are based on the same underlying model. However, if the natural-course estimates do not approximate the data well, then that is evidence of model misspecification, which needs to be investigated further.

Both the size of and contribution of specific mediators to a health disparity are dependent on the scale that the disparity is measured on. For example, a difference in mortality between two populations and the contribution of smoking to this difference may vary based on whether the disparity is measured as a mortality risk ratio, a survival risk ratio or an absolute difference in mortality rates. A major strength of our decomposition algorithm is that the researcher is not limited to one scale and can estimate and explain the disparity using multiple

measures. This is because the decomposition algorithm works by first generating pseudo-populations based around model-predicted values rather than by comparisons of model coefficients.

### Empirical example: the contribution of smoking to sex differences in mortality in South Korea

We now demonstrate the application of the approach that we outlined in the previous section to real data from the Korean Longitudinal Study of Aging. In the interest of providing a simple pedagogical example, we conduct a stylized analysis and thus the results should be interpreted cautiously. A more rigorous analysis that fully explores and accounts for the different sources of confounding and measurement error is outside the scope of this example. The simplified example also raises conceptual issues that we omit discussion of, such as whether some of the confounders may instead mediate the relationship between ever smoking and mortality. However, to lend some credence to the analysis, we note that the results of our example are in line with other literature on the contribution of smoking to sex differences in mortality.<sup>27</sup>

#### Data: Korean Longitudinal Study of Aging

We use data from the 2006–2012 waves of the Korean Longitudinal Study of Aging—a nationally representative survey of South Korean individuals aged  $\geq 45$  years.<sup>28</sup> We use data on adults aged  $\geq 50$  years from the baseline 2006 waves, using the subsequent waves for mortality follow-up. Our total sample consists of 7615 individuals with 42 405 person-years of follow-up. We convert our data from a person to person-age format, with one observation for every age lived in the survey, along with a dichotomous indicator for whether an individual survived through or died at that age. Individuals leave the survey through death, censoring from loss to follow-up before 2012, or from censoring at the end of the survey period in 2012.

#### Main variables: outcome, mediator, and confounders

Our outcome of interest is a dichotomous indicator for whether an individual died or survived to the next age and our primary mediator is a dichotomous indicator for whether an individual reported ever regularly smoking cigarettes. We adjust for the following potential confounders of the smoking–mortality relationship: age, how frequently an individual reported drinking alcohol, schooling, urbanicity, and marital status.

#### Step 0: Specify a summary measure, contrast, and counterfactual distribution

Our main summary measure is the age-adjusted 1-year risk of death (surviving to the next age). For this summary measure, our contrast of interest is the risk ratio of mortality for men relative to women. We construct this contrast using the following Poisson regression on person-year observations (adjusting for age using indicator variables for 5-year age groups):

$$\log(E[Y|Female, Age]) = \alpha_0 + (\alpha_1 \cdot Female) + \sum_i (\alpha_i \cdot Agegr_i)$$

where  $\alpha_1$  is our estimate of interest. We use a Poisson regression here to just estimate the summary contrast (the exponent of  $\alpha_1$ ) but could have alternatively directly estimated an age-standardized risk ratio from the data. Importantly, because we are interested in the observed difference between men and women (adjusting for just age), we do not add any confounders to this model.<sup>19</sup>

For this analysis, we set the smoking levels among men to be equal to those among women as our counterfactual scenario.

#### Step 1: Estimate relationships in the data (using regression models)

*Mediator model.* We model the probability of ever regularly smoking for men and women using the following logistic-regression model:

$$\text{logit}(E[Smk|Female, Age, C]) = \beta_0 + (\beta_1 \cdot Female) + (\beta_2 \cdot Age) + (\beta_3 \cdot Age \cdot Female) + \sum_i (\beta_{c_i} \cdot C_i)$$

Here, *Smk* is a binary variable for whether an individual self-reported ever regularly smoking, *Sex* is the indicator variable for female, *Age* is a continuous measurement of age and  $C_i$  are the confounders described previously. We use this model to estimate the group  $\rightarrow$  causes association pathway in [Figure 1B](#). We include the confounders in this model, not to adjust for confounding, but rather to allow us to predict and match the sex-specific smoking prevalence within confounder strata.

*Outcome model.* We model mortality as a function of smoking, sex and the confounders by fitting the following logistic-regression model:



**Figure 2** Flowchart for simulating the natural-course and counterfactual smoking and mortality values for a single male in the data. The regression estimates are based on the models described in the 'Methods' section.

```

Start loop b from 1 to B
                                                                    BOOTSTRAP

Draw a bootstrap sample of the wide, person-level, data
bootstrap.data.wide <- sample.with.replacement(empirical.data)

Reshape bootstrapped data to the person-month level
bootstrap.data.long <- reshape.long(bootstrap.data.wide)

Fit the outcome model
outcome.model <- logistic.regression(died ~ female + ever.smoke + female*ever.smoke + age + age*female + age*ever.smoke + confounders, data =
bootstrap.data.long)

Fit the mediator model
mediator.model <- logistic.regression(ever.smoke ~ female + age + female*age+ confounders, data = bootstrap.data.wide )

Start loop m from 1 to M
                                                                    MONTE CARLO

Make a copy of the data within each Monte Carlo loop
montecarlo <- bootstrap.data.long

Form the natural course estimates

Draw values of smoking from the model-predicted probabilities
montecarlo$ever.smoke <- binomial.draw(probability = predict(mediator.model, data = montecarlo))

Draw values of mortality from model-predicted probabilities and updated smoking values
montecarlo$died <- binomial.draw(probability = predict(outcome.model, data = montecarlo))

Estimate the outcome and contrast: Age-adjusted mortality risk ratio
natural.course.risk.ratio.mc[m] <- coef(poisson.regression(died ~ female + age.groups, data = montecarlo))

Form the counterfactual estimates

Make a dataset for just men
men.montecarlo.cf <- montecarlo[men]

Assign them the sex identifier of women so that the counterfactual smoking values are drawn from the female distribution
men.montecarlo.cf$female <- 1

Draw values of smoking again, this time from the female probabilities
men.montecarlo.cf$ever.smoke <- binomial.draw(probability = predict(mediator.model, data = men.montecarlo.cf))

Set the sex identifier back to men before predicting mortality
men.montecarlo.cf$female <- 0

Draw values of mortality again, this time with the counterfactual smoking values
men.montecarlo.cf$died <- binomial.draw(probability = predict(outcome.model, data = men.montecarlo.cf))

Form the counterfactual pseudopopulation by updating the male values
montecarlo[men] <- men.montecarlo.cf

Estimate the outcome and contrast: Age-adjusted mortality risk ratio
counterfactual.risk.ratio.mc[m] <- coef(poisson.regression(died ~ female + age.groups, data = montecarlo))

End m

Save the mean values across Monte Carlo loops in the bth place in a vector
natural.course.risk.ratio[b] <- mean(natural.course.risk.ratio.mc)
counterfactual.risk.ratio[b] <- mean(counterfactual.risk.ratio.mc)
End b

```

**Figure 3** Example code for estimating the contribution of smoking to sex differences in mortality in South Korea. For this example, we have a binomial mediator 'smoke' (ever-smoker), binomial outcome 'died' (death in a person-year), our summary measures and contrast is the age-adjusted mortality risk ratio and, for the counterfactual scenario, we assign men the smoking distribution of women. In the models, C represents covariates needed for exchangeability.

$$\begin{aligned} \text{logit}\left(E[Y|Female, Smk, Age, C]\right) = & \delta_0 + \left(\delta_1 \cdot Female\right) \\ & + \left(\delta_2 \cdot Smk\right) + \left(\delta_3 \cdot Female \cdot smk\right) + \left(\delta_4 \cdot Age\right) + \left(\delta_5 \right. \\ & \left. \cdot Age \cdot Female\right) + \left(\delta_6 \cdot Age \cdot Smk\right) + \sum_i (\delta_{c_i} \cdot C_i) \end{aligned}$$

We use this model to estimate the causes  $\rightarrow$  outcome effect pathway in [Figure 1B](#).

### Steps 2 and 3: simulation to form the natural-course and counterfactual pseudo-populations

Based on the results of the two models, we simulate the natural-course and counterfactual pseudo-populations for both men and women. In [Figure 2](#), we provide a step-by-step example of how to use the regression estimates to form the simulated values for a single male individual in the data. The pseudocode in [Figure 3](#) and R code in the [Supplementary Material](#), available as [Supplementary data](#)



at *IJE* online, demonstrate how to do this for all individuals in the data using common statistical software.

#### Step 4: Calculate and compare the contrasts of interest and determine the percent contribution of smoking

Once pseudo-populations have been created, the final step is to calculate the contrast of interest. We then estimate the contribution of smoking to sex differences in mortality by measuring how much the contrast changes between the natural-course and counterfactual worlds. All steps needed to estimate the decomposition are also shown as pseudo-code in [Figure 3](#). We also provide code for how to estimate the example in R using our function *cfdecomp* in the [Supplementary Material](#), available as [Supplementary data](#) at *IJE* online.

**Table 1** Descriptive characteristics of the sample at baseline, in adults aged  $\geq 50$  years, Korean Longitudinal Study of Aging, 2006

	Men		Women	
	Mean	SD	Mean	SD
Age (years)	66.2	9.0	67.4	9.9
	%	<i>n</i>	%	<i>n</i>
Marital status				
Never married	0.01	105	0.00	100
Married/partnered	0.93	17 147	0.64	15 350
Separated/divorced	0.02	349	0.02	499
Widowed	0.05	893	0.33	7962
Completed schooling				
None	0.09	1706	0.31	7299
Elementary or middle	0.45	8249	0.53	12 574
More than middle	0.46	8539	0.17	4038
Rural	0.27	4987	0.27	6534
Ever-smoker	0.61	11 276	0.04	1015
Alcohol consumption				
None/less than once a month	0.43	7868	0.87	20 808
One to several times a month	0.16	3040	0.08	2000
One to several times a week	0.28	5119	0.04	906
Most days of the week	0.05	935	0.00	113
Every day of the week	0.08	1532	0.00	84

**Table 2** Estimates of the contribution of smoking to the age-adjusted 1-year mortality risk ratio using the counterfactual decomposition method, Korean Longitudinal Study of Aging, 2006–2012

	Natural-course RR (95% CI)	Counterfactual RR (95% CI)	Percent contribution (95% CI)
Mortality risk ratio for men relative to women	1.89 (1.65, 2.14)	1.65 (1.38, 1.92)	28% (8%, 47%)

## Results

*Descriptive characteristics.* The mean age was 66.2 years for men and 67.4 years for women ([Table 1](#)). A greater share of men were currently married compared with women (93% compared with 64%) due to a much higher proportion of widowhood among women (33% compared with 5%). There were important health and socio-economic differences between men and women. Men were far more likely to smoke (61% compared with 4%) and drink regularly (proportion who reported drinking at least once a week: 41% compared with 4%). Men were also substantially more likely to have completed more than middle school (46% compared with 17%).

*Decomposition of the age-adjusted 1-year risk of mortality.* Men were 1.89 times [95% confidence interval (CI): 1.65, 2.14] more likely to die within 1 year of an interview compared with women (after adjusting for age) ([Table 2](#)). After setting men to have the same smoking distribution of women, this risk ratio reduced to 1.65 (95% CI: 1.38, 1.92). The resulting change corresponds to a  $(1 - 0.65/0.89) = 28\%$  (95% CI: 0.08, 0.47) contribution of smoking to sex differences in the age-adjusted 1-year risk of mortality.

## Discussion

We introduce a general yet easily applied procedure for implementing counterfactual decompositions using the parametric g-formula and Monte Carlo integration.<sup>19</sup> We demonstrate this approach by estimating the contribution of smoking to sex differences in mortality in South Korea by decomposing the age-adjusted mortality risk ratio for men relative to women. We find that the large smoking difference between men and women in South Korea explains 27% of the age-adjusted mortality risk ratio among adults aged  $\geq 50$  years.

The age-adjusted mortality risk could also be decomposed using closed-form decomposition equations.<sup>12,13,19</sup> The algorithm we outline does not replace closed-form decomposition approaches, but rather provides an alternative using simulations, which provides two main advantages. First, we can decompose summary measures based on any outcome distribution in the generalized linear model family without having to derive or use separate decomposition

equations depending on whether an outcome is binomially, Poisson, or normally distributed. Moving between outcome distributions simply requires changing the regression type used to model the outcome in the decomposition algorithm.

The second advantage of the simulation algorithm is that we can easily switch between different contrasts, since we effectively regenerated entire micropopulations for the observed and counterfactual worlds. For example, once natural-course and counterfactual pseudo-populations have been generated, we decomposed the risk ratio by estimating Poisson regressions of mortality on sex within both pseudo-populations and measuring how the risk ratio changes between the natural-course and counterfactual worlds. If we were instead interested in decomposing the odds ratio, we would simply switch from Poisson to logistic regressions and compare the odds ratios.

Despite these advantages, our algorithm comes with important trade-offs compared with existing decomposition implementations. Compared with the closed-form equations, our approach requires substantial computational power and time. This is not a trivial consideration and decompositions with large data sets may take hours to even days to complete even when considerable computational power is available. Furthermore, as with any method seeking to provide causal explanations, the causal validity of the decomposition results hinges on assumptions of exchangeability (also known as no unmeasured confounding), common support (positivity), and consistency. We discuss these three issues in more detail in [Supplementary Appendix 1](#), available as [Supplementary data](#) at *IJE* online, for interested readers.

## Conclusions

Decomposing the sources of differences in health and other outcomes is a key research endeavour in epidemiology and other population health sciences. We describe an implementation of the counterfactual decomposition that builds on and generalizes the rich existing body of work on decomposition methods in the health and social sciences. The approach provides a highly flexible and easily implemented way of estimating decompositions that are grounded in potential outcomes and counterfactual theory, and applicable to a wide range of population health questions.

## Supplementary data

[Supplementary data](#) are available at *IJE* online.

## Ethics approval

This study uses publicly available and de-identified secondary data, and was exempt from institutional review-board approval.

## Funding

N.S. receives funding from the Alexander von Humboldt Foundation.

## Data availability

Data are freely available (after registration) at [g2aging.org](https://g2aging.org).

## Conflicts of interest

None declared.

## References

1. Peyvandi S, Baer RJ, Moon-Grady AJ *et al*. Socioeconomic mediators of racial and ethnic disparities in congenital heart disease outcomes: a population-based study in California. *J Am Heart Assoc* 2018;7:e010342.
2. Sudharsanan N, Ho JY. Rural–urban differences in adult life expectancy in Indonesia: a parametric g-formula based decomposition approach. *Epidemiology* 2020;31:393.
3. Martikainen P, Mäkelä P, Peltonen R, Myrskylä M. Income differences in life expectancy: the changing contribution of harmful consumption of alcohol and smoking. *Epidemiology* 2014;25:182–90.
4. Imai K, Keele L, Tingley D. A general approach to causal mediation analysis. *Psychol Methods* 2010;15:309–34.
5. Lin S-H, Young J, Logan R, Tchetgen EJT, VanderWeele TJ. Parametric mediational g-formula approach to mediation analysis with time-varying exposures, mediators, and confounders. *Epidemiol Camb Mass* 2017;28:266.
6. Lin S-H, Young JG, Logan R, VanderWeele TJ. Mediation analysis for a survival outcome with time-varying exposures, mediators, and confounders. *Stat Med* 2017;36:4153–66.
7. De Stavola BL, Daniel RM, Ploubidis GB, Micali N. Mediation analysis with intermediate confounding: structural equation modeling viewed through the causal inference lens. *Am J Epidemiol* 2015;181:64–80.
8. Vansteelandt S, Daniel RM. Interventional effects for mediation analysis with multiple mediators. *Epidemiol Camb Mass* 2017;28:258.
9. VanderWeele TJ, Robinson WR. On causal interpretation of race in regressions adjusting for confounding and mediating variables. *Epidemiol Camb Mass* 2014;25:473.
10. Blinder AS. Wage discrimination: reduced form and structural estimates. *J Hum Resour* 1973;8:436–55.
11. Oaxaca R. Male-female wage differentials in urban labor markets. *Int Econ Rev* 1973;14:693–709.
12. Powers DA, Yun M-S7. Multivariate decomposition for hazard rate models. *Sociol Methodol* 2009;39:233–63.
13. Yun M-S. Decomposing differences in the first moment. *Econ Lett* 2004;82:275–80.

14. Machado JA, Mata J. Counterfactual decomposition of changes in wage distributions using quantile regression. *J Appl Econom* 2005;20:445–65.
15. Kitagawa EM. Components of a difference between two rates. *J Am Stat Assoc* 1955;50:1168–94.
16. Arriaga EE. Measuring and explaining the change in life expectancies. *Demography* 1984;21:83–96.
17. Horiuchi S, Wilmoth JR, Pletcher SD. A decomposition method based on a model of continuous change. *Demography* 2008;45:785–801.
18. Andreev EM, Shkolnikov VM, Begun AZ. Algorithm for decomposition of differences between aggregate demographic measures and its application to life expectancies, healthy life expectancies, parity-progression ratios and total fertility rates. *Demogr Res* 2002;7:499–522.
19. Jackson JW, VanderWeele TJ. Decomposition analysis to identify intervention targets for reducing disparities. *Epidemiology* 2018;29:825–35.
20. Nandi A, Glymour MM, Subramanian S. Association among socioeconomic status, health behaviors, and all-cause mortality in the United States. *Epidemiology* 2014;25:170–77.
21. Bijlsma MJ, Sudharsanan N, Li P. *cfdecomp: Counterfactual Decomposition: MC Integration of the G-Formula*. 2020. <https://cran.r-project.org/package=cfdecomp> (2 February 2021, date last accessed).
22. Bijlsma MJ, Wilson B. Modelling the socio-economic determinants of fertility: a mediation analysis using the parametric g-formula. *J R Stat Soc Ser A Stat Soc* 2020;183:493–513.
23. VanderWeele TJ, Vansteelandt S, Robins JM. Effect decomposition in the presence of an exposure-induced mediator-outcome confounder. *Epidemiol Camb Mass* 2014;25:300.
24. Young JG, Tchetgen Tchetgen EJ. Simulation from a known Cox MSM using standard parametric models for the g-formula. *Stat Med* 2014;33:1001–14.
25. Keil AP, Edwards JK, Richardson DR, Naimi AI, Cole SR. The parametric G-formula for time-to-event data: towards intuition with a worked example. *Epidemiol Camb Mass* 2014;25:889–97.
26. Westreich D, Cole SR, Young JG *et al*. The parametric g-formula to estimate the effect of highly active antiretroviral therapy on incident AIDS or death. *Stat Med* 2012;31:2000–09.
27. Preston SH, Wang H. Sex mortality differences in the United States: the role of cohort smoking patterns. *Demography* 2006;43:631–46.
28. Jang S-N. Korean Longitudinal Study of Ageing (KLoSA): overview of research design and contents. *Encycl Geropsychol* 2015;1–9.