



Neighbourhood-related socioeconomic perinatal health inequalities: An illustration of the mediational g-formula and considerations for the big data context

Lizbeth Burgos Ochoa¹ | Lindsey van der Meer¹ | Adja J. M. Waelput¹ |
Jasper V. Been^{1,2,3} | Loes C. M. Bertens¹

¹Department of Obstetrics and Gynaecology, Erasmus MC, University Medical Centre Rotterdam, Rotterdam, The Netherlands

²Division of Neonatology, Department of Paediatrics, Erasmus MC – Sophia Children's Hospital, University Medical Centre Rotterdam, Rotterdam, The Netherlands

³Department of Public Health, Erasmus MC, University Medical Centre Rotterdam, Rotterdam, The Netherlands

Correspondence

Lizbeth Burgos Ochoa, Department of Obstetrics and Gynaecology, Erasmus MC, University Medical Centre Rotterdam, Rotterdam, The Netherlands.
Email: l.burgoschoa@erasmusmc.nl

Funding information

Erasmus Initiative Smarter Choices for Better Health.

Abstract

Background: Advances in computing power have enabled the collection, linkage and processing of big data. Big data in conjunction with robust causal inference methods can be used to answer research questions regarding the mechanisms underlying an exposure–outcome relationship. The g-formula is a flexible approach to perform causal mediation analysis that is suited for the big data context. Although this approach has many advantages, it is underused in perinatal epidemiology and didactic explanation for its implementation is still limited.

Objective: The aim of this was to provide a didactic application of the mediational g-formula by means of perinatal health inequalities research.

Methods: The analytical procedure of the mediational g-formula is illustrated by investigating whether the relationship between neighbourhood socioeconomic status (SES) and small for gestational age (SGA) is mediated by neighbourhood social environment. Data on singleton births that occurred in the Netherlands between 2010 and 2017 ($n = 1,217,626$) were obtained from the Netherlands Perinatal Registry and linked to sociodemographic national registry data and neighbourhood-level data. The g-formula settings corresponded to a hypothetical improvement in neighbourhood SES from disadvantaged to non-disadvantaged.

Results: At the population level, a hypothetical improvement in neighbourhood SES resulted in a 6.3% (95% confidence interval [CI] 5.2, 7.5) relative reduction in the proportion of SGA, that is the total effect. The total effect was decomposed into the natural direct effect (5.6%, 95% CI 5.1, 6.1) and the natural indirect effect (0.7%, 95% CI 0.6, 0.9). In terms of the magnitude of mediation, it was observed the natural indirect effect accounted for 11.4% (95% CI 9.2, 13.6) of the total effect of neighbourhood SES on SGA.

Conclusions: The mediational g-formula is a flexible approach to perform causal mediation analysis that is suited for big data contexts in perinatal health research. Its

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial](https://creativecommons.org/licenses/by-nc/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

© 2023 The Authors. *Paediatric and Perinatal Epidemiology* published by John Wiley & Sons Ltd.

application can contribute to providing valuable insights for the development of policy and public health interventions.

KEYWORDS

birth outcomes, health inequalities, mediation analysis, neighbourhood

1 | BACKGROUND

1.1 | Big data in perinatal epidemiology

To design interventions aimed at improving perinatal health, causal knowledge on the effects of exposures of interest (and underlying mechanisms) on perinatal outcomes is necessary. The field of perinatal health research has generally focused on answering causal research questions using randomised controlled trials (RCT).¹ However, conducting RCTs in this field to investigate the effect of certain exposure on perinatal health is often unrealistic due to practical and ethical considerations. For example, in the study of health inequalities, it is unfeasible to randomly assign the population to advantaged and disadvantaged socioeconomic conditions. Moreover, questions regarding the underlying mechanisms cannot be answered using RCTs.² In these situations, researchers have supported the concept of causal inference with observational (big) data.

Recent advances in computing power enabled the collection, linkage and processing of large amounts of data from multiple sources, that is big data. Big data can refer to datasets with a large number of observations (e.g., population registry data) or datasets with a large number of variables (e.g., genomics data).³ While big data are typically not collected for research purposes, it can contribute to health research through its potential to link health records with multiple datasets or by covering a large number of observations (often entire populations). In exchange for these advantages, big data present the challenges of being potentially incomplete, inaccurate and computationally intensive to process. Furthermore, the observational nature of big data represents challenges for causal inference.⁴

1.2 | Causal inference and the parametric g-formula

Answering causal questions with big data requires both high-quality data and robust statistical methods. The Neyman-Rubin potential outcomes framework,^{5,6} provides conceptual definitions and supports analytic methods for estimating causal effects from observational data.⁷ This approach uses counterfactuals (i.e., 'what-if' scenarios) to define causal effects.⁸

The parametric g-formula,⁹ a technique embedded in the potential outcomes framework, was first introduced in 1986 by Robins. However, its widespread application only became feasible with increasing computational power.¹⁰ This technique is recognised as a unified and flexible causal inference approach that allows for

Synopsis

Study Question

Is the relationship between neighbourhood socioeconomic status (SES) and small for gestational age (SGA) mediated by neighbourhood social environment?

What is already known

The link between neighbourhood SES and adverse perinatal outcomes has been well established. However, little is known about the underlying mechanisms (mediators). One of the potential mechanisms is neighbourhood's social environment. Big data in conjunction with robust methods could be used to identify such underlying mechanisms.

What this study adds

This paper provides a didactic illustration of the mediational g-formula, a robust and flexible approach that can be used to perform causal mediation analysis in a big data context. In the example, we found that the natural indirect effect accounted for 11.4% of the effect of neighbourhood SES on SGA.

designing custom interventions, a property available in a few other methods.^{11,12} The g-formula was originally proposed for applications in settings with confounders affected by previous exposure and can naturally be extended to mediation analysis.¹²

Mediation analysis evaluates the relative magnitude of pathways by which an exposure influences an outcome.^{13,14} The most utilised approach to perform mediation analysis is the Baron and Kenny, traditional, approach.¹⁵ The traditional mediation approach has important shortcomings as it is prone to bias when exposure and mediator interact and when the outcome is non-linear (e.g., dichotomous).¹⁶⁻¹⁹ The parametric g-formula extends mediation analysis to settings involving non-linearities and interactions,^{10,12} and its estimates are easily understandable population-averaged effects.¹¹ The g-formula uses parametric regression models to predict outcomes under hypothetical intervention scenarios (counterfactuals), which are used to estimate mediation effects via Monte Carlo simulation. The parametric g-formula is referred to as mediational g-formula when used

for causal mediation analysis. For simplicity, in the remaining of the manuscript, we will refer to this approach as the g-formula.

While applications of the g-formula for mediation analysis have been increasing in recent years, it remains underused among substantive researchers and didactic explanation for its implementation is still limited.¹² We provide a didactic demonstration of the implementation of the g-formula by means of an example from perinatal health inequalities research. The demonstration in this paper corresponds to a simple scenario and is meant to provide a gentle introduction to the potential outcomes framework and the use of the g-formula in R software.²⁰

2 | METHODS

2.1 | Illustrative example research question and dataset

Compelling evidence shows a consistent link between neighbourhood socioeconomic status (SES) and perinatal outcomes.^{21,22} Although this relationship has been well established, little is known about the underlying mechanisms. One of the hypothesised pathways in the literature is neighbourhood social environment.^{23,24} While SES refers to the economic conditions of a neighbourhood, social environment is defined as the relationships and processes that exist between its residents along with the social composition of a neighbourhood in terms of, for example life stage.²⁵⁻²⁷ In our example, we use the g-formula to investigate whether the relationship between neighbourhood SES and perinatal health is mediated by social environment.

2.1.1 | Outcome

This paper focuses on small for gestational age (SGA) as the outcome, defined as birthweight below the 10th centile for gestational age and sex, according to national reference curves.²⁸ Data from the Netherlands Perinatal Registry (Perined) were acquired for singleton births at gestational ages between 24+0 and 41+6 weeks between 1 January 2010 and December 2017. Perined contains high-quality information on perinatal outcomes and maternal characteristics. The perinatal registry was linked by Statistics Netherlands (CBS) to several individual-level sociodemographic registries. Unfortunately, only live births could be linked by CBS.

2.1.2 | Exposure

Neighbourhood SES was quantified using the Neighbourhood Status Score by the Netherlands Institute of Social Research.²⁹ The SCP Status Score is a validated relative indicator of neighbourhood SES computed using factor analysis to summarise into a single score the following three characteristics: (i) the percentage of residents with a

low income; (ii) the percentage of inhabitants without a paid job; and (iii) the percentage of inhabitants with a low education level.²⁹ More information is available in Appendix S1, file 1.

2.1.3 | Mediator

The measure for the mediator corresponds to the Social Environment Score from the neighbourhood liveability assessment ('Leefbaarometer') by the Netherlands Ministry of the Interior.³⁰ The Social Environment Score (range, -50-50), one of the dimensions of the Leefbaarometer, provides a single score based on the following indicators: residential stability (number of relocations), life stage diversity of households (e.g., single, couples and family households), population density and social cohesion (more information in Appendix S1, file 1). The score has shown good internal and external validity.³⁰ Information on other neighbourhood-level characteristics was obtained from CBS.³² All neighbourhood-level data were linked to birth records using the mother's residential postcode and year of birth.

To facilitate the explanation of the g-formula approach, exposure, outcome and mediator variables were dichotomised. To create the exposure categories, quintiles of neighbourhood SES were first calculated. The disadvantaged neighbourhood SES category corresponds to the lowest quintile and the non-disadvantaged category refers to the remaining quintiles, thus resulting in two categories. The same approach was taken for the social environment categories.

2.1.4 | Ethics approval

According to Dutch law (WMO), no formal ethical review was required. Perined provided approval (19.13) for this research project.

2.2 | Mediation analysis using the parametric mediational g-formula

2.2.1 | Counterfactuals

Under the potential outcomes framework, mediation analysis defines causal effects as the difference between two counterfactual outcomes.⁵ Counterfactuals can be thought of as what would have happened under alternative histories.⁸ Thus, a counterfactual outcome refers to the outcome value that *would be* observed whether the exposure *would be* set to a certain value. Let Y denote the outcome of interest (SGA) and SES the exposure of interest (neighbourhood SES), which can take the (observed) values $SES = 1$ (disadvantaged SES) or $SES = 0$ (non-disadvantaged SES). We use upper case SES to denote the observed values of SES. If the exposure *would be* set to disadvantaged SES ($ses = 1$), the counterfactual outcome would be denoted as $Y_{ses=1}$, and if the exposure *would be* set to advantaged SES ($ses = 0$) the counterfactual outcome would

be $Y_{ses=0}$. Lowercase *ses* is used to denote 'set' values of SES. The effect of the exposure is defined at the population level as the difference between these two counterfactual outcomes, that is $E[Y_{ses=0} - Y_{ses=1}]$. Since these are counterfactual outcomes under alternative exposure levels, only one would be factual (observed).⁷ However, through the g-formula and identification conditions (section 2.2.3), observational data can be used to extract information about the unobserved counterfactual outcome.

Adding a mediator makes the definitions of counterfactuals more complex.⁸ For each value of the exposure, there is a counterfactual value for the mediator and one for the outcome. Let *M* denote our mediator variable (social environment) where $M_{ses=1}$ and $M_{ses=0}$ would be the counterfactual values of the mediator under both potential exposure values. If the value for the exposure would be set to *ses* = 1 and the mediator would take on the value that would naturally be observed under *ses* = 1, that is $M_{ses=1}$, the counterfactual outcome would be denoted as $Y_{ses=1 M_{ses=1}}$. Similarly, if the exposure would be set to *ses* = 0, the counterfactual outcome would be $Y_{ses=0 M_{ses=0}}$. These so-called *nested counterfactual outcomes* are used to define the total and mediated effects.

2.2.2 | Total and mediated effects

The counterfactual mediation approach outlines a natural direct effect (NDE) and a natural indirect effect (NIE) that add up to the total effect (TE).³¹ These effects are defined in Table 1. As mentioned earlier, there are two counterfactual scenarios in our example: (i) setting the neighbourhood SES value to disadvantaged SES (*ses* = 1) and (ii) setting the neighbourhood SES value to non-disadvantaged SES (*ses* = 0). The TE is the difference in outcomes of changing the exposure value from *ses* = 1 to *ses* = 0 (from disadvantaged to non-disadvantaged), defined as $E[Y_{ses=0 M_{ses=0}} - Y_{ses=1 M_{ses=1}}]$.¹⁰ We refer to this change as a hypothetical intervention on the exposure where neighbourhood SES was improved.

The NIE, that is the effect that operates through the mediator (social environment), is interpreted as the effect of changing the mediator value from $M_{ses=0}$ to $M_{ses=1}$, while holding the exposure value

constant to *ses* = 1, that is $E[Y_{ses=1 M_{ses=0}} - Y_{ses=1 M_{ses=1}}]$. The NDE is the effect from changing the exposure from *ses* = 0 to *ses* = 1 and in both cases letting the value of the mediator be at their potential level as in $M_{ses=0}$, that is $E[Y_{ses=0 M_{ses=0}} - Y_{ses=1 M_{ses=0}}]$. As seen above, the nested counterfactual $Y_{ses=1 M_{ses=0}}$ is introduced to be able to define the mediation effects. Using this counterfactual, we can interpret the NIE as the observed effect of changing the mediator as if one had changed the exposure but without actually changing the exposure itself. Likewise, the NDE effect is the effect of changing the exposure, but keeping the mediator fixed at whatever level it would be had the exposure not been changed.¹⁴

2.2.3 | Causal diagram and identification assumptions

To give the total and mediation effects a causal interpretation, we must make certain identification assumptions: consistency, positivity and exchangeability.³¹ These identification assumptions, described in Table 2, are not exclusive of the counterfactual framework (or the g-formula), but this framework made them explicit.

The causal diagram in Figure 1 represents the hypothesised relationships between exposure, mediator, outcome and confounding variables. In our example, the models account for exposure–outcome confounders, that is individual-level characteristics (e.g., maternal age, parity, ethnicity, household income and education), and area-level average home value. Additionally, the models accounted for area-level percentage of non-western migrants, not only an exposure–mediator confounder (which influences SES and social environment), but also a mediator–outcome confounder as it is related to perinatal outcomes.³³ More information on the confounders included in the model is available in Appendix S1, file 2. A sensitivity analysis was conducted to assess the impact of women moving to another neighbourhood during (or shortly prior to) their pregnancy (Appendix S1, file 2).

In recent years, researchers have proposed the use of single-world intervention graphs (SWIGs) as a unification of causal diagrams and the counterfactual approach.³⁴ In these graphs, single worlds are represented, for example, a world where *ses* = 0 separate from the world where *ses* = 1, are fully represented. If our main question would be related to the estimation of the total effect and not the decomposition of it, SWIGs could be used in a straightforward manner. However, to address our mediation research question, the effects defined in Table 1 have cross-world references making the use of SWIGs not feasible. We refer the interested reader to the work of Richardson and colleagues for more guidance on the use of SWIGs.³⁴

2.2.4 | The g-formula procedure

The total and mediated effects (Table 1) were estimated following the g-formula steps in Table 3 (also addressed elsewhere¹²). In step

TABLE 1 Effect definitions used in causal mediation analysis

Effect	Definition
Total effect (TE)	$E[Y_{ses=0} - Y_{ses=1}] = E[Y_{ses=0 M_{ses=0}} - Y_{ses=1 M_{ses=1}}]$
Natural indirect effect (NIE)	$E[Y_{ses=1 M_{ses=0}} - Y_{ses=1 M_{ses=1}}]$
Natural direct effect (NDE)	$E[Y_{ses=0 M_{ses=0}} - Y_{ses=1 M_{ses=0}}]$

Note: Where *Y* refers to the outcomes and *ses* to the 'set' values of the exposure, neighbourhood SES, where *ses* = 1 refers to the disadvantaged counterfactual scenario, whereas *ses* = 0 denotes the non-disadvantaged scenario. *M* refers to the mediator, that is social environment. Given that the effect definitions used for the g-formula refer to differences in outcome means, the formulas shown above are in the difference scale.

1, the observed data were used to fit suitable regression models (underlying models) for mediator and outcome variables. These models included the individual and area-level confounders described

TABLE 2 Causal identification assumptions

Consistency
This condition connects the counterfactuals with observed outcomes by assuming that the nested counterfactuals will take the observed values when the treatment and mediator are actively set to the values they would naturally have had in the absence of an intervention. ¹⁴ To meet the consistency assumption, the exposure and mediator must be well defined and there must not be multiple versions of either of them.
Positivity
It assumes that for every combination of covariates the probability of observing any of the exposure values is nonzero. Furthermore, it assumes that for every combination of covariates and exposure values the probability of observing any of the mediator values is also nonzero. ¹⁴
Exchangeability
It assumes that one could exchange groups without changing the outcome of the study. Groups would not be exchangeable in settings where there is selection bias and/or confounding. Selection on certain characteristics, for example selection on live births, can lead to bias due to conditioning on a collider, which opens a non-causal path between exposure and outcome. ⁴⁷ When selection on these characteristics is unavoidable, this bias can be reduced by adjusting for common causes of the collider variable and the outcome. ⁴⁸ In our illustrative example, our dataset contains live births only. The underlying models were adjusted for known common causes of stillbirths and SGA, that is maternal age, parity, education and income. In a sensitivity analysis, we also included maternal lifestyle factors and pre-existent conditions (see Appendix S1, file 2). Confounders are defined as covariates that are expected to be common causes of, for example the exposure and the outcome. Thus, to interpret the total effect as causal, we assume no uncontrolled confounding for the exposure–outcome relationship. Additionally, in mediation analysis, to identify the direct and indirect effects it is also necessary to account for confounding for the exposure–mediator and mediator–outcome relationships, including mediator–outcome confounding affected by the exposure (see Appendix S1, file 2). ⁴⁹

in section 2.2.3. The model for the outcome additionally included the mediator (social environment). These models may also include exposure–mediator interactions if required (we refer to other work for guidance³⁵). Parametric models, that is logistic regression, were used for the outcome and mediator. The odds ratios for the underlying model for the outcome can be found in Table S2. The g-formula has the benefit that in big data settings where there is a large number of candidate confounders, it can easily be combined with machine learning algorithms (e.g. the superlearner³⁶) to perform variable selection for the underlying models.³⁷

The model parameters from the first step are employed to obtain predicted probabilities for mediator and outcome variables. These predicted probabilities are used in step 2 for the Monte Carlo (MC) simulation where a dataset that resembles the observed data, natural course (NC), was simulated by keeping neighbourhood SES at its observed values. The mediator is simulated first, and then, its values are used in the model for the outcome. Using the same procedure, in step 3, datasets for the two counterfactual scenarios are simulated by fixing the exposure to the corresponding value ($ses = 0$ or $ses = 1$). Furthermore, in step 4 a mediation scenario, that is $Y_{ses=1}, M_{ses=0}$ (Table 1), was simulated to be able to estimate the NIE and NDE.

In step 5, the mean values for mediator and outcome are saved for all simulated scenarios, which represent the proportion of births with a given outcome (or mediator) in each scenario. The simulation process involves (randomly) drawing values from probability distributions and the exact values differ across draws. This variability is known as Monte Carlo error,³⁸ which can be reduced by repeating the simulation process (and mean values calculation) multiple times, that is iteratively (step 6). The number of MC iterations must be enough (30 iterations in our case) to have stable estimates, which can be checked with the R package *cfdecomp*.³⁹

In step 7, the mean outcome values saved across MC iterations were used to estimate the TE, NIE and NDE based on definitions from Table 1. These values are the point estimates of the effects. In step 8, the 95% confidence intervals of the effects are obtained via bootstrapping, that is sampling with the replacement of the same number of individuals in the dataset. In our example, we sampled clusters (neighbourhoods) instead of individuals to account for a

FIGURE 1 Conceptual directed acyclic graph for the relationship between the exposure (neighbourhood SES) and the outcome (small for gestational age) via a mediator variable (neighbourhood social environment).

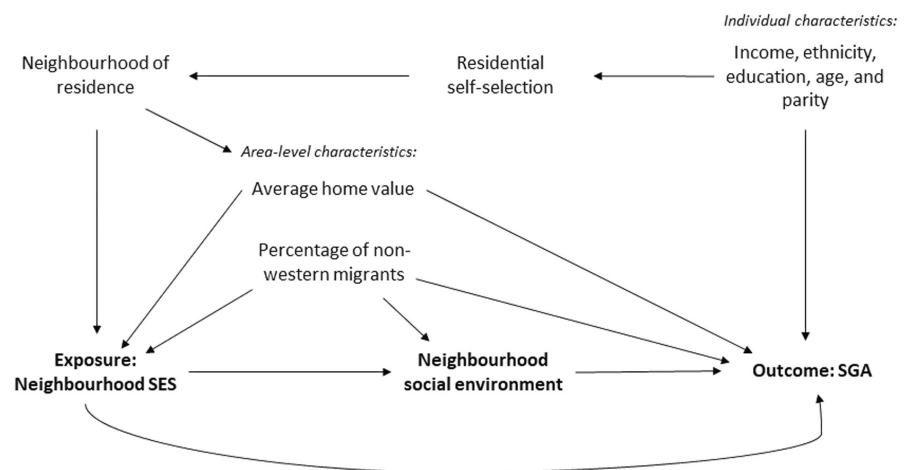


TABLE 3 Parametric mediational g-formula procedure

G-formula step-by-step procedure

1. Use the original data to fit the underlying models, that is suitable parametric models for mediator and outcome, that is a logistic regression model if the outcome is a dichotomous variable. These models include the confounders and the model for the outcome also includes the mediator. Exposure–mediator interactions are possible.
2. Use the model parameters from step 1 to predict probabilities for mediator and outcome. The predicted probabilities are used to draw new values from the probability distribution assumed when modelling mediator and outcome (e.g., binomial distribution for dichotomous variables) to simulate a new dataset without intervention, that is the natural course scenario (NC). The mediator is simulated first and then its values are used in the model for the outcome.
3. Next, using the dataset from step 1, simulate two datasets under the two counterfactual scenarios (CF). This is done by setting (fixing) the exposure to the corresponding value for each CF ($ses = 0$ or $ses = 1$) and following the same procedure as in step 2.
4. Additional to the CF scenarios a mediation scenario is simulated where neighbourhood SES is intervened as in $ses = 1$ but the mediator values will be derived from the $ses = 0$ scenario. This scenario is later used for the estimation of (natural) direct and indirect effects.
5. Save the average values for mediators and outcomes over the simulated scenarios. For dichotomous outcomes, the averages correspond to the proportion of cases with a given outcome (or mediator) in each scenario.
6. The simulations and calculation of the average values (steps 2–5) are repeated J times, where J is a number of iterations sufficient to produce stable estimates. This can be checked by producing stability plots, for example with the *cfdecomp* R package (a tutorial available via this link).³⁹
7. The average of the J (Monte Carlo) iterations is used to obtain the point estimates of the effects. The effects are estimated based on the definitions of Table 1: the total effect (TE) is obtained from the difference between the average values of the two counterfactual scenarios. The mediation scenario is used to obtain the natural direct effect (NDE) and the natural indirect effect (NIE). The NDE is the difference in average values between the CF where $ses = 0$ and the mediation scenario. Last, the difference between the average values for the mediation scenario and the CF where $ses = 1$ is the NIE.
8. The steps above are repeated K times to produce bootstrap confidence intervals for the effects, the estimated effect values are saved for each bootstrap iteration, where K is a large enough value (200+) to produce stable estimates (use stability plots). The confidence intervals are obtained as the 2.5th and 97.5th quantiles of the distribution.
9. The comparison between the observed means and the means under the NC (no intervention) scenario is used as a check against gross model misspecification. If the NC predictions are not close to observed values, then models for outcome and/or mediators are likely to be incorrectly specified.

Note: R code available in public repository (link in Appendix S1, file 7).

multilevel data structure. Similarly, to the MC iterations, we ran a sufficient number of bootstrap iterations to obtain stable estimates (250 iterations). The computation time of the g-formula depends on the number of observations. In big data settings, as these numbers

TABLE 4 G-formula mediation effects of neighbourhood SES improvement from disadvantaged to an advantaged category on small for gestational age births (percentage reduction)

Effect	Mean (95% confidence interval)
Total effect (TE)	6.3% (5.2, 7.5)
Natural indirect effect (NIE)	0.7% (0.6, 0.9)
Natural direct effect (NDE)	5.6% (5.1, 6.1)

increases, researchers may consider parallel computing or taking a random subset of the sample to perform the simulation.^{12,40}

The g-formula is prone to bias due to misspecification of the underlying models either by misspecifying the functional form (for mediator or outcome models) or by omitting confounders. In step 9, we performed a check against gross model misspecification where we compared the observed means (for the outcome and mediators) and the means under the simulated NC scenario.^{10,12} If the means for the NC scenario are not close to observed values, then outcome and/or mediator models are likely misspecified.

For interpretability, results for the effects are presented in relative terms, that is the percentage change in the proportion of births with a given outcome (see Appendix S1, file 3 for further explanation). To assess the extent to which the total effect of the exposure on the outcome operates through the mediator, the proportion mediated can be calculated. As pointed out in previous work,⁴¹ when the effects are used on the difference scale (i.e., additive scale; as in Table 1), the proportion mediated simply corresponds to the ratio of the natural indirect effect to the total effect, that is $PM = NIE/TE$.

3 | RESULTS

After the exclusion of non-linked births, multiple births, births with gestational age below 24+0 weeks or above 41+6 weeks, and cases with missing information (<2%), there were 1,217,626 births available for analysis. Due to the small percentage of missing data, no data imputation was conducted. Population summary characteristics and a flow diagram can be found in Figure S1 and Table S1. The natural course scenario yielded similar mean values to the ones from the observed dataset, (Table S5) meaning that gross model misspecification is unlikely to be an issue.

Table 4 shows the effects estimated using the g-formula. The absolute values for these effects are shown in Table S4. At the population level, a hypothetical improvement in neighbourhood SES from disadvantaged to non-disadvantaged resulted in a 6.3% (95% CI 5.2, 7.5) relative reduction in the proportion of SGA, that is the total effect. This effect was decomposed into direct and indirect effects as observed in Table 4. As a measure of the magnitude of the mediation, the proportion mediated was computed as specified in the previous section (please see Appendix S1, file 3 for more information). Thus, the natural indirect effect accounted for

11.4% (95% CI 9.2%, 13.6%) of the total effect of neighbourhood SES on SGA.

4 | COMMENT

4.1 | Principal findings

In this didactic demonstration of the mediational *g*-formula, we investigated whether neighbourhood social environment mediates the relationship between neighbourhood SES and SGA. The results showed that a hypothetical improvement in neighbourhood SES from disadvantaged to non-disadvantaged resulted in a 6.3% reduction in SGA births and that 11.4% of this total effect is mediated by neighbourhood social context.

4.2 | Strengths of the study

Regarding the analysis performed in the illustrative example, a first strength corresponds to the ability to link several high-quality national-level datasets, leading to information on over 1.2 million births available for analysis. The use of the *g*-formula allowed us to investigate one of the potential pathways driving the exposure–outcome relationship of interest in a setting with a dichotomous outcome, helping to overcome potential non-collapsibility issues.

4.3 | Limitations of the data

Foremost, our study is based on registry data, which makes it rather difficult to observe all potential confounders. For example, there might be unobserved individual-level characteristics, such as preferences, that influence both exposures to certain neighbourhood environments and perinatal health. Another limitation is that our dataset consisted of live births, which might lead to selection bias by conditioning on a collider. While we have followed a strategy to reduce this bias, this scenario may result in a violation of the exchangeability condition. A separate issue may come from the exposure and mediator variables being dichotomised. Although the categories are well defined, when dichotomising, for example the mediator, one value of the mediator measure corresponds to multiple values of the true mediator resulting in a violation of the consistency assumption.⁴² However, it has been argued that mediation effects can be interpreted even if the consistency assumption does not hold.⁴²

Another potential concern is measurement error. It is likely that, for example, the measure for the mediator is imprecise. Previous work has shown that, in the context of mediation analysis, measurement error can affect the direct and indirect effects resulting in bias towards the null for the indirect effect and bias away from the null for the direct effect.⁴³ Thus, it is likely that the proportion mediated is underestimated.⁴³ Finally, the assessment of mediation involves

an aspect of temporality where the exposure should be measured before the mediator, and this in turn is measured before the outcome. These conditions are relevant to prevent reverse causation and overadjustment. The model presented in the DAG reflects theoretical considerations in the study of neighbourhood health effects where the social environment is seen as a pathway for the effect of SES on health.²³ However, we cannot rule out that in this case social environment may also influence SES. To avoid this issue, ideally, one must use a measure of the exposure that temporally precedes the measure of the mediator as done in Burgos Ochoa et al.⁴⁴ However, in our real-life example, this was not feasible as only two reporting years for the mediator (2014 and 2017) were available.

The application of the *g*-formula approach also has shortcomings. The validity of the *g*-formula estimation is dependent on the validity of the underlying models used to create the simulated data. In the example, we found that observed and natural course means were practically equivalent. However, this check against gross model misspecification cannot fully rule out the presence of milder forms of this issue.¹⁰ Another challenge is that the *g*-formula is very computational-intensive.¹² While there are solutions available for very large datasets (section 2.2.4), this remains a concern for researchers in settings with computational power constraints.

4.4 | Interpretation

In this study, we observed that the hypothetical improvement in neighbourhood SES led to a 6.3% reduction in the proportion of SGA, which corresponds to a small but meaningful effect, particularly when compared to effect estimates found in previous studies.^{24,45} Regarding the proportion mediated, we observed that neighbourhood social environment accounted for 11.4% of the effect of neighbourhood SES on SGA. While this is a meaningful quantity, the results point that a large share of the effect remains unexplained and there is a need for further research on other potential mediators, for example crime rates or environmental pollution.

The *g*-formula, being a flexible approach, can be used in various scenarios in perinatal epidemiology. The *g*-formula can accommodate all types of outcomes of interest, for example continuous outcomes, such as birthweight, or survival outcomes like neonatal mortality. Furthermore, researchers in perinatal epidemiology are interested often in multiple underlying mechanisms, which might interact and even influence each other. The *g*-formula can handle multiple mediators at once without making the stringent assumption of them not being interrelated, as in other approaches.⁴⁶ Finally, longitudinal designs are frequently used in this field, and with them comes the challenge of time-varying exposures and confounders, and the issue of adjusting for confounders affected by previous exposure. The *g*-formula is suitable for these challenging settings.¹² Given the wide variety of potential applications, the *g*-formula can be considered a promising analytical approach in the field of perinatal health research.

5 | CONCLUSIONS

The mediational g-formula is a flexible approach to performing causal mediation analysis that is suited for big data contexts in perinatal epidemiology. This approach overcomes many of the limitations of traditional mediation analysis methods.

AUTHOR CONTRIBUTIONS

JVB and LCMB obtained funding for the study. LBO, JVB, LM and LCMB conceived the study. LBO and LCMB analysed the data. All authors were involved in interpreting the data. LBO and LM wrote the draft paper and LCMB supervised the writing. JVB and AJMV provided additional input at the writing stage. All authors read and approved the final version of the manuscript.

FUNDING INFORMATION

This work was supported by the Erasmus Initiative Smarter Choices for Better Health.

CONFLICT OF INTEREST

None declared.

DATA AVAILABILITY STATEMENT

This study is based on registry data from the Dutch Perinatal Registry (Perined) and microdata from Statistics Netherlands (<https://www.cbs.nl/nl-nl>). Access to linked electronic health and sociodemographic records requires approval from Statistics Netherlands following the procedure on their website.

ORCID

Lizbeth Burgos Ochoa  <https://orcid.org/0000-0002-8379-2749>

Loes C. M. Bertens  <https://orcid.org/0000-0003-0897-0709>

REFERENCES

- Snowden JM, Bovbjerg ML, Dissanayake M, Basso O. The curse of the perinatal epidemiologist: inferring causation amidst selection. *Curr Epidemiol Rep*. 2018;5:379-387.
- Lee H, Herbert RD, Lamb SE, Moseley AM, McAuley JH. Investigating causal mechanisms in randomised controlled trials. *Trials*. 2019;20:524.
- de Mauro A, Greco M, Grimaldi M. A formal definition of big data based on its essential features. *Library Rev*. 2016;65:122-135.
- Dolley S. Big data's role in precision public health. *Front Public Health*. 2018;6:68.
- Rubin DB. Estimating causal effects of treatments in randomized and nonrandomized studies. *J Educ Psychol*. 1974;66:688-701.
- Rubin DB. Causal inference using potential outcomes: Design, modeling, decisions. *J Am Stat Ass*. 2005;100(469):322-331.
- Hernán MA, Robins JM. *Causal Inference: What if?*. Chapman & Hall/CRC; 2020.
- Rudolph KE, Goin DE, Paksarian D, Crowder R, Merikangas KR, Stuart EA. Causal mediation analysis with observational data: considerations and illustration examining mechanisms linking neighborhood poverty to adolescent substance use. *Am J Epidemiol*. 2019;188:598-608.
- Robins J. A new approach to causal inference in mortality studies with a sustained exposure period—application to control of the healthy worker survivor effect. *Math Model*. 1986;7:1393-1512.
- 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. *Epidemiology*. 2014;25:889.
- Pitkänen J, Bijlsma MJ, Remes H, Aaltonen M, Martikainen P. The effect of low childhood income on self-harm in young adulthood: mediation by adolescent mental health, behavioural factors and school performance. *SSM-Population Health*. 2021;13:100756.
- Wang A, Arah OA. G-computation demonstration in causal mediation analysis. *Eur J Epidemiol*. 2015;30:1119-1127.
- VanderWeele TJ. Mediation analysis: a practitioner's guide. *Annu Rev Public Health*. 2016;37:17-32.
- Lange T, Hansen KW, Sørensen R, Galatius S. Applied mediation analyses: a review and tutorial. *Epidemiol Health*. 2017;39:39.
- Baron RM, Kenny DA. The moderator-mediator variable distinction in social psychological research: conceptual, strategic, and statistical considerations. *J Pers Soc Psychol*. 1986;51:1173-1182.
- Kaufman JS, MacLehose RF, Kaufman S. A further critique of the analytic strategy of adjusting for covariates to identify biologic mediation. *Epidemiol Perspect Inn*. 2004;1:1-13.
- Rijnhart JJM, Lamp SJ, Valente MJ, MacKinnon DP, Twisk JWR, Heymans MW. Mediation analysis methods used in observational research: a scoping review and recommendations. *BMC Med Res Methodol*. 2021;21:1-17.
- Cole SR, Hernán MA. Fallibility in estimating direct effects. *Int J Epidemiol*. 2002;31:163-165.
- MacKinnon DP, Lockwood CM, Brown CH, Wang W, Hoffman JM. The intermediate endpoint effect in logistic and probit regression. *Clin Trials*. 2007;4:499-513.
- R-Core-Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. 2021.
- Ncube CN, Enquobahrie DA, Albert SM, Herrick AL, Burke JG. Association of neighborhood context with offspring risk of preterm birth and low birthweight: a systematic review and meta-analysis of population-based studies. *Soc Sci Med*. 2016;153:156-164.
- Vos AA, Posthumus AG, Bonsel GJ, et al. Deprived neighbourhoods and adverse pregnancy outcomes: a systematic review and meta-analysis. *Acta Obstet Gynecol Scand*. 2014;93:181-740.
- van Ham M, Manley D, Bailey N, Simpson L, MacLennan D. Understanding neighbourhood dynamics: new insights for neighbourhood effects research. *Understanding Neighbourhood Dynamics*. Springer; 2012:1-21.
- Morenoff JD. Neighborhood mechanisms and the spatial dynamics of birth weight. *Am J Sociol*. 2003;108:976-1017.
- Kepper MM, Myers CA, Denstel KD, Hunter RF, Guan W, Broyles ST. The neighborhood social environment and physical activity: a systematic scoping review. *Int J Behav Nutr Phys Act*. 2019;16:1-14.
- Suglia SF, Shelton RC, Hsiao A, Wang YC, Rundle A, Link BG. Why the neighborhood social environment is critical in obesity prevention. *J Urban Health*. 2016;93:206-212.
- McNeill LH, Kreuter MW, Subramanian SV. Social environment and physical activity: a review of concepts and evidence. *Soc Sci Med*. 2006;63:1011-1022.
- Hoftiezer L, Hof MHP, Dijks-Elsinga J, Hogeveen M, Hukkelhoven CWPM, vanLingen RA. From population reference to national standard: new and improved birthweight charts. *Am J Obstet Gynecol*. 2019;220(4):383. e1-383. e17.
- Social and Cultural Planning Office (SCP). Socio-Economic Status by postcode area [Internet]. 2019. Available from: <https://bronn.en.zorggegevens.nl/Bron?naam=Sociaal-Economische-Status-perpostcodegebied>
- Mandemakers J, Leidelmeijer K, Burema F, Halbersma R, Middeldorp M, Veldkamp J. Leefbaarometer 3.0 Instrumentontwikkeling [Internet]. 2021. Available from: <https://>

- www.leefbaarometer.nl/resources/LBM3Instrumentontwikkeling.pdf
31. Robins JM, Greenland S. Identifiability and exchangeability for direct and indirect effects. *Epidemiology*. 1992;3:143-155.
 32. Central Bureau for Statistics. Wijk en Buurt Statistieken. 2020.
 33. Schölicherich VLN, Erdem Ö, Borsboom G, et al. The association of neighborhood social capital and ethnic (minority) density with pregnancy outcomes in The Netherlands. *PLoS One*. 2014;9:e95873.
 34. Richardson TS, Robins JM. *Single World Intervention Graphs (SWIGs): A Unification of the Counterfactual and Graphical Approaches to Causality*. Center for the Statistics and the Social Sciences. University of Washington Series. Working Paper 2013; 2013:128.
 35. Vander Weele T. When to include an exposure-mediator interaction. *Explanation in Causal Inference: Methods for Mediation and Interaction*. Oxford University Press; 2015:45-47.
 36. van der Laan MJ, Polley EC, Hubbard AE. Super learner. *Stat Appl Genet Mol Biol*. 2007;6:25.
 37. Snowden JM, Rose S, Mortimer KM. Implementation of G-computation on a simulated data set: demonstration of a causal inference technique. *Am J Epidemiol*. 2011;173:731-738.
 38. Koehler E, Brown E, Haneuse SJ-PA. On the assessment of Monte Carlo error in simulation-based statistical analyses. *Am Statist*. 2009;63:155-162.
 39. Bijlsma MJ, Sudharsanan N. Cfdecomp: counterfactual decomposition: MC integration of the G-formula. R Package Version 0.4.0. 2021.
 40. Daniel RM, de Stavola BL, Cousens SN. Gformula: estimating causal effects in the presence of time-varying confounding or mediation using the g-computation formula. *Stata J*. 2011;11:479-517.
 41. Vanderweele TJ. *Explanation in causal inference. Methods for mediation and interaction*. Oxford University Press; 2015.
 42. Vanderweele TJ. Mediation analysis with multiple versions of the mediator. *Epidemiology*. 2012;23:454-463.
 43. VanderWeele TJ, Valeri L, Ogburn EL. The role of measurement error and misclassification in mediation analysis. *Epidemiology*. 2012;23:561.
 44. Burgos Ochoa L, Bijlsma MJ, Steegers EAP, Been J, LCM B. Does Neighbourhood crime mediate the relationship between Neighbourhood socioeconomic status and birth outcomes? An application of the mediational G-formula. *Eur J Pub Heal*;32. <https://doi.org/10.1093/eurpub/ckac130.041>. In press.
 45. Masi CM, Hawkey LC, Piotrowski ZH, Pickett KE. Neighborhood economic disadvantage, violent crime, group density, and pregnancy outcomes in a diverse, urban population. *Soc Sci Med*. 2007;65:2440-2457.
 46. VanderWeele T, Vansteelandt S. Mediation analysis with multiple mediators. *Epidemiol Methods*. 2014;2:95-115.
 47. Infante-Rivard C, Cusson A. Reflection on modern methods: selection bias—a review of recent developments. *Int J Epidemiol*. 2018;47:1714-1722.
 48. Neophytou AM, Kioumourtzoglou M-A, Goin DE, Darwin KC, Casey JA. Educational note: addressing special cases of bias that frequently occur in perinatal epidemiology. *Int J Epidemiol*. 2021;50:337-345.
 49. Greenland S, Robins JM. Identifiability, exchangeability and confounding revisited. *Epidemiologic Perspectives & Innovations*. 2009;6:1-9.

SUPPORTING INFORMATION

Additional supporting information can be found online in the Supporting Information section at the end of this article.

How to cite this article: Ochoa LB, van der Meer L, Waelpu AJM, Been JV, Bertens LCM. Neighbourhood-related socioeconomic perinatal health inequalities: An illustration of the mediational g-formula and considerations for the big data context. *Paediatr Perinat Epidemiol*. 2023;00:1-9. doi:[10.1111/ppe.12954](https://doi.org/10.1111/ppe.12954)