

Uncertainty reduction and environmental justice in air pollution epidemiology: the importance of minority representation

Mariana Alifa¹, Stefano Castruccio¹, Diogo Bolster¹, Mercedes A Bravo², and Paola Crippa¹

¹University of Notre Dame

²Duke University

June 1, 2023

Abstract

Ambient air pollution is an increasing threat to society, with rising numbers of adverse outcomes and exposure inequalities across the globe. Reducing uncertainty in health outcomes models and exposure disparity studies is therefore essential to develop policies effective in protecting the most affected places and populations. This study uses the concept of information entropy to study tradeoffs in mortality uncertainty reduction from increasing input data of air pollution versus health outcomes. We study a case scenario for short-term mortality from fine particulate matter (PM2.5) in North Carolina for 2001-2016, employing a case-crossover design with inputs from an individual-level mortality dataset and high-resolution gridded datasets of PM2.5 and weather covariates. We find a significant association between mortality and PM2.5, and the information tradeoffs indicate that in this case increasing information from mortality may reduce model uncertainty at a faster rate than increasing information from air pollution. We also find that Non-Hispanic Black (NHB) residents tend to live in relatively more polluted census tracts, and that the mean PM2.5 for NHB cases in the mortality model is significantly higher than that of Non-Hispanic White (NHW) cases. The distinct distribution of PM2.5 for NHB cases results in a relatively higher information value, and therefore faster uncertainty reduction, for new NHB cases introduced into the mortality model. This newfound influence of exposure disparities in the rate of uncertainty reduction highlights the importance of minority representation in environmental research as a quantitative advantage to produce more confident estimates of the true effects of environmental pollution.

Hosted file

963027_0_art_file_11004548_rpvxvy.docx available at <https://authorea.com/users/620448/articles/644451-uncertainty-reduction-and-environmental-justice-in-air-pollution-epidemiology-the-importance-of-minority-representation>

Hosted file

963027_0_supp_11004549_rgvxvy.docx available at <https://authorea.com/users/620448/articles/644451-uncertainty-reduction-and-environmental-justice-in-air-pollution-epidemiology-the-importance-of-minority-representation>

26 **Abstract**

27 Ambient air pollution is an increasing threat to society, with rising numbers of adverse
28 outcomes and exposure inequalities across the globe. Reducing uncertainty in health outcomes
29 models and exposure disparity studies is therefore essential to develop policies effective in
30 protecting the most affected places and populations. This study uses the concept of information
31 entropy to study tradeoffs in mortality uncertainty reduction from increasing input data of air
32 pollution versus health outcomes. We study a case scenario for short-term mortality from fine
33 particulate matter (PM_{2.5}) in North Carolina for 2001-2016, employing a case-crossover design
34 with inputs from an individual-level mortality dataset and high-resolution gridded datasets of
35 PM_{2.5} and weather covariates. We find a significant association between mortality and PM_{2.5}, and
36 the information tradeoffs indicate that in this case increasing information from mortality may
37 reduce model uncertainty at a faster rate than increasing information from air pollution. We also
38 find that Non-Hispanic Black (NHB) residents tend to live in relatively more polluted census
39 tracts, and that the mean PM_{2.5} for NHB cases in the mortality model is significantly higher than
40 that of Non-Hispanic White (NHW) cases. The distinct distribution of PM_{2.5} for NHB cases
41 results in a relatively higher information value, and therefore faster uncertainty reduction, for
42 new NHB cases introduced into the mortality model. This newfound influence of exposure
43 disparities in the rate of uncertainty reduction highlights the importance of minority
44 representation in environmental research as a quantitative advantage to produce more confident
45 estimates of the true effects of environmental pollution.

46 **1. Introduction**

47 Air pollution is an increasing threat to today’s society. Data from the Global Burden of
48 Disease study ranked ambient pollution from PM_{2.5} as the 5th leading global mortality risk factor
49 in 2015, causing 4.2 million deaths and 103.1 million disability-adjusted life years due to health
50 impacts such as lung cancer, lower respiratory infection, chronic obstructive pulmonary disease,
51 cerebrovascular disease, and ischemic heart disease (Cohen et al., 2017). A recent update for this
52 study (Fuller et al., 2022) reports a rise in ambient pollution attributable deaths to 4.5 million in
53 2019, a 7% increase since 2015 and a 66% increase since 2000, revealing that, despite increased
54 awareness and attempts at remediation of this problem, our efforts have so far been insufficient
55 in protecting society from the harms of ambient pollution.

56 The United States stands out as a successful case of continued efforts to curb air pollutant
57 emissions. The Clean Air Act required in 1970 that the Environmental Protection Agency (EPA)
58 set National Ambient Air Quality Standard (NAAQS) for “criteria pollutants” and establish a
59 network of ambient pollution monitoring stations to assess compliance to these standards. The
60 first NAAQS specifically for PM_{2.5} was issued in 1997 (once monitors were advanced enough to
61 measure particles of this size), setting the standard for annual mean concentration at 15 µg/m³
62 (EPA, 1997). However, subsequent findings of harmful health effects at air pollution
63 concentrations that blend into background levels have prompted the continual lowering of
64 NAAQS (McClellan, 2002). The standard for PM_{2.5} was lowered to 12 µg/m³ in 2012 (EPA,
65 2013), and a proposal issued in January of 2023 is now currently underway to further lower the
66 NAAQS to 9-10 µg/m³ (EPA, 2023). Although these nationwide measures have been effective in
67 reducing overall levels of air pollution, they have not been as successful in curbing demographic
68 and socioeconomic inequalities in relative exposure (Colmer et al., 2020; Liu et al., 2021).

69 Extensive research has found demographic and/or socioeconomic disparities in exposure
70 to PM_{2.5} and other air pollutants across different regions of the world (Hajat et al., 2015). In the
71 United States, multiple studies have found that people of color have been systematically exposed
72 to higher levels of air pollution (Colmer et al., 2020; Liu et al., 2021; Tessum et al., 2021). These
73 racial disparities are not only found across different income levels, urbanicity levels, and
74 emission types (Liu et al., 2021; Tessum et al., 2021), but they have also persisted despite the
75 nationwide decreasing trend in air pollution seen in the last four decades, with studies identifying
76 that the relatively most polluted census tracts in present day are largely the same census tracts
77 that were most polluted in the 80s and the 90s (Colmer et al., 2020; Liu et al., 2021).

78 In light of this lack of progress in addressing both air pollution-related health outcomes at
79 the global level and pollution exposure disparities at the national level, it is essential to develop
80 policies that will effectively target the places and populations most affected by ambient air
81 pollution. However, one of the multiple challenges to effective policy is the uncertainty affecting
82 ambient pollution health impact assessments (HIAs) used to guide AQ standards from local and
83 national (EPA, 2019; EU, 2008) to global (WHO, 2006) levels. These studies integrate multiple
84 sources of information such as, among others, air pollution concentrations and related population
85 exposure, physiological responses to pollution exposure, and their variation by individual-level
86 factors (such as gender, age, body mass, race, etc.) as well as residential factors (such as
87 proximity to water bodies or green spaces). Each of these sources of information involved in the
88 air pollution HIA may introduce several different kinds of uncertainty into the final assessment
89 model (Nethery & Dominici, 2019).

90 Among the many possible sources of uncertainty in HIAs, this study focuses on
91 uncertainty stemming from incomplete knowledge of the pollution and/or health impact

92 scenarios, caused by data scarcity in the input information. When there is a recognized scarcity
93 in observational data precluding the full characterization of the pollution-exposure-effects
94 scenario, action can be taken to augment the available input datasets to increase our knowledge
95 of the problem and gain confidence in the results of the final assessment. Solutions to the
96 problem of data scarcity have been indeed addressed extensively in both the air pollution and the
97 epidemiology fields.

98 Air pollution research has proposed different approaches to data assimilation for better
99 risk characterization, mainly by supplementing ground observations from official monitoring
100 stations (for example, those from the United States' Environmental Protection Agency, EPA)
101 with other sources of data, such as citizen-science observations (Bonas & Castruccio, 2021; Shen
102 et al., 2021), satellite observations of atmospheric and aerosol properties (Van Donkelaar et al.,
103 2021; Van Donkelaar et al., 2015; Zani et al., 2020), chemical transport models, or CTMs (Giani,
104 Anav, et al., 2020; Giani, Castruccio, et al., 2020), and/or dispersion models (Bates et al., 2018).
105 In cases where ground-based pollution data is sparse, CTMs able to reproduce monitored
106 pollutant concentrations have also been used to make robust assessments of the region's
107 pollution risks (Mead et al., 2018). Therefore, several studies have focused on localized
108 downscaling of existing CTMs to achieve finer resolution in areas of interest (Tessum et al.,
109 2017) or in the implementation of higher-resolution CTMs for a more accurate representation of
110 meteorological, chemical and aerosol properties (Crippa et al., 2019).

111 Previous work has also focused on assessing epidemiological uncertainty. For example,
112 meta-analyses of epidemiological studies combine multiple previous studies' results for
113 robustness (Atkinson et al., 2014; Pope et al., 2020). Another approach (Burnett et al., 2014)
114 developed an integrated exposure-response model by combining epidemiological data from

115 multiple PM_{2.5} sources, such as ambient air pollution, active and second hand tobacco smoke,
116 and household solid cooking fuel. A recent study (Coffman et al., 2020) derived distributions
117 from existing epidemiological data to model uncertainty in the exposure-response curve at low
118 levels of PM_{2.5}, for which data is usually sparse. Other studies have performed disaggregation of
119 exposure data with the goal of improving health effect estimation in future epidemiological
120 studies (Beckx et al., 2009; Breen et al., 2020).

121 Data scarcity in air pollution epidemiology studies also has environmental justice
122 implications. Studies of air pollution epidemiology have been traditionally based on ambient air
123 pollution monitoring data from the US Environmental Protection Agency (EPA), resulting in an
124 urban bias in the assessment (Bell et al., 2004; Dominici et al., 2006) since the EPA prioritizes
125 monitor placements in population-dense areas (Bravo et al., 2012; Miranda et al., 2011). Even
126 within relatively-urbanized counties, minority populations have been found to live closer to
127 sources of air pollution but further away from monitoring stations (Stuart et al., 2009). Recent
128 research has therefore leveraged the use of satellite data, land use regression, and air quality
129 models to expand and diversify the spatial area and thus, population, for which PM_{2.5} exposures
130 and health effects can be estimated (Ha et al., 2014; Hyder et al., 2014; Kloog et al., 2012; Qian
131 et al., 2019).

132 Although the problem of data scarcity has been extensively studied as it relates to air
133 pollution, epidemiology, and environmental justice, there remains a need for more
134 interdisciplinary research linking the findings from all these fields under a single framework. We
135 began addressing this need in a previous study (Alifa et al., 2022) where we adapted a
136 methodology proposed in the hydrology field (De Barros & Rubin, 2008; De Barros et al., 2009)
137 to create a novel framework that identifies the most efficient pathway to reduce uncertainty in

138 estimates of air pollution-associated health risks. The studies in hydrology (De Barros & Rubin,
139 2008; De Barros et al., 2009) had explored the concept of uncertainty tradeoffs in the modeling
140 of the health effects of groundwater contaminants combining the concept of information entropy
141 with Bayesian inference methods; Our subsequent study (Alifa et al., 2022) adapted this
142 framework for frequentist inference to study the effect of data increase on the reduction of air
143 pollution mortality uncertainty, measured through the metric of information entropy, and
144 visualize the tradeoffs in the resulting uncertainty of the mortality model depending on the kind
145 of input data gained. The two cases presented in that study (Alifa et al., 2022), one with artificial
146 data for PM_{2.5} and mortality data used in a long-term exposure model, and one with real time-
147 series data used in a short-term exposure model, demonstrated the applicability of the method for
148 aiding stakeholders in choosing the most efficient pathway for HIA uncertainty reduction when
149 limited resources (e.g. time, money, computational power) prevent them from investing in
150 improvements for both pollution and health outcomes data.

151 We now seek to explore this framework further by applying it to a more complex case
152 scenario involving spatio-temporal data. We use a case-crossover model design (Jaakkola, 2003)
153 to investigate the association of short-term PM_{2.5} exposure with mortality in North Carolina for
154 the years 2001-2016, through the use of individual-level mortality data and high-resolution
155 gridded datasets of PM_{2.5} and weather covariates. This study aims to not only illustrate the
156 usefulness of our information entropy tradeoff methodology to generate more robust impact
157 assessments, but also to gain new knowledge of the influence of socio-demographic inequalities
158 in the dynamics of uncertainty reduction.

159 The rest of the study is structured as follows: section 2 describes the datasets and
160 methods used to study exposure disparities, pollution-mortality associations, and uncertainty

161 tradeoffs from changes in input information. Section 3 presents the study results, and section 4
162 concludes with a discussion of the results' implications and dialogue with recent literature.

163 **2. Methods**

164 **2.1 Data**

165 Mortality data

166 We use individual-level mortality data for North Carolina from 2001 to 2016. The data
167 was obtained from the North Carolina State Center for Health Statistics, Vital statistics
168 department. Our analysis utilizes each participant's date of death, residential location, and
169 race/ethnicity. We studied total mortality (all causes of death except external causes,
170 International Classification of Diseases, ICD10, A00-R99). Other individual characteristics not
171 analyzed in this work are also included in the mortality dataset, such as sex, age at death,
172 education, and marital status. Additional analysis of the correlation of air pollution mortality
173 with these individual-level variables, as well as that of residential and environmental variables,
174 has been performed elsewhere (Son et al., 2020).

175 Air pollution data

176 We use daily gridded data from a 1km model of PM_{2.5} concentration (Di et al., 2021).
177 This ensemble-based model utilizes machine learning algorithms and multiple variables from
178 monitoring stations from the Environmental Protection Agency (EPA), satellite measurements,
179 land use terms, chemical transport model output, and others, to predict daily PM_{2.5} for the entire
180 United States. More details about model development and evaluation are available elsewhere (Di
181 et al., 2019). The exposure assigned to each participant is based on the 1km gridcell that contains
182 their residential location.

183 Weather data

184 We include daily gridded data on mean temperature and dewpoint temperature as
185 covariates in our mortality modeling. Inclusion of these covariates is common practice in air
186 pollution-epidemiology studies (e.g., (Nhung et al., 2017; Son et al., 2020)) to control for
187 weather-related mortality. These data are obtained on a 4×4km grid from the Parameter-elevation
188 Regressions on Independent Slopes Model (PRISM), which combines ground-based
189 measurement station data with a digital elevation model to create gridded climate products for
190 the U.S. Additional details are available elsewhere (Daly et al., 2008; PRISM Climate Group,
191 2004). Similarly to the air pollution data, each participant is assigned the weather data of the grid
192 cell containing their residence.

193 Census data

194 We utilize US census data on race for the analysis of disparities in air pollution exposure.
195 We chose the data for 2010 since this census year falls around the middle of the range of our
196 analysis (2001-2016). A comparison with 2020 census data determined that although North
197 Carolina’s population is increasing, the changes in racial composition and spatial distribution of
198 the population are small enough for the results of our study to not be affected by the choice of
199 census year.

200 **2.2 Exposure disparities**

201 The 2010 US census reports 21.2% of the population of North Carolina was NHB,
202 making them the largest racial minority in the state. Therefore, we focus our study of PM_{2.5}
203 exposure disparities on the NHB population.

204 We derive the average PM_{2.5} concentration between 2001 and 2016 for each census tract
205 in the state and compare these to the tract’s %NHB using quantile regression (Koenker & Bassett

206 Jr, 1978; Koenker & Hallock, 2001). Quantile regression estimates the conditional quantile(s) of
207 interest of the response variable (in this case, $PM_{2.5}$) as a linear combination of the predictor
208 variable (in this case, %NHB). We model the 10th, 25th, 50th, 75th, and 90th percentile $PM_{2.5}$ using
209 data from the 1405 census tracts in the state with NHB residents. Ordinary linear regression, in
210 contrast, estimates the conditional mean of the response variable, only giving information about
211 the relationship between air pollution levels and the percentage of NHB residents for the
212 “average” census tract. Using quantile regression provides more comprehensive results, allowing
213 us to study this relationship for the more and least polluted census tracts, as well as the median
214 census tracts, thus exploring racial inequalities in exposure at different relative exposure levels.

215 In addition to state-wide results, we also investigate exposure disparities for the two most
216 populated counties in the state: Mecklenburg County (population 923,427 in the 2010 census,
217 50.5% Non-Hispanic White (NHW) and 30.2% NHB) and Wake County (population 906,969 in
218 the 2010 census, 62.2% NHW and 20.4% NHB). We report quantile regression results for each
219 county, and we also compare the density function of the %NHB population in the least polluted
220 census tracts in each county, determined as those with average $PM_{2.5}$ in the 1st quartile, to density
221 function of %NHB in the most polluted census tracts (those with average $PM_{2.5}$ in the 4th
222 quartile). This comparison of density functions provides an assessment of the differences in the
223 racial distribution of the population between the most polluted and least polluted census tracts in
224 the county.

225 **2.3 Mortality modeling**

226 We model the association between $PM_{2.5}$ and short-term mortality with a case-crossover
227 design. This model uses each individual as their own control, eliminating the need to control for
228 individual-level characteristics and thus greatly reducing the number of necessary covariates for

229 good model specification. This low number of covariates presents an advantage for our goal of
230 isolating the influence of increasing input data for a specific variable (in this study, either for
231 PM_{2.5} or mortality) on the uncertainty reduction of the epidemiology model. For a different type
232 of model requiring more individual-level controls, the epistemic uncertainty introduced by a high
233 number of covariates could obscure the uncertainty reduction achieved by any single variable's
234 information gain. We select control days based on the same day of the week of the same month
235 of the individual's death. Each case day therefore has more than one control, and we allow for
236 bi-directional sampling of controls (selection of control days both before and after the
237 individual's death) to control for bias from temporal trends in the pollution data (Navidi, 1998).
238 Temperature and dewpoint temperature are also incorporated as covariates in the model.

239 The choice to investigate the pollution-mortality association in the short-term is
240 motivated by the type of health data available for this study. We use a dataset where cases have
241 been selected based on health outcome (in this case, mortality), making the data suitable for a
242 short-term study using a case-control design and further, for a case-crossover design since we do
243 not have data on other individuals who did not experience the outcome of interest (Belbasis &
244 Bellou, 2018; Jaakkola, 2003). Since air pollution has been widely recognized to have both
245 short-term and long-term effects, the same information tradeoffs methodology presented here
246 could be applied to a different epidemiology model in the presence of health data suitable for a
247 long-term study. For example, a long-term study could be performed using a cohort design,
248 where participants are selected based on their degree of exposure to air pollution and placed into
249 the "exposed" or "unexposed" group, and then health outcomes for these groups are observed
250 and compared over a specified period of time (Belbasis & Bellou, 2018).

251 The coefficients of the case-crossover model are fit using conditional logistic regression
 252 (Pampel, 2020). If we describe mortality Y_i as following a Bernoulli distribution (equation (1a)),
 253 where Y_i can be equal to 1 for the day of death or 0 for the control day(s), and the probability that
 254 $Y_i = 1$ is P , then we can model the logged-odds of P as a linear relationship between our
 255 predictors of interest (equation(1b)):

$$Y_i \sim \text{Bernoulli}(P); \quad (1a)$$

$$\ln\left(\frac{P}{1-P}\right) = \alpha + \beta PM_{2.5} + \gamma T_t + \delta D_t, \quad (1b)$$

256 where α is the intercept and β is the fitted coefficient describing the association of $PM_{2.5}$ with
 257 mortality, also called exposure coefficient. We will focus on β for the study of uncertainty
 258 reduction in the case-crossover model (additional details are provided in section 2.4). The
 259 coefficients γ and δ describe the association of temperature (T) and dewpoint temperature (D),
 260 respectively. Solving for the odds by exponentiating equation (1b) gives us the expression:

$$\frac{P}{1-P} = e^\alpha \times e^{\beta PM_{2.5}} \times e^{\gamma T} \times e^{\delta D}, \quad (1)$$

261 where each exponent term can be interpreted as the odds ratio (OR) for the association of each
 262 covariate with mortality. Our main interest lies in the second exponent on the right-hand side,
 263 $e^{\beta PM_{2.5}}$. This term represents the OR for a $PM_{2.5}$ increment of $1 \mu\text{g}/\text{m}^3$, which we will refer to as
 264 OR_1 . For consistency with common practice in reporting of epidemiology results, we will report
 265 the OR for a $PM_{2.5}$ increment of $10 \mu\text{g}/\text{m}^3$ (OR_{10}) which can be derived from OR_1 as:

$$OR_{10} = e^{\beta \times 10} = (e^\beta)^{10} = (OR_1)^{10}. \quad (2)$$

266 We initially examine the association of mortality with $PM_{2.5}$ at multiple lags: lag0, lag1,
 267 and lag 2 (meaning the $PM_{2.5}$ on the day of death, 1 day before death, and 2 days before death,
 268 respectively). We also analyze two cumulative lags: lag01 (the cumulative effect of lags 0 and 1)
 269 and lag02 (cumulative effect of lags 0, 1, and 2), by fitting mortality against the average of the

270 PM_{2.5} levels at the lags of interest. Then we perform stratified analysis to investigate differences
271 in effects between the NHB and NHW populations at the aforementioned PM_{2.5} lags. Since this
272 stratified analysis performs multiple tests on subsets of the same dataset, we adjust its results for
273 multiplicity by using the Bonferroni correction (Chen et al., 2017; Hochberg & Tamhane, 1987).
274 Based on the results of the full model and the stratified analysis, we will select a single lag of
275 PM_{2.5} for further investigation of uncertainty tradeoffs. The temperature and dewpoint
276 temperature covariates have the same lag as the PM_{2.5} in each model fit.

277 **2.4 Information change and uncertainty tradeoffs**

278 This study adopts the uncertainty tradeoffs methodology developed in (Alifa et al., 2022)
279 for the study of a realistic case scenario through the use of spatio-temporal data on pollution,
280 mortality, and demographics. We will study how fitting the case-crossover model described in
281 2.3 with changing input information on mortality and air pollution (Y_i and PM_{2.5} in equation (1
282)), respectively) affects the uncertainty of the pollution-mortality coefficient, β , in the model fit.
283 We will also take advantage of the demographic information included in the mortality dataset to
284 investigate racial differences in uncertainty reduction from improved health data.

285 Uncertainty quantification of the mortality model

286 We use the metric of information entropy to characterize the uncertainty of our estimate
287 for the exposure coefficient, $\hat{\beta}$. Since we can assume $\hat{\beta}$ is a continuous random variable, its
288 entropy can be defined as (Christakos, 2012):

$$H(\hat{\beta}) = - \int_{-\infty}^{\infty} f(\hat{\beta}) \ln(f(\hat{\beta})) d\hat{\beta}, \quad (3)$$

289 where $f(\hat{\beta})$ is the probability density function (PDF) of the estimate. As more input information
290 is acquired for the model in equation (1)), the inference becomes more accurate such that $\hat{\beta} \rightarrow \beta$

291 in probability, which results in a reduction of $H(\hat{\beta})$. Our previous publication (Alifa et al., 2022)
292 demonstrated several methods for deriving entropy both parametrically and non-parametrically.
293 For this study, we derive $H(\hat{\beta})$ parametrically from the standard error of the exposure coefficient,
294 $\hat{\sigma}_{\beta}^2$, output from the conditional logistic regression fit. Assuming $\hat{\beta}$ to be asymptotically normal,
295 we use the closed form equation for the entropy of a normal distribution,

$$H(\hat{\beta}) = \frac{1}{2} \log(2\pi e \hat{\sigma}_{\beta}^2). \quad (4)$$

296 Additionally, the relative entropy $\Delta H_{\hat{\beta}}$ is a useful metric to compare the uncertainty of
297 different information stages. We can define the vector $\Delta H_{\hat{\beta}}$ as:

$$\Delta H_{\hat{\beta}} = \mathbf{H}_{\hat{\beta}} - H_{\hat{\beta},\text{ref}}, \quad (5)$$

298 where $\mathbf{H}_{\hat{\beta}}$ is a vector containing $H(\hat{\beta})$ for different stages of information, and $H_{\hat{\beta},\text{ref}}$ is the
299 entropy for the information stage selected as reference. For this study we order the elements of
300 $\mathbf{H}_{\hat{\beta}}$ from those computed with least to most information, and select the stage with most
301 information as our reference, resulting in a $\Delta H_{\hat{\beta}}$ that decreases towards 0.

302 Change in air pollution information

303 We generate different stages of air pollution information by upscaling the original 1km
304 PM_{2.5} model to two coarser resolutions, 6km and 12km. We then fit the model in equation (1)
305 with the three different resolutions and compare $H(\hat{\beta})$ for the three cases. These different stages
306 of information simulate a situation where stakeholders are currently operating with coarse-
307 resolution output such as that from the EPA's Community Multiscale Air Quality Model
308 (CMAQ, 12km resolution) or other similar gridded products, and want to explore the information
309 benefits of downscaling their data to higher resolutions.

310 Change in mortality information

311 To change the amount of input mortality information, we fit equation (1) with varying
312 number of mortality records. This simulates a case where stakeholders are interested in
313 investigating the benefit of augmenting the health outcomes dataset used for their assessment,
314 due to known or suspected missing cases in said dataset. We will investigate the effect of racial
315 bias in the missing data by comparing the uncertainty reduction when cases are missing only
316 from the NHW population versus cases missing only from the NHB population. We choose these
317 two subpopulations for comparison since in the 2010 US census the racial majority in North
318 Carolina was NHW with 65.2% of the population, while the largest racial minority was NHB,
319 conforming 21.2% of the population. Since NHB cases represented about 20% of the study
320 population, this is the maximum number of missing cases we explore for both races. Therefore,
321 we initially fit the model with ~80% of the total mortality data, where the ~20% of missing cases
322 are either all NHW or NHB patients. Then we increase the number of patients and repeat the fit
323 again with ~90% of data, and lastly with 100% data coverage. We select missing cases at random
324 from the pool of participants of the race of interest, and repeat each model fit 100 times to obtain
325 ensemble results from which we compute the mean and 95% CI of $H(\hat{\beta})$ at each information
326 stage.

327 Information yield curves

328 Information yield curves (Alifa et al., 2022; De Barros & Rubin, 2008; De Barros et al.,
329 2009) are a graphical device designed to display the tradeoffs in uncertainty reduction between
330 information gain in air pollution and health data. This tool plots together, in mirror image, the
331 separate effects of information increase for each of these datasets on the uncertainty reduction of
332 $\hat{\beta}$, enabling decision-makers to visualize the most efficient pathway to improve their assessment
333 in their particular case scenario. In our previous study (Alifa et al., 2022) the changes in input

334 data were first associated with changes in uncertainty for separate pollution and health models
335 which when brought together would propagate to the final mortality uncertainty. Therefore, the
336 information yield curve compared the changes in entropy for the separate pollution and health
337 models (in the x axis) to the final change in entropy of the pollution-mortality assessment (in the
338 y axis). The nature of the datasets in this current study requires a modification of the previous
339 method by associating the changes in information for the input datasets directly with the changes
340 in the final uncertainty of the case-crossover model fit. This results in an x-axis of qualitative
341 nature, since there is no common unit to compare increased number of mortality records to
342 increased resolution of the PM_{2.5} grid. However, decision-makers taking advantage of this
343 method in the future would be able to find a common metric for information increase from each
344 dataset given their particular case scenario, such as cost of added data or time for data
345 computation/procurement.

346 **3. Results**

347 **3.1 Descriptive statistics**

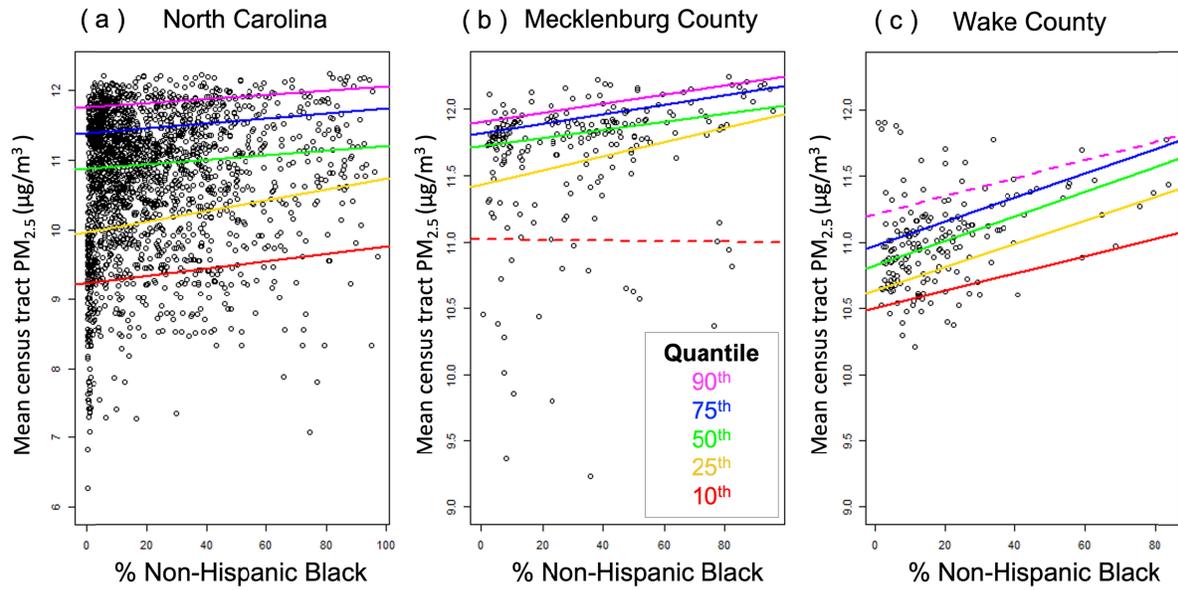
348 The mortality model had input of a total of 1,065,699 cases with 3,621,521 controls (3.40
349 controls per case). These cases contained more females than males (52.1% vs 47.9%), and the
350 majority of deaths were from people older than 65 years old (75.4%). Most cases were Non-
351 Hispanic White (77.4%), while the second most cases were Non-Hispanic Black (20.4%). Table
352 S1 shows the full demographics of the mortality data used in the model.

353 The median of the PM_{2.5} in the model was 9.5 µg/m³, with lower bound (5th percentile) of
354 3.8 µg/m³ and upper bound (95th percentile) of 21.5 µg/m³. These quantiles varied by less than
355 0.1 µg/m³ when recomputed separately for case days and control days. The median temperature

356 was 15.7°C, with 5th and 95th percentiles of 0.7°C and 27.4°C, respectively. The median
357 dewpoint temperature was 10.5°C and its 5th and 95th percentiles were -8.4°C and 21.9°C,
358 respectively.

359 **3.2 Exposure disparities**

360 The quantile regression for the whole state shows a significant, positive correlation
361 between average PM_{2.5} and percent NHB population across all the quantiles modeled (Figure 1,
362 panel a). This indicates that more polluted census tracts tend to have a higher percentage of NHB
363 population across the entire state, regardless of the relative exposure level. Localized results
364 from Mecklenburg and Wake counties (Figure 1, panels b and c) show the same significant,
365 positive association for most quantiles studied. Figure 2 also shows that in both these counties,
366 the majority of the least-polluted census tracts (those ranked in quartile 1 using average PM_{2.5} as
367 criteria) have a low percentage of NHB population, while the most polluted tracts (ranked in
368 quartile 4) tend to have comparatively higher percentages of NHB residents.



369

370 *Figure 1. Quantile regression between census tract average $PM_{2.5}$ (years 2001-2016) and*

371 *census tract percent of Non-Hispanic Black population for (a) all census tracts in North*

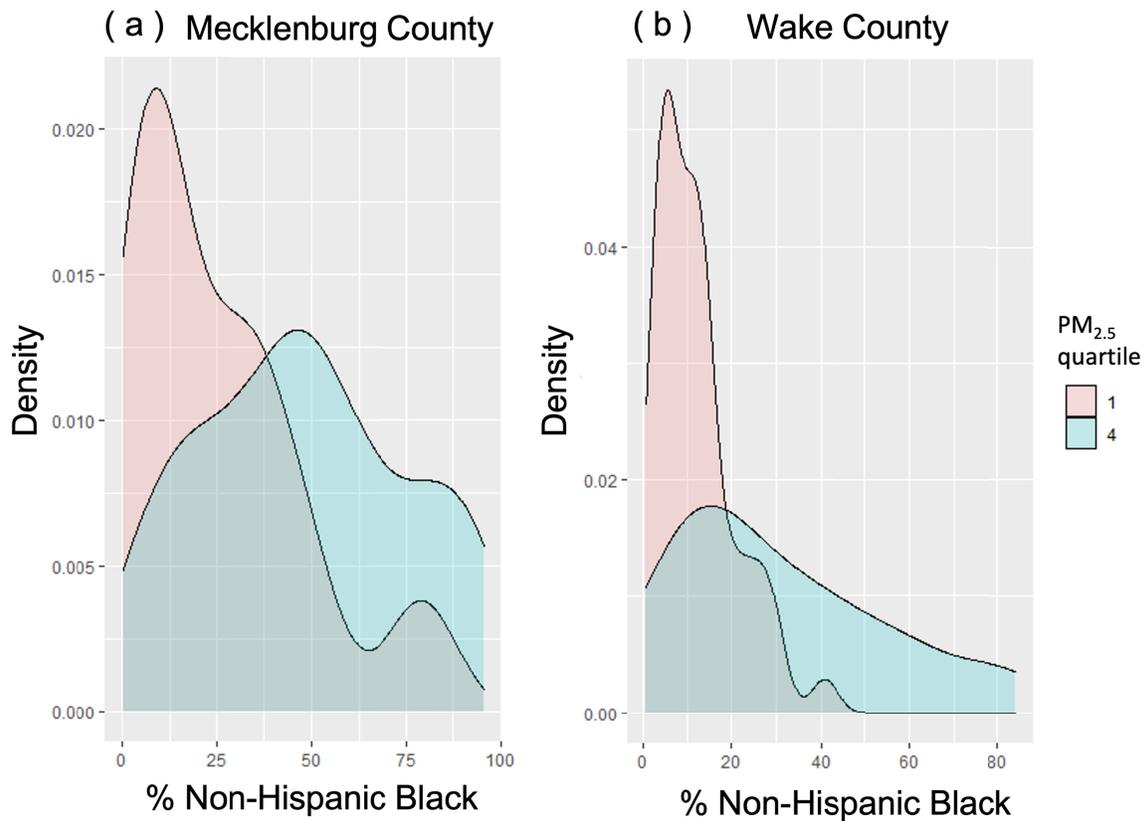
372 *Carolina, (b) census tracts in Mecklenburg County, and (c) census tracts in Wake County. The*

373 *inset in panel (b) provides a color reference for the quantiles plotted. Non-statistically*

374 *significant results are represented with dashed lines. Note the y-axis scale in panel (a) is*

375 *different from that in panels (b) and (c).*

376



377

378 *Figure 2. Density of percent Non-Hispanic Black for census tracts with average PM_{2.5} in*
 379 *the first quartile (red) and in the fourth quartile (blue), for (a) Mecklenburg County and (b)*
 380 *Wake County.*

381 **3.3 Mortality model**

382 We first present the results of the case-crossover model computed with the full record of
 383 mortality and using data from the highest resolution PM_{2.5} gridded data (1km). We will later
 384 compare the changes in uncertainty for that model when fit with less data, by either reducing the
 385 number of mortality cases in the model or by using data from coarser PM_{2.5} grids. All the model
 386 fits are performed with the same (4x4km) datasets for temperature and dewpoint temperature
 387 taken at the same temporal lags as the PM_{2.5} data.

388 Table 1 reports the odds ratios for a 10 $\mu\text{g}/\text{m}^3$ increase in $\text{PM}_{2.5}$ (OR_{10}) and its 95%
 389 confidence intervals for the five different lags investigated. The significant associations observed
 390 were, in descending magnitude: for lag01, $\text{OR}_{10} = 1.016$ (95% CI 1.011–1.021); lag02, $\text{OR}_{10} =$
 391 1.016 (95% CI 1.010–1.022); lag0, $\text{OR}_{10} = 1.013$ (95% CI 1.009–1.018), and lag1, $\text{OR}_{10} = 1.012$
 392 (95% CI 1.007–1.017). The association for lag2 was not statistically significant.

393 Our results were very similar to those of a previous study that used the same model
 394 design and mortality data (Son et al., 2020), with minor (and statistically non-significant)
 395 differences attributable to differences in sources and averaging techniques for the pollution and
 396 temperature data (comparison can be found in Figure S1).

397 *Table 1. Odds Ratios and 95% confidence intervals for the association of $\text{PM}_{2.5}$ with*
 398 *mortality at different lags. Non-significant results are colored in grey.*

Lag	OR_{10}
Lag0	1.013 (1.009 - 1.018)
Lag1	1.012 (1.007 - 1.017)
Lag2	1.004 (0.999 - 1.008)
Lag01	1.016 (1.011 - 1.021)
Lag02	1.016 (1.010 - 1.022)

399
 400 We also fit the case crossover models separately for the NHW and NHB cases to
 401 investigate effect differences between these population groups. Table 2 shows the OR_{10} and the
 402 (multiplicity adjusted) 95% confidence interval for each lag and race. The association between
 403 $\text{PM}_{2.5}$ and short-term mortality was significant in the NHW population for all lags except Lag2,
 404 the same lags where the association was also significant when the whole study population was
 405 represented (Table 1). This is a sensible result since the majority of the mortality cases studied
 406 come from the NHW population (77.4%). The results for the NHB population present wider
 407 confidence intervals, associated to the relatively lower number of cases that were used to fit the

408 model since only 20.4% of the study population is NHB, making the multiplicity-adjusted results
 409 for NHB not statistically significant. We will use the Lag1 model for subsequent analysis since it
 410 was the lag with the closest to significant association for NHB.

411 *Table 2. Odds Ratios and 95% confidence intervals for the association of PM_{2.5} with*
 412 *mortality at different lags. Non-significant results are colored in grey.*

Lag	OR₁₀ NHW	OR₁₀ NHB
Lag0	1.015 (1.010 - 1.020)	1.006 (0.992 - 1.021)
Lag1	1.013 (1.007 - 1.018)	1.010 (0.999 - 1.022)
Lag2	1.005 (0.998 - 1.011)	no effect
Lag01	1.018 (1.012 - 1.024)	1.011 (0.998 - 1.025)
Lag02	1.018 (1.011 - 1.025)	1.010 (0.994 - 1.026)

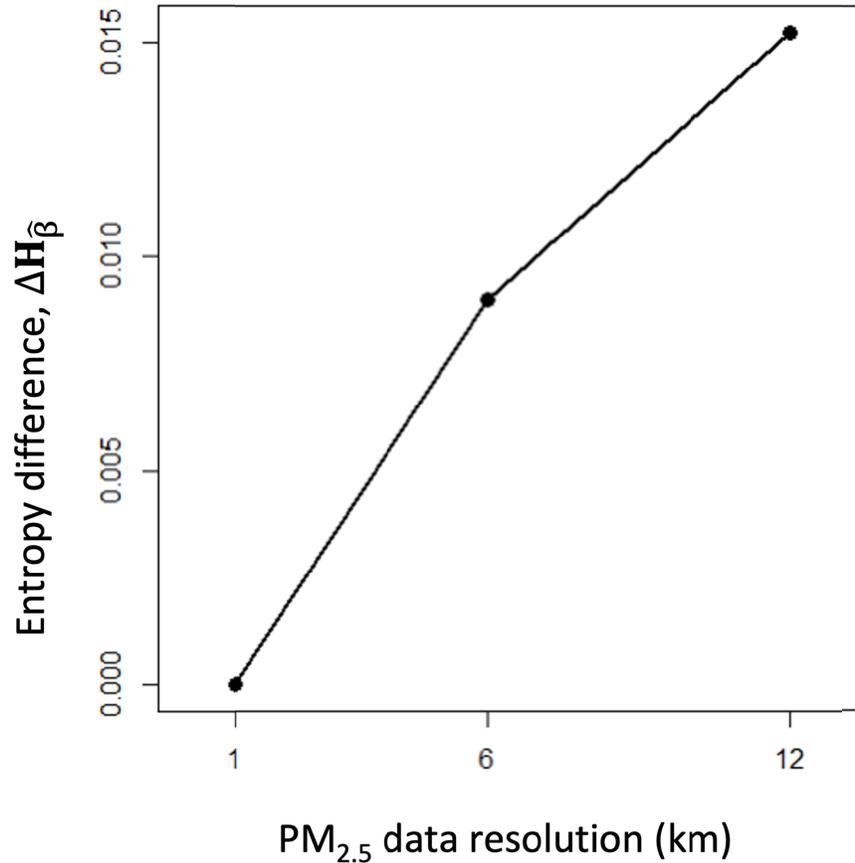
413

414 **3.4 Uncertainty tradeoffs from information changes**

415 To study uncertainty tradeoffs, we fit the model in equation (1) with varying input of
 416 either PM_{2.5} data or mortality data (Y_i), in order to compare each of these datasets' influence in
 417 the final uncertainty of the case-crossover model, measured through the entropy of the exposure
 418 coefficient β , as explained in section 2.4.

419 First, we isolate the influence of changing air pollution data on the case-crossover
 420 model's uncertainty reduction. To achieve this, we fit the model with the full record of mortality
 421 data while varying PM_{2.5} data, by fitting the model three times with PM_{2.5} data of different
 422 resolutions (1km, 6km, and 12km). Figure 3 shows that fitting the model with finer resolution
 423 PM_{2.5} data results in lower uncertainty of β . Since the PM_{2.5} exposure is assigned based on each
 424 individual's gridcell of residence, a coarser grid may result in more deaths that happened the
 425 same day falling within the same gridcell, causing multiple cases to have identical PM_{2.5} data.
 426 Although weather covariate data may still be different for each case (since these are always on

427 the same 4km grid) making the cases sharing $PM_{2.5}$ data still likely distinct, the repeated
428 sampling of the same $PM_{2.5}$ values does not provide new information to the model, therefore
429 reducing the information value of the air pollution data input.



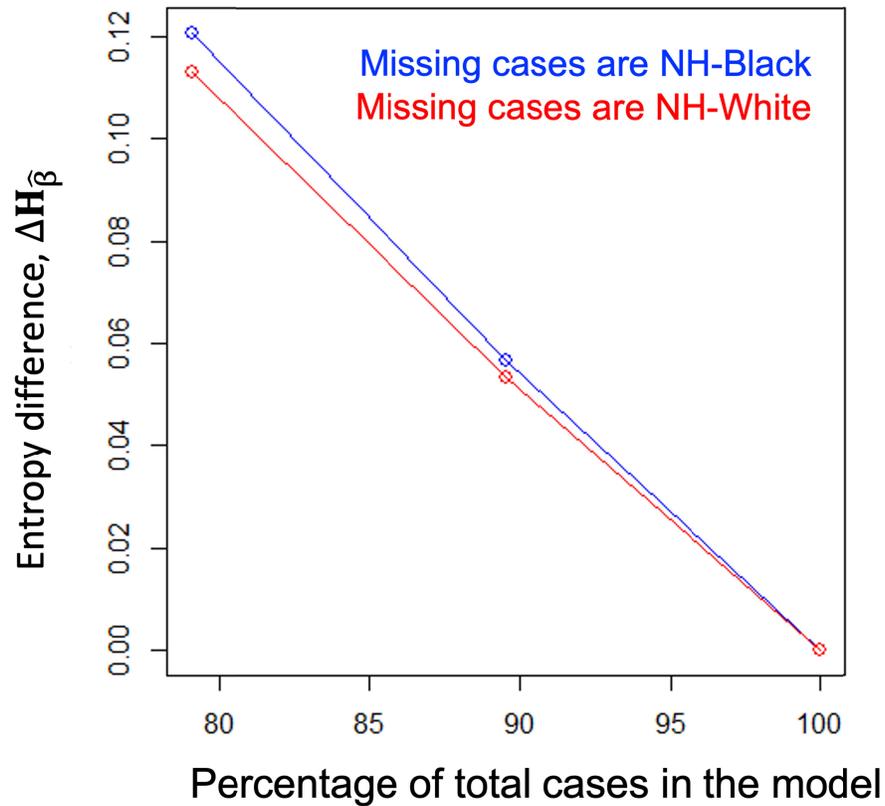
430

431 *Figure 3. Entropy changes for the estimate of the exposure coefficient $\hat{\beta}$ for case-*
432 *crossover model fit with $PM_{2.5}$ data of different spatial resolutions.*

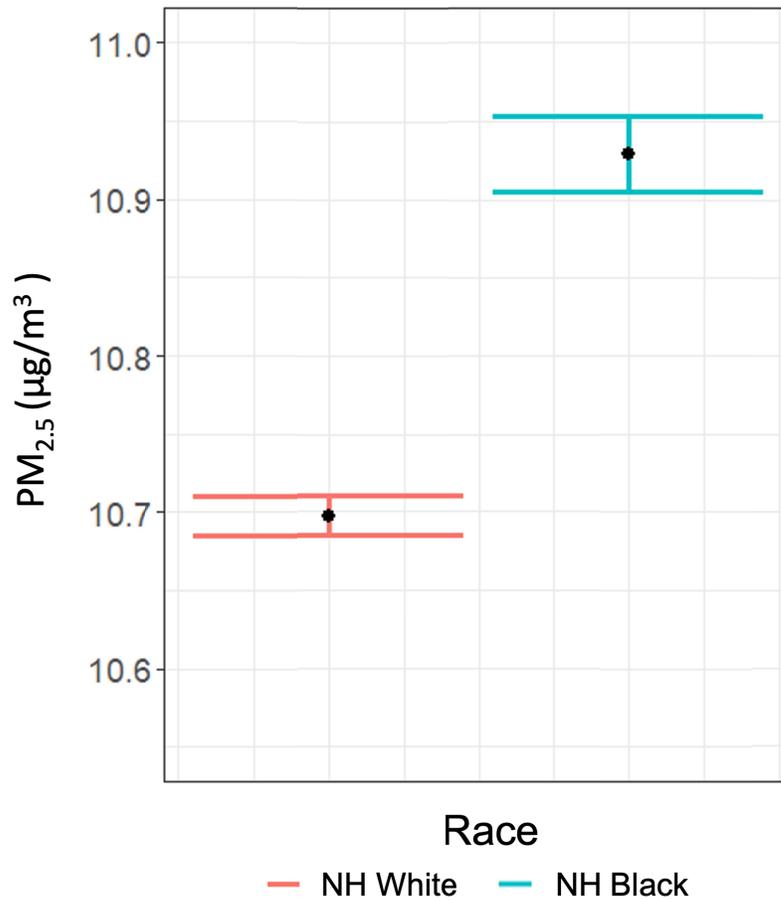
433 Then, we isolate the effect of changing mortality data in the uncertainty of the case-
434 crossover model. To achieve this, we fit the model with the highest-resolution $PM_{2.5}$ data (1km)
435 while varying the number of mortality cases input into the model. We do this analysis twice,
436 selecting the missing cases to be either all from the NHW population or the NHB population, in

437 order to investigate the effect of racial bias in the uncertainty reduction dynamics of health data.
438 Since NHB cases represented approximately 20% of the study population, this is the maximum
439 number of missing cases we explore for both races. Therefore, we initially fit the model with
440 ~80% of data, and we then increase the number of cases to ~90% and finally to 100% data
441 coverage. Figure 4 shows that while increasing the number of mortality cases reduces uncertainty
442 in the model for both scenarios, the slope of uncertainty reduction is steeper when the new cases
443 introduced are from the NHB population. The exposure disparities experienced by the NHB
444 population shown in section 2.2 may be related to this difference, since differential exposure of a
445 subpopulation may lead to a higher diversity of pollution data input in the model. This
446 hypothesis is confirmed by the differences in the distribution of the mean of the Lag1 $PM_{2.5}$ data
447 associated with cases and controls from the NHB population versus that one associated to the
448 NHW population (Figure 5). The 95% confidence intervals between both distributions do not
449 cross, making the mean $PM_{2.5}$ associated with NHB individuals statistically different from that of
450 NHW individuals. At the lowest stage of information the model is fit with ~80% of the data, the
451 majority of which comes from NHW individuals, so adding more data from NHW individuals
452 will introduce samples from the $PM_{2.5}$ distribution that is already known the most. In contrast,
453 new data from NHB individuals introduces information from a distribution of $PM_{2.5}$ that is
454 different from the majority distribution, providing new information to the model and generating a
455 faster uncertainty reduction. This result is not caused by the higher magnitude of the mean $PM_{2.5}$
456 for NHB shown in Figure 5, but by the fact that the NHB are a minority population with a
457 statistically different $PM_{2.5}$ exposure distribution from that of the NHW population. Therefore,
458 uncertainty reduction should have been steeper with new NHB data even if this subpopulation

459 was exposed to less pollution than the NHW population, as long as the mean $PM_{2.5}$ between
460 subpopulations remained statistically different.



461
462 *Figure 4. Entropy changes for the estimate of the exposure coefficient $\hat{\beta}$ for case-*
463 *crossover model fit when more information is acquired for NHB cases only (blue series) or NHW*
464 *cases only (red series).*



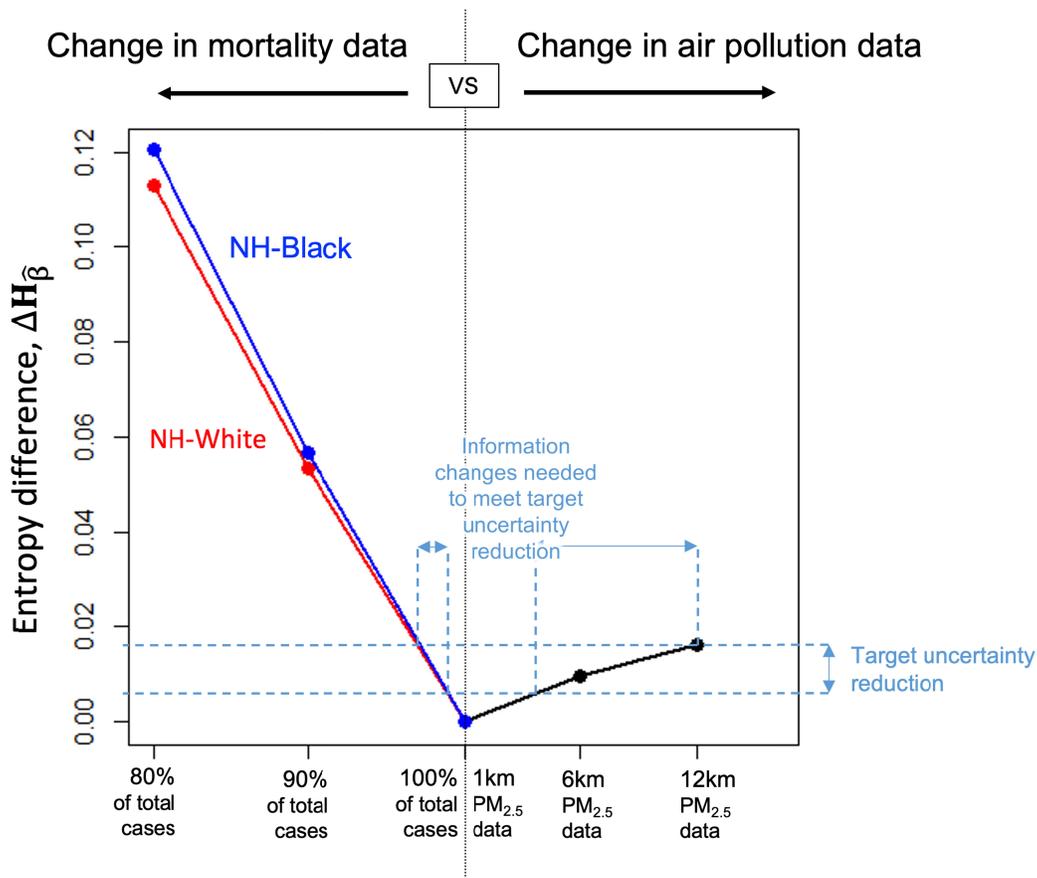
465

466 *Figure 5. Mean of the lag-1 PM_{2.5} associated with NHW cases (red) and NHB cases*
 467 *(blue) in the case-crossover model (equation (1)) computed with state-wide data, and its 95%*
 468 *confidence interval.*

469 **3.5 Information yield curve**

470 While we showed in section 3.4 that increasing air pollution information and health
 471 effects information both reduce the uncertainty in the final mortality estimate, their contribution
 472 to uncertainty reduction is not equal. The information yield curve in Figure 6 compares the
 473 individual effects of information gain from each dataset in the model's uncertainty reduction.
 474 The dashed light-blue lines illustrate a graphical interpretation that can be used for decision-

475 making purposes. If for a case scenario of interest, the target for mortality uncertainty reduction
476 is $\Delta H_{\hat{\beta}}$ as indicated by the horizontal dashed lines, the change in the x axis required for the data
477 in each side can be compared to find the most efficient pathway for uncertainty reduction. In the
478 case below, increasing health data seems to reduce the uncertainty in the model more efficiently,
479 since the same $\Delta H_{\hat{\beta}}$ can be achieved with a smaller change in x. However, the figure below
480 presents a qualitative x-axis, as there is no common basis of comparison between increasing
481 patient data and downscaling pollution model resolution. For a real-world scenario, stakeholders
482 would be able to apply a common metric to these data improvements, such as cost or time,
483 making the x-axis quantitative and potentially altering the decision-making outcomes presented
484 here.
485



486

487 *Figure 6. Information yield curve comparing the effect of information gain in mortality*
 488 *(left side) versus air pollution (right side) on the uncertainty reduction of the exposure coefficient*
 489 *in the case crossover model. The dashed light-blue lines provide graphical interpretation of the*
 490 *information yield curve by illustrating the different data increases necessary to achieve a fixed*
 491 *risk uncertainty reduction.*

492

493 **4. Discussion and conclusion**

494 The results of this study illustrate the usefulness of our information entropy tradeoff
495 methodology to not only generate more robust impact assessments, but also to gain new
496 knowledge about the role of data from minority populations in the dynamics of uncertainty
497 reduction.

498 We found associations between short-term PM_{2.5} exposure and mortality for years 2001-
499 2016 in North Carolina that were statistically significant and consistent with a previous study of
500 the same mortality dataset (Son et al., 2020), despite the state’s relatively low and decreasing air
501 pollution levels. North Carolina had a state-wide average PM_{2.5} concentration of 13.5 µg/m³ in
502 2002, and state-wide decreases in concentrations resulted in the whole state presenting annual
503 mean PM_{2.5} below the EPA’s standard of 12 µg/m³ by 2016 (Bravo et al., 2022). Despite this
504 improving trend in pollution concentrations, our findings add to the mounting evidence that
505 particulate matter has detectable health effects even at pollution levels formerly seen as safe,
506 motivating ongoing updates of air quality guidelines such as the EPA’s proposal in January of
507 2023 to reduce the PM_{2.5} standard to between 9 and 10 µg/m³.

508 We also explored tradeoffs between data increases in air pollution or health outcomes in
509 the uncertainty reduction of the case-crossover model used to investigate the pollution-mortality
510 relationship. The information yield curve presented in Figure 6 compared the different
511 uncertainty reduction effects of augmenting information in air pollution and health data. While
512 both data types reduce uncertainty in the case-crossover model when information is increased,
513 the effect of new data for mortality resulted in a steeper rate of uncertainty reduction. One
514 qualification of this outcome is that information increase was done by different methods for each
515 dataset, making the comparison of information change merely qualitative as there is no common

516 variable in the x-axis of the information yield curve. If this method were applied to a scenario
517 where information increases are associated to costs, time, or, as done in our previous study (Alifa
518 et al., 2022), pollution/health model uncertainties, the comparison could be done qualitatively
519 and the decision-making outcomes of the information yield curve may change. The goal of this
520 work is not to provide an absolute answer to the choice between investing in pollution versus
521 health information, but to develop a framework applicable to any data set and environmental
522 exposure scenario used in any epidemiological model.

523 The positive relationship between average $PM_{2.5}$ and %NHB population found at the
524 census tract level through quantile regression is consistent with previous findings of disparities in
525 exposure for the NHB population in both nationwide (Miranda et al., 2011; Tessum et al., 2021;
526 Woo et al., 2019); and regional (Bravo et al., 2016; Servadio et al., 2019; Stuart et al., 2009)
527 studies. Our study of Mecklenburg and Wake counties further illustrated the presence of this
528 inequality for the most populated areas of the state, which experience relatively higher levels of
529 air pollution. However, the state-wide positive association found with respect to all the
530 concentration quantiles also reveals that exposure inequalities can be detected not only among
531 counties such as Mecklenburg and Wake with high emissions (placed in the high $PM_{2.5}$
532 quartiles), but also among counties with lower emissions (those in the low $PM_{2.5}$ quartiles),
533 indicating that these racial inequalities may be independent from the relative difference in
534 pollution levels between counties that have different emission types or levels of urbanicity,
535 agreeing with recent nationwide findings (Liu et al., 2021; Tessum et al., 2021). These findings
536 of exposure disparities are not reflected in the results of the stratified case crossover model,
537 possibly due to the relatively low $PM_{2.5}$ levels in the state that result in relatively small
538 magnitude of exposure disparities.

539 A key finding of this paper is that disparities in $PM_{2.5}$ exposure can affect model
540 uncertainty reduction. If exposure from a certain minority subpopulation (in this case, the NHB
541 population) is significantly different than that of the majority population, as shown in Figure 5,
542 then data from this minority have relatively higher information value resulting in a faster rate of
543 uncertainty reduction in the mortality model (Figure 4). The authors hypothesize that this result
544 is transferrable to the study of any minority subpopulation (by race, income, residential location,
545 etc.) that experiences a different exposure from the majority, implying that minority
546 representation in environmental research benefits not only the minorities in question, but also the
547 researchers and stakeholders performing the research. In a situation where there is a known or
548 suspected environmental exposure difference between sub-populations, ensuring the
549 representation of all groups in the data used for the environmental impact assessment will result
550 in a wider sampling of the problem's information space, providing the quantitative advantage of
551 reduced uncertainty. Since minority groups have been found to be both over-exposed and at
552 times under-monitored (Stuart et al., 2009), the application of this framework will also provide
553 researchers with increased awareness of both exposure and information disparities by design,
554 contributing to the ongoing work of environmental justice.

555 There still remain multiple interesting opportunities for future expansion of the
556 uncertainty reduction framework proposed in our first study (Alifa et al., 2022) and further
557 expanded in this present work. One possible next step in future work is considering a case
558 scenario where the assessment goes from an initial baseline of comparatively scarce pollution,
559 epidemiology, or demographic information to subsequent stages of more information, via data
560 augmentation methods such as assimilation, disaggregation, and/or downscaling. This work
561 would require the integration of multiple datasets (e.g., by combing air pollution monitoring

562 station data, gridded CTM output, and area-based demographic and health outcomes data),
563 introducing new kinds of epistemic uncertainties, such as those stemming from errors in
564 pollution and exposure measurements, model specification, data aggregation, and extrapolation
565 of exposure-response functions, among others (Nethery & Dominici, 2019). These uncertainties
566 are different from the one addressed in our framework in that they increase monotonically with
567 the increase of input data, having the potential to obscure any uncertainty reduction from
568 information gain if the epistemic errors in the data are too high (Rao, 2005). For this reason, our
569 work so far has taken advantage of full datasets and simulated information scarcity by modeling
570 only subsets of this data, which has allowed us to explore the proposed framework without
571 having to deal with the epistemic uncertainties introduced by data assimilation errors.

572 The choice of North Carolina for this case study was prompted by the unique availability
573 of high-resolution mortality data, but the relatively low PM_{2.5} levels in the state prevented us
574 from incorporating true data assimilation into this project, since the noise introduced by multiple
575 PM_{2.5} data sources would have been greater than the signal of the PM_{2.5} data itself. This
576 limitation speaks to the wider issue of data scarcity in air pollution, health outcomes, and
577 demographics for the regions of the world that are most in need of epidemiology and exposure
578 disparities studies.

579 The framework developed here could still be useful, however, for a case of interest where
580 there is availability of pollution data only. As mentioned in the introduction, multiple methods to
581 augment air pollution observations through assimilation of other datasets such as CTMs, satellite
582 data, citizen-science observational networks have been devised in recent years. In a scenario
583 where stakeholders want to augment their observational network but are unsure of which method
584 to choose for the task, studying the information entropy tradeoffs between different data

585 assimilation methods may be an efficient way to inform a decision. Furthermore, if demographic
586 data is also available (such as census data), stakeholders would be able to investigate how
587 information increases from different air pollution sources have different effects in the uncertainty
588 of the estimates of exposure inequalities between different subpopulations, and whether focusing
589 on augmenting data in regions with high versus low concentrations of minority populations
590 yields different effects in uncertainty reduction.

591 As the scientific community continues efforts to improve characterization of
592 environmental exposure effects for overlooked areas and populations around the world, the
593 framework presented here gives researchers a new opportunity to elevate minority representation
594 from a qualitative afternote in a study's discussion section to a centerpiece of the study's design,
595 aiding a quantitatively more accurate analysis and producing confident estimates of the true
596 effects of environmental pollution.

597

598 **Acknowledgements**

599 This publication is based upon work supported by the Lucy Family Institute for Data &
600 Society at the University of Notre Dame, grant number 22006.

601

602 **Open Research**

603 The detailed death records data were obtained from the Children's Environmental Health
604 Initiative (CEHI) at Notre Dame (Children's Environmental Health Initiative, 2020). These data
605 are governed by data use agreements with data providers and protocols reviewed and approved
606 by the Institutional Review Board (IRB) at the University of Notre Dame. The data may be
607 accessed through a collaboration request to CEHI: <https://www.cehidatahub.org/collaborate>. The

608 1km gridded air pollution data was obtained from NASA’s SEDAC (Di et al., 2021) and can be
609 downloaded here: [https://sedac.ciesin.columbia.edu/data/set/aqdh-pm2-5-concentrations-](https://sedac.ciesin.columbia.edu/data/set/aqdh-pm2-5-concentrations-contiguous-us-1-km-2000-2016/data-download)
610 [contiguous-us-1-km-2000-2016/data-download](https://sedac.ciesin.columbia.edu/data/set/aqdh-pm2-5-concentrations-contiguous-us-1-km-2000-2016/data-download). The 4km gridded temperature and dewpoint
611 temperature was obtained from the PRISM Climate Group at Oregon State University (PRISM
612 Climate Group, 2004) and can be downloaded here: <https://prism.oregonstate.edu/downloads/>.
613 The 2010 census data can be downloaded from the Census Bureau, <https://data.census.gov/>. All
614 analyses were performed using R Statistical Software (v 4.2.3, R Core Team, 2023).

615 **5. References**

616

617 Alifa, M., Castruccio, S., Bolster, D., Bravo, M., & Crippa, P. (2022). Information entropy
618 tradeoffs for efficient uncertainty reduction in estimates of air pollution mortality.
619 *Environmental Research*, 212, 113587.

620 Atkinson, R., Kang, S., Anderson, H., Mills, I., & Walton, H. (2014). Epidemiological time
621 series studies of PM_{2.5} and daily mortality and hospital admissions: a systematic review
622 and meta-analysis. *Thorax*, 69(7), 660-665.

623 Bates, J. T., Pennington, A. F., Zhai, X., Friberg, M. D., Metcalf, F., Darrow, L., . . . Russell, A.
624 (2018). Application and evaluation of two model fusion approaches to obtain ambient air
625 pollutant concentrations at a fine spatial resolution (250m) in Atlanta. *Environmental*
626 *Modelling & Software*, 109, 182-190.

627 Beckx, C., Panis, L. I., Uljee, I., Arentze, T., Janssens, D., & Wets, G. (2009). Disaggregation of
628 nation-wide dynamic population exposure estimates in The Netherlands: Applications of
629 activity-based transport models. *Atmospheric Environment*, 43(34), 5454-5462.

630 Belbasis, L., & Bellou, V. (2018). Introduction to epidemiological studies. *Genetic*
631 *epidemiology: methods and protocols*, 1-6.

632 Bell, M. L., McDermott, A., Zeger, S. L., Samet, J. M., & Dominici, F. (2004). Ozone and short-
633 term mortality in 95 US urban communities, 1987-2000. *Jama*, 292(19), 2372-2378.

634 Bonas, M., & Castruccio, S. (2021). Calibration of Spatial Forecasts from Citizen Science Urban
635 Air Pollution Data with Sparse Recurrent Neural Networks. *arXiv preprint*
636 *arXiv:2105.02971*.

637 Bravo, M. A., Anthopolos, R., Bell, M. L., & Miranda, M. L. (2016). Racial isolation and
638 exposure to airborne particulate matter and ozone in understudied US populations:
639 Environmental justice applications of downscaled numerical model output. *Environment*
640 *international*, *92*, 247-255.

641 Bravo, M. A., Fuentes, M., Zhang, Y., Burr, M. J., & Bell, M. L. (2012). Comparison of
642 exposure estimation methods for air pollutants: ambient monitoring data and regional air
643 quality simulation. *Environmental research*, *116*, 1-10.

644 Bravo, M. A., Warren, J. L., Leong, M. C., Deziel, N. C., Kimbro, R. T., Bell, M. L., & Miranda,
645 M. L. (2022). Where is air quality improving, and who benefits? A study of PM_{2.5} and
646 ozone over 15 years. *American Journal of Epidemiology*, *191*(7), 1258-1269.

647 Breen, M., Chang, S. Y., Breen, M., Xu, Y., Isakov, V., Arunachalam, S., . . . Devlin, R. (2020).
648 Fine-scale modeling of individual exposures to ambient PM_{2.5}, EC, NO_x, and CO for
649 the coronary artery disease and environmental exposure (CADEE) study. *Atmosphere*,
650 *11*(1), 65.

651 Burnett, R., Pope III, C. A., Ezzati, M., Olives, C., Lim, S. S., Mehta, S., . . . Brauer, M. (2014).
652 An integrated risk function for estimating the global burden of disease attributable to
653 ambient fine particulate matter exposure. *Environmental health perspectives*, *122*(4),
654 397-403.

655 Chen, S.-Y., Feng, Z., & Yi, X. (2017). A general introduction to adjustment for multiple
656 comparisons. *Journal of thoracic disease*, *9*(6), 1725.

657 Children's Environmental Health Initiative. (2020). *North Carolina Detailed Death Records*
658 *during the period 2000 - 2017. [Data set].*
659 CEHI. https://doi.org/10.25614/ddrgeo_2000_2017.

660 Christakos, G. (2012). *Random field models in earth sciences*. Courier Corporation.

661 Coffman, E., Burnett, R. T., & Sacks, J. D. (2020). Quantitative Characterization of Uncertainty
662 in the Concentration–Response Relationship between Long-Term PM_{2.5} Exposure and
663 Mortality at Low Concentrations. *Environmental Science & Technology*, *54*(16), 10191-
664 10200.

665 Cohen, A. J., Brauer, M., Burnett, R., Anderson, H. R., Frostad, J., Estep, K., . . . Dandona, R.
666 (2017). Estimates and 25-year trends of the global burden of disease attributable to
667 ambient air pollution: an analysis of data from the Global Burden of Diseases Study
668 2015. *The Lancet*, *389*(10082), 1907-1918.

669 Colmer, J., Hardman, I., Shimshack, J., & Voorheis, J. (2020). Disparities in PM_{2.5} air pollution
670 in the United States. *Science*, *369*(6503), 575-578.

671 Crippa, P., Sullivan, R., Thota, A., & Pryor, S. (2019). Sensitivity of simulated aerosol properties
672 over eastern North America to WRF-Chem parameterizations. *Journal of Geophysical*
673 *Research: Atmospheres*, *124*(6), 3365-3383.

674 Daly, C., Halbleib, M., Smith, J. I., Gibson, W. P., Doggett, M. K., Taylor, G. H., . . . Pasteris, P.
675 P. (2008). Physiographically sensitive mapping of climatological temperature and
676 precipitation across the conterminous United States. *International Journal of*
677 *Climatology: a Journal of the Royal Meteorological Society*, *28*(15), 2031-2064.

678 De Barros, F., & Rubin, Y. (2008). A risk-driven approach for subsurface site characterization.
679 *Water resources research*, *44*(1).

680 De Barros, F., Rubin, Y., & Maxwell, R. M. (2009). The concept of comparative information
681 yield curves and its application to risk-based site characterization. *Water Resources*
682 *Research*, *45*(6).

683 Di, Q., Amini, H., Shi, L., Kloog, I., Silvern, R., Kelly, J., . . . Schwartz, J. (2019). An Ensemble-
684 based Model of PM_{2.5} Concentration Across the Contiguous United States with High
685 Spatiotemporal Resolution. *Environment International*, 130, 104909.

686 Di, Q., Wei, Y., Shtein, A., Hultquist, C., Xing, X., Amini, H., . . . Mickley, L. J. (2021). *Daily
687 and Annual PM_{2.5} Concentrations for the Contiguous United States, 1-km Grids, v1
688 (2000 - 2016)* NASA Socioeconomic Data and Applications Center (SEDAC).

689 Dominici, F., Peng, R. D., Bell, M. L., Pham, L., McDermott, A., Zeger, S. L., & Samet, J. M.
690 (2006). Fine particulate air pollution and hospital admission for cardiovascular and
691 respiratory diseases. *Jama*, 295(10), 1127-1134.

692 EPA, U. S. (1997). National ambient air quality standards for particulate matter: Final rule. *Fed.
693 Regist.*, 62(138), 38,651-638,701.

694 EPA, U. S. (2013). National Ambient Air Quality Standards for Particulate Matter. *Fed. Regist.*,
695 78(10), 3,086-083,287.

696 EPA, U. S. (2019). *Integrated Science Assessment (ISA) for Particulate Matter*

697 EPA, U. S. (2023). Reconsideration of the National Ambient Air Quality Standards for
698 Particulate Matter. *Fed. Regist.*, 88(18), 5,558-555,719.

699 EU. (2008). Directive 2008/50/EC of the European Parliament and of the Council of 21 May
700 2008 on ambient air quality and cleaner air for Europe. *Official Journal of the European
701 Union*.

702 Fuller, R., Landrigan, P. J., Balakrishnan, K., Bathan, G., Bose-O'Reilly, S., Brauer, M., . . .
703 Corra, L. (2022). Pollution and health: a progress update. *The Lancet Planetary Health*.

704 Giani, P., Anav, A., De Marco, A., Feng, Z., & Crippa, P. (2020). Exploring sources of
705 uncertainty in premature mortality estimates from fine particulate matter: the case of
706 China. *Environmental Research Letters*, *15*(6), 064027.

707 Giani, P., Castruccio, S., Anav, A., Howard, D., Hu, W., & Crippa, P. (2020). Short-term and
708 long-term health impacts of air pollution reductions from COVID-19 lockdowns in China
709 and Europe: a modelling study. *The Lancet Planetary Health*, *4*(10), e474-e482.

710 Ha, S., Hui, H., Roussos-Ross, D., Haidong, K., Roth, J., & Xu, X. (2014). Th effects of air
711 pollution on adverse birth outcomes. *Environmental research*, *134*, 198-204.

712 Hajat, A., Hsia, C., & O'Neill, M. S. (2015). Socioeconomic disparities and air pollution
713 exposure: a global review. *Current environmental health reports*, *2*(4), 440-450.

714 Hochberg, Y., & Tamhane, A. C. (1987). *Multiple comparison procedures*. John Wiley & Sons,
715 Inc.

716 Hyder, A., Lee, H. J., Ebisu, K., Koutrakis, P., Belanger, K., & Bell, M. L. (2014). PM_{2.5}
717 exposure and birth outcomes: Use of satellite- and monitor-based data. *Epidemiology*,
718 *25*(1), 58-67.

719 Jaakkola, J. (2003). Case-crossover design in air pollution epidemiology. *European Respiratory*
720 *Journal*, *21*(40 suppl), 81s-85s.

721 Kloog, I., Melly, S. J., Ridgway, W. L., Coull, B. A., & J., S. (2012). Using new satellite based
722 exposure methods to study the association between pregnancy PM_{2.5} exposure, premature
723 birth and birth weight in Massachusetts. *Environmental Health*
724 *18*(11).

725 Koenker, R., & Bassett Jr, G. (1978). Regression quantiles. *Econometrica: journal of the*
726 *Econometric Society*, 33-50.

727 Koenker, R., & Hallock, K. F. (2001). Quantile regression. *Journal of economic perspectives*,
728 *15*(4), 143-156.

729 Liu, J., Clark, L. P., Bechle, M. J., Hajat, A., Kim, S.-Y., Robinson, A. L., . . . Marshall, J. D.
730 (2021). Disparities in air pollution exposure in the United States by race/ethnicity and
731 income, 1990–2010. *Environmental Health Perspectives*, *129*(12), 127005.

732 McClellan, R. O. (2002). Setting ambient air quality standards for particulate matter. *Toxicology*,
733 *181*, 329-347.

734 Mead, M. I., Castruccio, S., Latif, M. T., Nadzir, M. S. M., Dominick, D., Thota, A., & Crippa,
735 P. (2018). Impact of the 2015 wildfires on Malaysian air quality and exposure: a
736 comparative study of observed and modeled data. *Environmental Research Letters*, *13*(4),
737 044023.

738 Miranda, M. L., Edwards, S. E., Keating, M. H., & Paul, C. J. (2011). Making the environmental
739 justice grade: the relative burden of air pollution exposure in the United States.
740 *International journal of environmental research and public health*, *8*(6), 1755-1771.

741 Navidi, W. (1998). Bidirectional case-crossover designs for exposures with time trends.
742 *Biometrics*, 596-605.

743 Nethery, R. C., & Dominici, F. (2019). Estimating pollution-attributable mortality at the regional
744 and global scales: challenges in uncertainty estimation and causal inference. *European*
745 *heart journal*, *40*(20), 1597-1599.

746 Nhung, N. T. T., Amini, H., Schindler, C., Joss, M. K., Dien, T. M., Probst-Hensch, N., . . .
747 Künzli, N. (2017). Short-term association between ambient air pollution and pneumonia
748 in children: A systematic review and meta-analysis of time-series and case-crossover
749 studies. *Environmental Pollution*, *230*, 1000-1008.

750 Pampel, F. C. (2020). *Logistic regression: A primer*. Sage publications.

751 Pope, C. A., Coleman, N., Pond, Z. A., & Burnett, R. T. (2020). Fine particulate air pollution and
752 human mortality: 25+ years of cohort studies. *Environmental research*, *183*, 108924.

753 PRISM Climate Group. (2004). Oregon State University. <https://prism.oregonstate.edu/>. Data
754 accessed September 2021.

755 Qian, D., Amini, H., Shi, L., Kloog, I., Silvern, R., Kelly, J., . . . Lyapustin, A. (2019). An
756 ensemble-based model of PM_{2.5} concentration across the contiguous United States with
757 high spatiotemporal resolution. *Environment international*, *130*, 104909.

758 R Core Team, R. (2023). R: A language and environment for statistical computing. R
759 Foundation for Statistical Computing, Vienna, Austria. URL [https://www.R-](https://www.R-project.org/)
760 [project.org/](https://www.R-project.org/).

761 Rao, K. S. (2005). Uncertainty analysis in atmospheric dispersion modeling. *Pure and applied*
762 *geophysics*, *162*(10), 1893-1917.

763 Servadio, J. L., Lawal, A. S., Davis, T., Bates, J., Russell, A. G., Ramaswami, A., . . . Botchwey,
764 N. (2019). Demographic inequities in health outcomes and air pollution exposure in the
765 Atlanta area and its relationship to urban infrastructure. *Journal of Urban Health*, *96*,
766 219-234.

767 Shen, P., Crippa, P., & Castruccio, S. (2021). Assessing Urban Mortality from Wildfires with a
768 Citizen Science Network. *Air Quality, Atmosphere & Health.*, *Under review*.

769 Son, J.-Y., Lane, K. J., Miranda, M. L., & Bell, M. L. (2020). Health disparities attributable to
770 air pollutant exposure in North Carolina: Influence of residential environmental and
771 social factors. *Health & place*, *62*, 102287.

772 Stuart, A. L., Mudhasakul, S., & Sriwatanapongse, W. (2009). The social distribution of
773 neighborhood-scale air pollution and monitoring protection. *Journal of the Air & Waste*
774 *Management Association*, 59(5), 591-602.

775 Tessum, C. W., Hill, J. D., & Marshall, J. D. (2017). InMAP: A model for air pollution
776 interventions. *PloS one*, 12(4), e0176131.

777 Tessum, C. W., Paolella, D. A., Chambliss, S. E., Apte, J. S., Hill, J. D., & Marshall, J. D.
778 (2021). PM_{2.5} pollutants disproportionately and systemically affect people of color in the
779 United States. *Science Advances*, 7(18), eabf4491.

780 Van Donkelaar, A., Hammer, M. S., Bindle, L., Brauer, M., Brook, J. R., Garay, M. J., . . . Lee,
781 C. (2021). Monthly global estimates of fine particulate matter and their uncertainty.
782 *Environmental Science & Technology*, 55(22), 15287-15300.

783 Van Donkelaar, A., Martin, R. V., Brauer, M., & Boys, B. L. (2015). Use of satellite
784 observations for long-term exposure assessment of global concentrations of fine
785 particulate matter. *Environmental health perspectives*, 123(2), 135-143.

786 WHO. (2006). *Air quality guidelines for particulate matter, ozone, nitrogen dioxide and sulfur*
787 *dioxide - Global update 2005 - Summary of risk assessment*. W. H. Organization.

788 Woo, B., Kravitz-Wirtz, N., Sass, V., Crowder, K., Teixeira, S., & Takeuchi, D. T. (2019).
789 Residential segregation and racial/ethnic disparities in ambient air pollution. *Race and*
790 *social problems*, 11, 60-67.

791 Zani, N. B., Lonati, G., Mead, M., Latif, M., & Crippa, P. (2020). Long-term satellite-based
792 estimates of air quality and premature mortality in Equatorial Asia through deep neural
793 networks. *Environmental Research Letters*, 15(10), 104088.

794