11-1-2016

# The br2 – weighting Method for Estimating the Effects of Air Pollution on Population Health

Goran Krstic
*Fraser Health Authority, New Westminster, BC*, goran.krstic@fraserhealth.ca

Nikolas S. Krstic
*University of British Columbia, Vancouver, BC*

Mauricio Zambrano-Bigiarini
*Universidad de La Frontera, Temuco, Chile*

# The br2 – weighting Method for Estimating the Effects of Air Pollution on Population Health

# The *br²*–weighting Method for Estimating the Effects of Air Pollution on Population Health

**Goran Krstić**
Fraser Health Authority
New Westminster, British Columbia

**Nikolas S. Krstić**
University of British Columbia
Vancouver, British Columbia

**Mauricio Zambrano-Bigiarini**
Universidad de La Frontera
Temuco, Chile

Uncertainties, limitations and biases may impede the correct application of concentration-response linear functions to estimate the effects of air pollution exposure on population health. The reliability of a prediction depends largely on the strength of the linear correlation between the studied variables. This work proposes the joint use of the coefficient of determination, $r^2$, with the regression slope, $b$, as an improved measure of the strength of the linear relation between air pollution and its effects on population health. The proposed $br^2$-weighting method offers more reliable inferences about the potential effects of air pollution on population health, and can be applied universally to other fields of research.

*Keywords:*      Linear regression coefficients, uncertainty analysis, concentration-response function, air pollution, population health

## Introduction

Inherent uncertainties associated with the application of relative risks (RR), hazard ratios (HR) and concentration-response (C-R) functions derived from the epidemiological studies on air pollution exposure vs. population mortality/morbidity have been discussed in the published literature (Burnett et al., 2014; Fann, Gilmore, &Walker, 2013; Fann et al., 2011; Krewski et al., 2009; Environmental Protection Agency, 2006; Post, Watts, Al-Hussainy, & Neubig, 2005; Lipfert & Wyzga, 1995). Considering that confounding factors not controlled or accounted for could affect our ability to predict reliably the effects attributed to a variable of interest (e.g., effects of $PM_{2.5}$ on population health), epidemiological studies often include adjustments for potential impacts from various environmental, behavioral, genetic, and socio-economic health risk factors.

*Dr. Krstić is a Toxicologist and Human Health Risk Assessment Specialist. Email him at: goran.krstic@fraserhealth.ca. Mr. N. Krstić is a student in the Department of Statistics. Dr. Zambrano-Bigiarini is a Professor of Hydrology.*

The coefficient of correlation ($r$) has been developed in its current format by Pearson in 1895 (Rodgers & Nicewander, 1988). The squared value of $r$ is defined as the coefficient of determination ($r^2$), which provides an estimated proportion of the variation in a dependent/response variable $y$ that could be explained by the variation in an independent/explanatory variable $x$. In linear least squares regression with an estimated intercept term, the $r^2$ can be calculated with the following equation:

$$r^2 = \left( \frac{\sum_{i=1}^{n}(O_i - \bar{O})(P_i - \bar{P})}{\sqrt{\sum_{i=1}^{n}(O_i - \bar{O})^2} \sqrt{\sum_{i=1}^{n}(P_i - \bar{P})^2}} \right)^2 \tag{1}$$

where $O$ are the observed and P the predicted values (Krause, Boyle, & Bäse, 2005).

When used for regression between an environmental risk factor vs. population health, the $r^2$ provides a statistical estimate of how well the regression line approximates the real observations. The $r^2$ provides an estimate of the combined dispersion against the single dispersion of the observed and predicted series, with values in the range 0 to 1, where $r^2 = 1$ indicates a perfect linear correlation (i.e., the dispersion values of the observation and the prediction are equal) and $r^2 = 0$ indicates absence of a linear correlation between the studied variables. Refer to Rodgers and Nicewander, (1988) for a set of different ways to express $r$ and conversely the $r^2$.

The coefficient of determination ($r^2$) is sensitive to outliers and extreme dataset values, which may lead to a "bias toward the extreme events if correlation-based measures are employed in model evaluation" (Legates & McCabe, 1999, p. 234). Arnold et al. (2012) indicated the use of $r^2$ without the regression coefficients could be associated with an over-estimation bias and that "if $r^2$ is the primary statistical measure, it should always be used with slope and intercept to ensure that means are reasonable (slope = 1) and bias is low" (p. 1495).

The study by Pope, Ezzati, and Dockery, (2009) could be used as an example to illustrate the importance of $r^2$-value as well as the slope in predicting the effects of PM$_{2.5}$ on population health. Pope, Ezzati, and Dockery (2009, 2012) suggested a reduction in PM$_{2.5}$ concentration observed over the period 1980s – 2000s is responsible for a statistically significant improvement of life expectancy in the metropolitan areas of the United States. However, the observed correlation

with/without the influential observations is very weak ($r^2 \sim 0.05$) and Pope et al. (2012) acknowledged that "given that there are other determinates of life expectancy that may have changed in correlation with changes in air pollution" (p. 234) their analyses "cannot fully eliminate the potential of some residual confounding" (p. 234). This indicates in statistical terms that only approximately 5% of the variation in a change of life expectancy could be explained by the variation in a change of PM$_{2.5}$ concentration and that the remaining 95% could be attributed to a set of selected explanatory variables including income and proxy smoking or other environmental, behavioral, genetic and socio-economic health risk factors not controlled or accounted for in the presented study (e.g., medical practice improvement, public health expenditure change, ambient air temperature).

The focus of the current study is on improving the interpretation of statistical linear regression analyses between air pollution vs. population health. Krause et al., (2005) introduced the application of the regression slope (*b*) as a weighing factor of the coefficient of determination ($r^2$) to address potential under- or over-estimates of model predictions. The proposed method has been used extensively by other researchers in the field of hydrology (Malagò, Pagliero, Bouraoui, & Franchini , 2014; Feaster et al., 2014; Arnold et al., 2012; Zambrano-Bigiarini, 2010; Bellocchi, Rivington, Donatelli, & Matthews, 2009). However, application of this approach in the field of environmental health has been limited (Krstić, 2012; Young & Xia, 2013).

## Methodology

In a comparison of different efficiency criteria for hydrological model assessment, Krause et al., (2005) consider that $r^2$ alone may be limited in its ability to explain the relationship between the response and the explanatory variables, as it quantifies only the dispersion, where "a model which systematically over- or under-predicts all the time will still result in good $r^2$ values close to 1.0 even if all predictions were wrong" (p. 90). Hence, they argue that "for a proper model assessment the gradient b should always be discussed together with $r^2$" (p. 90), and proposed the following model of a weighted coefficient of determination ($_wr^2$) (Krause et al., 2005):

$$_wr^2 = \begin{cases} |b| \cdot r^2 & \text{for } b \leq 1 \\ |b|^{-1} \cdot r^2 & \text{for } b > 1 \end{cases} \qquad (2)$$

The weighted coefficient of determination ($_wr^2$) quantifies under- or over-predictions from both the $r^2$ and the slope or gradient of the regression line ($b$) for a more comprehensive representation of the variable dynamics and model results. In a recently developed R package (R Core Team, 2015) on goodness-of-fit functions for comparison of simulated and observed hydrological time series ("*hydroGOF*"), Zambrano-Bigiarini (2014) indicates "the br$^2$ coefficient allows accounting for the discrepancy in the magnitude of two signals (depicted by 'b') as well as their dynamics (depicted by r$^2$)" (p. 6). Hence, the commutative product of $|b|$ and $r^2$ presented above in (2) can be considered also from the opposite perspective, where $r^2$ is used for weighting the slope/gradient ($b$) to take into account the strength of the linear correlation between the studied variables.

For example, a weak correlation model (e.g., $r^2 < 0.1$) cannot be considered the same as a model with near perfect correlation (i.e., $r^2$ value close to 1.0), which should be taken into account for the interpretation of linear regression analyses by adjusting the slope/gradient ($b$) accordingly:

$$_wb = \begin{cases} r^2 \cdot |b| & \text{for } b \le 1 \\ r^2 \cdot |b|^{-1} & \text{for } b > 1 \end{cases} \tag{3}$$

where $_wb$ represents a weighted slope/gradient ($b$) of the regression line. If $r^2 = 1.0$, in a hypothetical situation of a perfect linear correlation, then $_wb = |b|$ or $_wb = |b|^{-1}$ (i.e., $r^2$ – neutral).

In case of $|b| \le 1$, the limit of $r^2 |b|$ equals 0 if both $|b|$ and $r^2$ approach 0. The same result for the limit of $r^2 |b|$ is obtained if $|b| \to 0$ and $r^2 \to 1$ as well as if $|b| \to 1$ and $r^2 \to 0$:

$$\lim_{\left(|b|,r^2\right) \to (0,0)} \left(r^2 \cdot |b|\right) = \lim_{\left(|b|,r^2\right) \to (0,1)} \left(r^2 \cdot |b|\right) = \lim_{\left(|b|,r^2\right) \to (1,0)} \left(r^2 \cdot |b|\right) = 0 \tag{4}$$

The limit of $r^2 |b|$ equals 1 when both $|b|$ and $r^2$ approach 1:

$$\lim_{\left(|b|,r^2\right) \to (1,1)} \left(r^2 \cdot |b|\right) = 1 \tag{5}$$

In case of $|b| > 1$, the limit of $r^2 |b|^{-1}$ equals 0 if $|b| \to 1$ and $r^2 \to 0$ or if $|b| \to \infty$ and $r^2 \to 0$ or if $|b| \to \infty$ and $r^2 \to 1$:

$$\lim_{\left(|b|,r^2\right)\to(1,0)}\left(r^2\cdot|b|^{-1}\right)=\lim_{\left(|b|,r^2\right)\to(\infty,0)}\left(r^2\cdot|b|^{-1}\right)=\lim_{\left(|b|,r^2\right)\to(\infty,1)}\left(r^2\cdot|b|^{-1}\right)=0 \qquad (6)$$

As in the case of $_wb=r^2\,|b|$, the limit of $r^2\,|b|^{-1}$ equals 1 when both $|b|$ and $r^2$ approach 1:

$$\lim_{\left(|b|,r^2\right)\to(1,1)}\left(r^2\cdot|b|^{-1}\right)=1 \qquad (7)$$

The results of linear regression analyses models can be used to make predictions about the effects of exposure to environmental/socio-economic factors on population health. A linear dose-response model or a linear concentration-response (C-R) function is typically assumed:

$$y=a+bx, \qquad (8)$$

where $y$ is the dependent/response variable, $x$ – independent/explanatory variable, $a$ – the $y$-axis intercept, and $b$ – the slope/gradient of the line. However, it needs to be taken into consideration that the reliability of a prediction made with the aforementioned model depends largely on the strength of the linear correlation between the studied variables, where $r^2$–values greater than ~ 0.5 indicate a strong relationship with high reliability and $r^2$–values less than ~ 0.1 indicate a weak relationship with low reliability of model predictions. This is where the weighted slope/gradient ($_wb$) can be used for a more robust procedure to assess the potential effects of exposure to environmental and/or socio-economic factors on population health.

Using the methodology for particulate matter risk analysis described by the U.S. Environmental Protection Agency (US EPA), Environmental Protection Authority of Victoria (Australia) developed the equations for dose-response or concentration-response (C-R) functions. The authors estimate health outcome changes and calculate the health-endpoint-specific effect coefficient ($\beta$) on the basis of available dose-response data (Burgers & Walsh, 2002).

The C-R functions can be estimated from epidemiological studies using a *Poisson regression* where the natural base logarithm of a health endpoint or an effect is presented as a linear function of air pollution (e.g., $PM_{2.5}$) concentration (Environmental Protection Agency, 2010a):

$$y_1 = y_0 \cdot e^{\beta x_1}, \tag{9}$$

where $y_1$ is the incidence rate of a specific health endpoint of interest at the ambient air pollution concentration ($x_1$), $e$ – the base of natural logarithm (ln or $\log_e$), $\beta$ – the health effect coefficient of ambient air pollution derived from epidemiological studies, and $y_o$ – the baseline incidence rate in hypothetical absence of ambient air pollution, provided that there is no threshold concentration (i.e., level of air pollution below which there is no significant health effect).

The change in the number of cases for a specific health endpoint (e.g., lung cancer incidence or mortality rate) $\Delta y = y_1 - y_o$ or $y_1 = \Delta y + y_o$, corresponding to a given change in ambient air pollution levels relative to the background ($\Delta x = x_1 - x_o$ or $x_1 = \Delta x + x_o$), can be calculated from the C-R function in (9) presented above using the following equation:

$$\Delta y = y_o \left( e^{\beta(\Delta x + x_o)} - 1 \right), \tag{10}$$

where $\beta$ is the health-endpoint-specific effect coefficient representing an incremental change in the health outcome associated with a unit change in air pollution ($\Delta x$). In a hypothetical situation where the background air pollution $x_o = 0$, (10) can be presented as following:

$$\Delta y = y_o \left( e^{\beta \Delta x} - 1 \right) \quad \text{or} \quad \Delta y = y_o \left( RR_{\Delta x} - 1 \right) \tag{11}$$

where the term $e^{\beta \Delta x}$ is also known as the relative risk ($RR_{\Delta x}$) associated with the change in $\Delta x$. If $e^{\beta \Delta x} = RR_{\Delta x}$ then $\beta \Delta x = \ln(RR_{\Delta x})$, and $\beta = \ln(RR_{\Delta x})/\Delta x$.

The percentage change in the number of cases of a given health endpoint ($z_p$), corresponding to a given change in air pollution concentration ($\Delta x$), can be calculated from (Burgers & Walsh, 2002):

$$z_p = \frac{(y_1 - y_o)}{y_o} \cdot 100 \quad \text{or} \quad z_p = \frac{\Delta y}{y_o} \cdot 100 \tag{12}$$

Then, combining and rearranging (11) and (12) provides the equation to calculate $\beta$ for different health endpoints on the basis of available dose-response data from epidemiological studies for a 1 µg/m$^3$ change in air pollution:

$$e^{\beta \Delta x} = 1 + z_p / 100 \quad \text{and} \quad \beta = \frac{\ln\left(1 + z_p / 100\right)}{\Delta x} \tag{13}$$

Finally, an impact of air pollution on a health endpoint can be calculated from the following equation (Fann et al., 2011):

$$\Delta y = y_o \left(e^{\beta \Delta x} - 1\right) \cdot pop, \tag{14}$$

where *pop* is population size of a particular group exposed to air pollution.

Case study data used in the current paper are obtained from Vinikoor-Imler, Davis, and Luben (2011), the National Center for Environmental Assessment of the U.S. EPA, who studied an association between air pollution and population health in North Carolina. They reported the following slopes for PM$_{2.5}$ vs. lung cancer mortality and incidence after adjusting for the neighborhood socio-economic status and the prevalence of cigarette smoking: $b = 0.96$ per 1 μg/m$^3$ PM$_{2.5}$ for lung cancer mortality (95% CI: 0.34, 1.59, $p$-value $< 0.01$; $r^2 = 0.18$; y-axis intercept, $a = 40.96$) and $b = 1.35$ per 1 μg/m$^3$ PM$_{2.5}$ for lung cancer incidence (95% CI: 0.36, 2.35, $p$-value 0.01; $r^2 = 0.09$; y-axis intercept, $a = 44.36$).

## Results

### Case Study Worked Example Calculations: Lung Cancer Mortality

Vinikoor-Imler et al., (2011) provided an adjusted slope of 0.96 lung cancer mortality cases per 100,000 population per 1 μg/m$^3$ change in PM$_{2.5}$ ($b = 0.96 \cdot 10^{-5}$), a y-axis intercept ($a$) or an estimated baseline lung cancer mortality rate at $x_o = 0$ of 40.96 cases per 100,000 population ($y_o = 40.96 \cdot 10^{-5}$), and lung cancer mortality rate per 100,000 population associated with an incremental 1 μg/m$^3$ increase in PM$_{2.5}$ ($y_1 = 0.96 \cdot 10^{-5} + 40.96 \cdot 10^{-5} = 41.92 \cdot 10^{-5}$). Using (12) the value of $z_p$ is calculated at 2.344%. Considering that $y_1 = bx_1 + a$ and $y_o = bx_o + a$, the same calculation can be obtained on the basis of the relationship: $y_1 - y_o = (bx_1 + a) - (bx_o + a)$ or $\Delta y = b\Delta x$, where if $\Delta x = 1$ μg/m$^3$ then $\Delta y = b$ (i.e., 0.96 cases per 100.000 population per 1 μg/m$^3$):

$$z_p = \frac{b}{y_o} \cdot 100 \quad . \tag{15}$$

The C-R coefficient $\beta$ can be then calculated using (13):

$$\beta = \frac{\ln\left(1+2.344/100\right)}{1\left(\mu g/m^3\right)} = 0.0232\left(\mu g/m^3\right)^{-1}.$$

On the basis of the analysis presented by Vinikoor-Imler et al., (2011), using (14), it is estimated that incremental 10 µg/m³ increase in PM$_{2.5}$ concentration could be associated with additional 10.68 cases of lung cancer mortality per 100,000 population (i.e., 34,710 additional cases in ~325 million U.S. population).

In the following estimate of the coefficient $\beta$, the slope of the regression line ($b$) is adjusted for the observed strength of the association between PM$_{2.5}$ exposure and lung cancer mortality ($r^2$) using (13) and (15) with (3), where $|b| = 0.96 \cdot 10^{-5}$ and $r^2 = 0.18$ for a weighted slope/gradient $_wb = 1.728 \cdot 10^{-6}$ per µg/m³ and where $\Delta x = 1$ µg/m³ for $\Delta y = {}_wb$:

$$z_w = \frac{|b| \cdot r^2}{y_o} \cdot 100 \ \text{ or } \ z_w = \frac{_wb}{y_o} \cdot 100$$

$$z_w = \frac{0.96 \cdot 10^{-5} \cdot 0.18}{40.96 \cdot 10^{-5}} \cdot 100 = 0.422\% \ .$$

(16)

A weighted coefficient $\beta_w$ can be then calculated using a weighted percentage increase in the number of cases of a given health endpoint $z_w$ in the following equation:

$$\beta_w = \frac{\ln\left(1+z_w/100\right)}{\Delta x}$$

$$\beta_w = \frac{\ln\left(1+0.422/100\right)}{1\left(\mu g/m^3\right)} = 0.0042\left(\mu g/m^3\right)^{-1}$$

(17)

Hence, adjusting for the neighborhood socio-economic status, cigarette smoking, and the $r^2$ between PM$_{2.5}$ concentration and lung cancer mortality yields a weighted C-R coefficient $\beta_w$ of 0.0042 per µg/m³. Using (14), it is estimated that an incremental 10 µg/m³ increase in PM$_{2.5}$ concentration could be associated with additional 1.76 cases of lung cancer mortality per 100,000 population or 5,720 additional cases if applied to ~325 million U.S. population, which is much lower

than the 34,710 additional cases foreseen by using the unadjusted slope coefficient *b*.

## Case Study Worked Example Calculations: Lung Cancer Incidence

Using the approach described above and the data from Vinikoor-Imler et al., (2011), a weighted C-R coefficient $\beta_w$ is calculated for the cancer incidence where the slope/gradient $b > 1$ (i.e., b = 1.35), $r^2 = 0.09$ and an estimated baseline lung cancer incidence rate $y_o = 44.36$ per 100,000 population at $x_o = 0$. Hence, from (3) a weighted slope/gradient is $_wb = r^2 \cdot /b/^{-1} = 0.09 \cdot 0.7407 \cdot 10^{-5} = 6.666 \cdot 10^{-7}$ per µg/m³ and $z_w$ can be calculated using a modified version of (16) to reflect that $b > 1$:

$$z_w = \frac{r^2 \cdot |b|^{-1}}{y_o} \cdot 100$$

$$z_w = \frac{0.09 \cdot 0.7407 \cdot 10^{-5}}{44.36 \cdot 10^{-5}} \cdot 100 = 0.1503\%$$

(18)

A weighted C-R coefficient $\beta_w$ is calculated using (17):

$$\beta_w = \frac{\ln\left(1 + 0.1503/100\right)}{1\left(\mu g/m^3\right)} = 0.0015\left(\mu g/m^3\right)^{-1}$$

then, using (14), an incremental 10 µg/m³ increase in PM$_{2.5}$ concentration could be associated with additional 0.67 cases of lung cancer incidence per 100,000 population or 2,178 additional cases if applied to ~325 million U.S. population.

## Discussion and Conclusion

Some of the key uncertainties and limitations of currently accepted approach in assessing the effects of air pollution on population health stem from the quality and reliability of epidemiological studies (e.g., study design, exposure assessment, confounding factors, statistical model assumptions, risk characterization, potential errors and biases). The assumptions required for a valid least-squares regression are often not possible to satisfy completely in epidemiological study designs. It should be emphasized that regression coefficient/slope *b* becomes meaningless and should not be used to make linear inferences/predictions if the $r^2$ approaches

0 (e.g., $r^2 < 0.1$) even in situations where it may appear to be statistically significant.

It is also important to consider available evidence for a plausible biological mechanism of toxicity and for a slope and shape of the dose-response relationship at low to very low levels of air pollution (Vedal, Brauer, White, & Petkau, 2003). There is no universal agreement among the researchers for an assumed linear no-threshold effect of air pollution on population health. Specifically regarding $PM_{2.5}$-related mortality the U.S. EPA indicated "a review of the time-series and cohort studies may lead to the conclusion that although a threshold is not apparent at commonly observed concentrations, one may exist at lower levels" (Environmental Protection Agency, 2010b, p. 23). Uncertainties associated with the evidence for and likelihood of causality should be acknowledged. In addition, there is variability in the estimated C-R functions and the magnitude of potential effects of air pollution on population health as reported by different research groups (Environmental Protection Agency, 2010a).

The described methodological approach, first proposed by Krause et al., (2005) in the context of hydrology, was applied by Krstić, (2012) and accepted by Young & Xia, (2013) from the National Institute of Statistical Sciences (NISS) to adjust the predicted population health effects in the context of ambient air pollution. The analyses presented in the current paper on the basis of epidemiological and environmental data from Vinikoor-Imler et al., (2011) showed that inclusion of the $r^2$ in the calculation is expected to yield better estimates of the predicted effects of air pollution on population health, which reflect more accurately the strength of the real linear correlation between the air pollution and the specified population health endpoint.

The proposed $br^2$-weighting method is sensitive to extreme values of both $|b|$ and $r^2$ where model prediction reliability increases if $|b|$ and $r^2$ approach 1 and decreases if $|b|$ departs from 1 in either direction (i.e., $|b| \to \infty$ or $|b| \to 0$) and/or if $r^2$ departs from 1 and approaches 0. The method identifies situations of maximum prediction ability as those of $|b| \leq 1$ as well as for $|b| > 1$, provided that both $|b|$ and $r^2$ are close to 1. This is in agreement with theoretical/ideal conditions in linear regression where a perfect correlation requires that $r = 1$, $|b| = 1$ and $y$-intercept $a = 0$ if the relationship between the studied variables is truly linear in nature, resulting in a 45° angle for the regression line as the best fit of the least-squares estimator (Nau, 2014; Legendre, 2014).

The least-squares regression coefficient $b$ is considered as an unbiased prediction estimator under the assumptions of a perfect correlation between the studied variables (Legendre, & Legendre, 1998). The estimated $r^2$-values closer to

1 allow more direct and reliable application of $b$ in making inferences and predictions. On the other hand, $r^2$-values closer to 0 indicate a necessity to adjust the slope $b$ for the observed reduction in model prediction ability. In situations of very low $r^2$-values, it becomes increasingly more likely even for the 95% confidence interval of the slope $b$ not to include the ideal 45° angle line of the best regression fit (Mesplé et al., 1996; Legendre, 2014).

The presented analyses illustrate the importance of weighting the slope of the regression ($b$) by the coefficient of determination ($r^2$) to obtain more reliable inferences in projecting potential effects of air pollution on population health. The proposed $br^2$-weighting method could be applied universally in studies of other environmental, behavioral, genetic or socio-economic risk factors for more comprehensive health impact estimates with lower potential bias and better decision-making.

## References

Arnold, J. G., Moriasi, D. N., Gassman, P. W., Abbaspour, K. C., White, M. J., Srinivasan, R., … Jha, M. K. (2012). SWAT: Model use, calibration, and validation. *Transactions of the ASABE, 55*(4), 1491-1508. doi: 10.13031/2013.42256

Bellocchi, G., Rivington, M., Donatelli, M., & Matthews, K. (2009). Validation of biophysical models: Issues and methodologies. A review. *Agronomy for Sustainable Development, 30*(1), 109-130. doi: 10.1051/agro/2009001.

Burgers, M., & Walsh, S. (2002). *Exposure assessment and risk characterisation for the development of PM2.5 standard*. Environment Protection Authority of Victoria. Retrieved from: http://www.scew.gov.au/system/files/resources/9947318f-af8c-0b24-d928-04e4d3a4b25c/files/aaq-pm25-rpt-exposure-assessment-and-risk-characterisation-final-200209.pdf (accessed in May 2014).

Burnett, R. T., Pope, C. A., III, Ezzati, M., Olives, C., Lim, S. S., Mehta, S, … Cohen, A. (2014). An integrated risk function for estimating the global burden of disease attributable to ambient fine particulate matter exposure. *Environmental Health Perspectives, 122*(4), 397-403. doi: 10.1289/ehp.1307049.

Environmental Protection Agency. (2006). *Expanded expert judgment assessment of the concentration-response relationship between PM2.5 exposure and mortality*. Research Triangle Park, NC: U.S. Environmental Protection

Agency. Retrieved from:
http://www.epa.gov/ttn/ecas/regdata/Uncertainty/pm_ee_report.pdf (accessed in May 2014).

Environmental Protection Agency. (2010a). *Quantitative health risk assessment for particulate matter* (EPA Publication No. 452/R-10-005). Research Triangle Park, NC: U.S. Environmental Protection Agency. Retrieved from: http://www.epa.gov/ttn/naaqs/standards/pm/data/PM_RA_FINAL_June_2010.pdf (accessed in May 2014).

Environmental Protection Agency. (2010b). *Summary of expert opinions on the existence of a threshold in the concentration-response function for PM$_{2.5}$-related mortality* (Technical Support Document (TSD)). Research Triangle Park, NC: U.S. Environmental Protection Agency. Retrieved from: http://www.epa.gov/ttn/ecas/regdata/Benefits/thresholdstsd.pdf (accessed in July 2014).

Fann, N., Gilmore, E. A., & Walker, K. (2013). *Characterizing the long-term PM$_{2.5}$ concentration response function: comparing the strengths and weaknesses of research synthesis approaches* (Working paper prepared for: Methods for Research Synthesis: A Cross-Disciplinary Workshop). Harvard Center for Risk Analysis. Retrieved from: http://www.hsph.harvard.edu/hcra/files/2013/09/Fann-Gilmore-Walker-Sept-20131.pdf (accessed in May 2014). (I found it at https://cdn1.sph.harvard.edu/wp-content/uploads/sites/1273/2013/09/Fann-Gilmore-Walker-Sept-20131.pdf)

Fann, N., Lamson, A. D., Anenberg, S. C., Wesson, K., Risley, D., & Hubbell, B. J. (2011). Estimating the national public health burden associated with exposure to ambient PM2.5 and ozone. *Risk Analysis, 32*(1), 81-95. doi: 10.1111/j.1539-6924.2011.01630.x.

Feaster, T. D., Benedict, S. T., Clark, J. M., Bradley, P. M., & Conrads, P. A. (2014). *Scaling up watershed model parameters—Flow and load simulations of the Edisto River Basin, South Carolina, 2007-09* (Scientific Investigations Report No. 2014–5104). U.S. Geological Survey. doi: 10.3133/sir20145104.

Krause, P., Boyle, D. P., & Bäse, F. (2005). Comparison of different efficiency criteria for hydrological model assessment. *Advances in Geosciences, 5*, 89-97.

Krewski, D., Jerrett, M., Burnett, R. T., Ma, R., Hughes, E., Shi, Y., … Thun, M. J. (2009). *Extended follow-up and spatial analysis of the American Cancer Society study linking particulate air pollution and mortality* (Research

Report 140). Boston, MA: HEI. Retrieved from:
http://pubs.healtheffects.org/getfile.php?u=478 (accessed in May 2014).

Krstić, G. (2012). A reanalysis of fine particulate matter air pollution versus life expectancy in the United States. *Journal of the Air & Waste Management Association, 62*(9), 989-991. doi:10.1080/10962247.2012.697445.

Legates, D. R., & McCabe, G. J., Jr. (1999). Evaluating the use of "goodness-of-fit" measures in hydrologic and hydroclimatic model validation. *Water Resources Research, 35*(1), 233-241.doi: 10.1029/1998WR900018

Legendre, P. & Legendre, L. (1998). *Numerical ecology* (2nd ed.). Amsterdam: Elsevier.

Legendre, P. (2014). *Model II regression user's guide, R edition*. The Comprehensive R Archive Network (CRAN). Retrieved from: http://cran.r-project.org/web/packages/lmodel2/vignettes/mod2user.pdf (accessed in December 2014).

Lipfert, F. W., & Wyzga, R. E. (1995). Air pollution and mortality: Issues and uncertainties. *Journal of the Air & Waste Management Association, 45*(12), 949-966. doi: 10.1080/10473289.1995.10467427

Malagò, A., Pagliero, L., Bouraoui, F., & Franchini, M. (2014). Comparing calibrated parameter sets of the SWAT model for the Scandinavian and Iberian Peninsulas. *Hydrological Sciences Journal, 60*(5), 949-967. doi: 10.1080/02626667.2014.978332

Mesplé, F., Troussellier, M., Casellas, C., & Legendre, P. (1996). Evaluation of simple statistical criteria to qualify a simulation. *Ecological Modelling, 88*(1/3), 9-18. doi: 10.1016/0304-3800(95)00033-X

Nau, R. F. (2014). *Introduction to linear regression analysis*. Durham, NC: Fuqua School of Business, Duke University. Retrieved from: http://people.duke.edu/~rnau/regintro.htm (accessed in December 2014).

Pope, C. A., III, Ezzati, M., & Dockery, D. W. (2009). Fine-particulate air pollution and life expectancy in the United States. *The New England Journal of Medicine*, 360, 376-386. doi: 10.1056/NEJMsa0805646

Pope, C. A., III, Ezzati, M., & Dockery, D. W. (2012). Validity of observational studies in accountability analyses: The case of air pollution and life expectancy. *Air Quality, Atmosphere & Health, 5*(2), 231-235. doi: 10.1007/s11869-010-0130-3

Post, E., Watts, K., Al-Hussainy, E., & Neubig, E. (2005). Particulate matter health risk assessment for selected urban areas (EPA Publication No. 452/R-05-

007A). Rockville, MD: U.S. Environmental Protection Agency. Retrieved from: http://www.epa.gov/ttn/naaqs/standards/pm/data/PMrisk20051220.pdf (accessed in July 2014).

R Core Team. (2015). R: A language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing. Retrieved from: http://www.R-project.org/.

Rodgers, J. L., & Nicewander, W. A. (1988). Thirteen ways to look at the correlation coefficient. *The American Statistician, 42*(1), 59-66. doi: 10.1080/00031305.1988.10475524

Vedal, S., Brauer, M., White, R., & Petkau, J. (2003). Air pollution and daily mortality in a city with low levels of pollution. *Environmental Health Perspectives, 111*(1), 45-51.

Vinikoor-Imler, L. C., Davis, J. A., & Luben, T. J. (2011). An ecologic analysis of county-level $PM_{2.5}$ concentrations and lung cancer incidence and mortality. *International Journal of Environmental Research and Public Health, 8*(6), 1865-1871. doi: 10.3390/ijerph8061865

Young, S. S., & Xia, J. Q. (2013). Assessing geographic heterogeneity and variable importance in an air pollution data set. *Statistical Analysis and Data Mining, 6*(4), 375-386. doi: 10.1002/sam.11202

Zambrano-Bigiarini, M. (2010). *On the effects of hydrological uncertainty in assessing the impacts of climate change on water resources* (Doctoral thesis). University of Trento, Trento, Italy. Retrieved from: http://eprints-phd.biblio.unitn.it/284/1/MZB-PhD_Thesis-UT-05Ago2010.pdf (accessed in December 2014).

Zambrano-Bigiarini, M. (2014). *Package "hydroGOF": Goodness-of-fit functions for comparison of simulated and observed hydrological time series*. Comprehensive R Archive Network. Retrieved from: http://cran.r-project.org/web/packages/hydroGOF/hydroGOF.pdf (accessed in December 2014)