

Model Uncertainty
and
Health Effect Studies for Particulate Matter

Merlise Clyde



NRCSE

Technical Report Series

NRCSE-TRS No. 027

Model Uncertainty and Health Effect Studies for Particulate Matter

Merlise Clyde

Institute of Statistics and Decision Sciences
210A Old Chemistry Box 90251
Duke University
Durham, NC, 27708-0251

email: clyde@isds.duke.edu

phone: (919) 681-8440

July 23, 1999

This work was supported by NSF grants DMS-96.26135 and DMS-97.33013, and was completed while the author was a visitor at the National Research Center for Statistics and the Environment, University of Washington. Although the research described in this article has been funded in part by the United States Environmental Protection Agency through agreement CR825173-01-0 to the University of Washington, it has as not been subjected to the Agency's required peer and policy review and therefore does not necessarily reflect the views of the Agency and no official endorsement should be inferred.

Model Uncertainty and Health Effect Studies for Particulate Matter

Summary

There are many aspects of model choice that are involved in health effect studies of particulate matter and other pollutants. Some of these choices concern which pollutants and confounding variables should be included in the model, what type of lag structure for the covariates should be used, which interactions need to be considered, and how to model nonlinear trends. Because of the large number of potential variables, model selection is often used to find a parsimonious model. Different model selection strategies may lead to very different models and conclusions for the same set of data. As variable selection may involve numerous tests of hypotheses, the resulting significance levels may be called into question, and there is the concern that the positive associations are a result of multiple testing. Bayesian model averaging is an alternative that can be used to combine inferences from multiple models and incorporate model uncertainty. This paper presents objective prior distributions for Bayesian model averaging in generalized linear models so that Bayesian model selection corresponds to standard methods of model selection, such as the Akaike Information Criterion (AIC) or Bayes Information Criterion (BIC), and inferences within a model are based on standard maximum likelihood estimation. These methods allow non-Bayesians to describe the level of uncertainty due to model selection, and can be used to combine inferences by averaging over a wider class of models using readily available summary statistics from standard model fitting programs. Using Bayesian Model Averaging and objective prior distributions, we re-analyze data from Birmingham, AL and illustrate the role of model uncertainty in inferences about the effect of particulate matter on elderly mortality.

KEYWORDS: AIC; BIC; Bayesian Model Averaging; Jeffreys' Prior; Model Selection; Noninformative Priors; Poisson Regression; PM10

1 Introduction

Statistical analyses of the effect of air pollution on human mortality have been performed for a multitude of cities and a number of authors have found statistically significant relationships between increased mortality in the elderly population (and other health outcomes) and increases in particulate matter (PM). Partly on the basis of such epidemiological studies, the U.S. Environmental Protection Agency in 1997 proposed new stricter standards for PM₁₀ (particulate matter with aerodynamic diameter less than 10 microns). One concern raised by the 1998 National Research Council report on “Research Priorities for Airborne Particulate Matter” is whether the positive associations between particulate matter and mortality (or other health outcomes) are an artifact of model selection due to multiple hypothesis testing.

There are many aspects of model choice that are involved in health effect studies of particulate matter and other co-pollutants. Some of these choices concern which pollutants and confounding variables should be included in the model, what type of lag structure for the covariates should be used, which interactions need to be considered, and what adjustments should be made in the multiple time-series for long-term trends and seasonality. Generalized additive models (GAMs) for Poisson data are often used to model daily mortality, adjusting for nonlinear trends using smoothing splines or other semi-parametric approaches, and including smoothed functions and lags of meteorological variables and pollution variables. Because of the large number of potential variables, it is neither practical nor desirable to include every possible covariate and model selection is often used to find a parsimonious model, balancing reducing autocorrelation and overdispersion against over-fitting the data and inflation of standard errors. This is often done in a highly exploratory fashion, and different model selection strategies may lead to different models and conclusions about the magnitude of relative risks associated with changes in particulate matter. For example, in analyses for Birmingham, AL, Schwartz (1993) found the best model had a relative risk of 1.11, based on 100 $\mu\text{g}/\text{m}^3$ increase in PM₁₀, while in contrast, Davis *et al.* (1996) failed to find any consistent PM₁₀ effect, with estimates of relative risks in the range 1.02-1.05. Davis *et al.* (1996) and Smith *et al.* (1997) found that results were highly sensitive to the choice of meteorological variables and lags of PM₁₀ that were included in the model.

For making inferences, the selected “best” model is often treated as if it were the true model. This procedure ignores the uncertainty involved in model selection, and may lead to overconfident predictions (Draper 1995) and policy decisions that are riskier than one thinks they are (Hodges 1987). One may also find “significant” spurious effects, while the meaning of reported significance levels for the “best” model is also questionable (Viallefont *et al.* 1998). Model uncertainty often outweighs other sources of uncertainty (Hoeting *et al.* 1999), but is typically ignored in general practice.

Bayesian Model Averaging (BMA) using hierarchical models provides a coherent approach for combining predictions and inferences from multiple models, and often leads to improved predictive performance and reduced frequentist risk (Clyde and George 1998, Lamon and Clyde 1998, Hoeting *et al.* 1999). With BMA, predictions and inferences are based on a set of models rather than a single model. For example, predictions are obtained by forming a weighted average of predictions over the different models, where the weights depend on the degree to which the data support each model. All variables are used, but coefficients for variables that are less important are shrunk towards zero. For the health effects models, model averaging can be used to incorporate uncertainty about which of several plausible pollution and meteorological variables are related to mortality, incorporate uncertainty about lag structures in pollution and meteorological variables, and can be used to model nonlinear trends if the data support such effects.

While model averaging is straightforward to implement in theory, model averaging requires specification of prior distributions for parameters within models and prior weights for each model. While subjective prior distributions can be used to incorporate previous knowledge about health effects of particulate matter, this can be controversial because of questions of whose prior beliefs are being represented. In this situation, it may be desirable to have reference analyses based on “non-informative” prior distributions to supplement analyses based on subjective prior distributions. Also because of the complexity of generalized additive models and the number of confounding variables, even carefully elicited prior distributions may have unforeseen consequences on model selection. As part of an overall sensitivity analyses, objective prior distributions play an important role. By incorporating model uncertainty and considering a range of objective prior distributions, we can, perhaps,

increase our confidence that positive associations are not an artifact of model selection.

In this paper, data from Birmingham, AL are re-analyzed using Bayesian Model Averaging (BMA) in conjunction with generalized additive models to assess the impact of model uncertainty on estimates of relative risks due to changes in PM_{10} . As in Schwartz (1993) and Davis *et al.* (1996), the response variable is non-accidental mortality. Additional information on the data and variables is given in section 2. In section 3, we describe the hierarchical Poisson regression model for model averaging. In section 4, we present a class of objective prior distributions for generalized linear models (GLMs) so that 1) Bayesian model selection can be calibrated to standard methods of model selection, such as AIC (Akaike Information Criterion; Akaike 1973), BIC (Bayes Information Criterion; Schwarz 1978) and other methods, and 2) inferences within a model are based on maximum likelihood theory. Such methods can be implemented using readily available summary statistics from standard model fitting programs. Model averaging using objective prior distributions provides a bridge between classical and Bayesian methods of estimation, and presents a natural framework so that non-Bayesians can make inferences from multiple models via model averaging. In Section 5, we apply Bayesian model averaging with several objective prior distributions to the Birmingham data. We construct posterior distributions for relative risks based on a $100 \mu g/m^3$ increase in PM_{10} . These distributions incorporate model uncertainty as well as parameter uncertainty. We also conduct a small model validation study to compare model selection and model averaging under BIC and AIC prior distributions.

2 Variables

The data used in this analysis were originally constructed by Davis *et al.* (1996) and were based on daily measurements from 1985–1988 of mortality (from the National Center for Health Statistics), PM_{10} (from the U.S. Environmental Protection Agency, EPA), and meteorology variables (from the U.S. National Climatic Data Center in Ashville, NC). Variable names and descriptions are given in Table 1.

[Table 1 here]

The response variable for this analysis is daily elderly non-accidental mortality, which is

defined as the total number of deaths on a given day of all individuals age 65 and older, excluding deaths attributed to accidental causes. While Schwartz (1993) used total non-accidental mortality, the number of non-accidental deaths in individuals under 65 averages around 2 per day, and should not impact conclusions greatly.

PM₁₀ data are available from the EPA's aerometric data base for monitors in 8 locations in Jefferson County, AL, which contains the metropolitan area of Birmingham (Figure 1). These consist of readings from a daily monitor located in Birmingham (monitor ID 0023), a daily monitor in Leeds (monitor ID 1010), and several monitors (ID's 0002, 0012, 0026, 2003, 3003, 6002) throughout the county that collected data every 6 days. Schwartz (1993) and Davis *et al.* (1996) used the daily mean of PM₁₀ data from all available PM₁₀ monitors within the metropolitan area to construct a daily area-wide measure of PM₁₀. Because the daily monitor in Leeds collected data only for the latter half of the time period and exhibited lower values on average than the daily monitor in Birmingham (Figure 2), the daily area-wide average varies substantially depending on whether the Leeds data are included. To avoid a bias in the area-wide average because the Leeds data are not missing at random during the first half of the time period, we used PM₁₀ data from the daily monitor within Birmingham (Monitor ID 0023), which started operation in August, 1985 and was in operation until the end of 1988. We will explore the difference in results based on using the area-wide average (pma) versus a single daily monitor 0023 (pm). Uncertainty in which monitors are representative of population exposure is an important question and an open area for research.

[Figures 1 and 2 here]

After exploratory modeling, Schwartz (1993) reported that the average of PM₁₀ from the three previous days was the best predictor of mortality. Part of the difference in the results between Schwartz and Davis *et al.* can be explained by how the three-day averages were constructed. While Schwartz used the three previous days, Davis *et al.* used the current day, and the previous two days. Smith *et al.* (1997) investigated using individual lags rather than 3-day averages of PM₁₀. In order to take into account uncertainty in the lag structure, we constructed up to three day lags of PM₁₀, where lag 0 is the current day.

The meteorological variables and lags (up to 2 days) used are as defined by Davis *et al.* (1996) and are summarized in Table 1. Several of the variables originally listed by

Davis *et al.* (1996) are deterministic functions of the other variables (Cooling and Heating Degree Days, and Apparent Temperature for Heat/Cold Stress) and have been excluded from consideration. As low wind speeds are confounded with high pollution events, but not expected to be associated with mortality (Schwartz 1993), average wind speed was also eliminated from consideration.

The mortality data have a strong seasonal component that may be related to flu epidemics or other causes and longer term changes in population size. As long term trends and yearly variation in the covariates are confounded with long term and yearly variation in mortality, it is necessary to remove this source of variation from the analysis prior to assessing the effect of PM_{10} . This is often done by adding a smooth function of time to the model using smoothing splines (Smith *et al.* 1997) or sine-cosine functions (Schwartz 1993). The Time Series Working Group at the NRCSE Particulate Methodology Workshop held in Seattle, Oct 19-22,1998 (URL <http://www.nrcse.washington.edu/events/pm-workshop.html>) concluded “While the method used to remove longer-term variation is unlikely to be important, the choice of which time scales to include and which to exclude may influence the results. Removing too little information exposes the analysis to confounding by season, removing too much reduces the power of the analysis and may exclude important health effects.” In general, large scale variation that should be removed is on a time scale of one month or greater.

We have found that using cubic splines or thin-plate splines leads to little difference in the smoothed long term trend, but the number of knots does have a large impact on smoothness of the unknown trend. In what follows we report results using a thin-plate spline basis. To construct this basis, we selected 30 knots at equally spaced time points over the length of the sampling period. Let k_j denote a knot at location j , where $0 < k_j < n$. The j th basis element evaluated at the point t is constructed as

$$b_j(t) = (t - k_j)^2 \log(|t - k_j|).$$

A function $f(t)$, where t is the time index, representing the unknown trend can be represented as

$$f(t) = \alpha_0 + \sum_{j=0}^K \alpha_j b_j(t).$$

Removing knot j is equivalent to setting an α_j to zero. As the dimension or number of knots is unknown, this is also a variable selection problem; model uncertainty in the number of knots can be addressed using BMA. For more discussion of Bayesian approaches to function estimation using thin-plate or other radial bases see Holmes and Mallick (1997).

3 Hierarchical Poisson Regression Model

Daily mortality reflects counts which are usually modeled with a Poisson or over-dispersed Poisson distribution. We will let Y_i denote the non-accidental elderly mortality for day i , $i = 1, \dots, n = 1247$ and let $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p]$ denote the $(n \times p)$ design matrix based on all variables under consideration. The design matrix \mathbf{X} can include basis terms for smoothing splines to model nonlinear trends, meteorological variables, such as temperature and humidity, pollution variables, lags of meteorological and pollution variables, seasonal indicators, or any other known confounders. We will focus on the set of variables listed in Table 1.

Under the full model, we assume that the observations $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)'$ are independent Poisson random variables with means $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)'$ and that the means are related to the covariates via a link function,

$$g(\boldsymbol{\mu}) = \mathbf{X}\boldsymbol{\beta}$$

where the canonical link function g is the log link. In the present context of variable selection, models correspond to different probability specifications for the data, so that under the m th model (\mathcal{M}_m)

$$\mathcal{M}_m : \quad \mathbf{Y} \sim \text{Poisson}(\boldsymbol{\mu}) \quad \log(\boldsymbol{\mu}) = \mathbf{X}_m\boldsymbol{\beta}_m$$

where \mathbf{X}_m is the design matrix under model \mathcal{M}_m and $\boldsymbol{\beta}_m$ is the vector of regression coefficients for model \mathcal{M}_m . The set of possible models is given by $\mathcal{S} = \{\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_M\}$ and includes all subsets of potential covariates, or $M = 2^p$ models.

3.1 Hierarchical Model

Model uncertainty can be formally accounted for by building an expanded model that encompasses all models in \mathcal{S} . In constructing the hierarchical model, it is convenient to introduce

a p dimensional vector $\boldsymbol{\gamma}$ of indicator variables where γ_j equals one if variable \mathbf{x}_j is included under \mathcal{M}_m . To link \mathcal{M}_m and $\boldsymbol{\gamma}$, we take $\boldsymbol{\gamma}$ to be the binary representation of m .

The hierarchical model is defined in three stages. The first stage of the hierarchy defines the distribution for the data in terms of all variables:

$$\mathbf{Y}|\boldsymbol{\beta}, \mathcal{M}_m \sim \text{Poisson}(\exp(\mathbf{X}\boldsymbol{\beta})) \quad (1)$$

with density $f(\mathbf{Y}|\boldsymbol{\beta}, \mathcal{M}_m) = f(\mathbf{Y}|\boldsymbol{\beta})$. In the expanded model, variables are eliminated by allowing coefficients to be exactly 0. This is achieved by allowing point masses at zero in the prior distribution for $\boldsymbol{\beta}$ given \mathcal{M}_m in the second stage. Coefficients for variables not in \mathcal{M}_m are identically zero as specified through distributions $\delta_0(\beta_j)$ that are degenerate at zero, while coefficients for variables included under \mathcal{M}_m ($\boldsymbol{\beta}_m$) have non-degenerate prior distributions, so that the joint distribution for $\boldsymbol{\beta}$ given \mathcal{M}_m is

$$p(\boldsymbol{\beta}|\mathcal{M}_m) = p(\boldsymbol{\beta}_m|\mathcal{M}_m) \prod_{j=1}^p \delta_0(\beta_j)^{1-\gamma_j} \quad (2)$$

(here $\boldsymbol{\beta}_m$ corresponds to the elements of $\boldsymbol{\beta}$ where γ equals 1). The last stage of the hierarchical model assigns prior weights to each of the models,

$$\mathcal{M}_m \sim \pi(\mathcal{M}_m) \quad (3)$$

which are often taken to be uniform *a priori*. By collapsing the last two stages, the marginal prior distribution for $\boldsymbol{\beta}$ is a mixture of point mass and continuous distributions defined on $(\{0\} \cup (-\infty, \infty))^p$. Likewise, the posterior distribution is also a mixture distribution, reflecting model uncertainty.

3.2 Posterior Distributions

Using Bayes Theorem, the posterior probability of model \mathcal{M}_m is

$$\pi(\mathcal{M}_m|\mathbf{Y}) = \frac{f(\mathbf{Y}|\mathcal{M}_m)\pi(\mathcal{M}_m)}{\sum_{k=1}^M f(\mathbf{Y}|\mathcal{M}_k)\pi(\mathcal{M}_k)} \quad (4)$$

where the marginal distribution of the data is

$$f(\mathbf{Y}|\mathcal{M}_m) = \int f(\mathbf{Y}|\boldsymbol{\beta}, \mathcal{M}_m)p(\boldsymbol{\beta}|\mathcal{M}_m)d\boldsymbol{\beta} \quad (5)$$

given model \mathcal{M}_m , which provides a measure of how much the data support each model.

Under BMA, the distribution of quantities of interest, Δ , such as future mortality or relative risks, is represented as a mixture distribution,

$$f(\Delta|\mathbf{Y}) = \sum_{m=1}^M f(\Delta_m|\mathbf{Y}, \mathcal{M}_m)\pi(\mathcal{M}_m|\mathbf{Y}) \quad (6)$$

where the model specific distributions are weighted by the posterior model probabilities.

3.3 Prior Distributions for BMA

A key component to model averaging is the prior distribution for the parameters and models; how should we choose $p(\boldsymbol{\beta}_m|\mathcal{M}_m)$ and $\pi(\mathcal{M}_m)$? Subjective prior distributions may be difficult to elicit in large problems, especially when there are complicated interactions among variables. Robustness of the prior is also a concern. In complicated models, aspects of the prior distribution that have not received careful attention may lead to undesirable behaviour in the posterior distribution (Berger 1985). In linear regression, $p(\boldsymbol{\beta}_m|\mathcal{M}_m)$ is often based on a normal distribution with mean zero, and prior covariance $\tau(\mathbf{X}_m'\mathbf{X}_m)^{-1}$ and $\pi(\mathcal{M}_m)$ is a uniform distribution over \mathcal{S} , so that all models are equally likely *a priori*. Similar priors are also often used in generalized linear models (Raftery 1996). In linear regression, this form is often selected out of convenience because the calculations under conjugate distributions lead to closed form solutions for the posterior distributions. While τ is often based on considerations of the range of $\boldsymbol{\beta}$, τ also has a strong impact on model selection, which is often not taken into consideration. In fact, proper, but vague “non-informative” priors obtained by taking τ large, can lead to Bayes factors favoring the null model, even in situations where the parameters estimates may be far from 0.

While a careful subjective Bayesian analysis is ideal, we argue that even when subjective information is available, it may be desirable to present results based on objective prior distributions to accompany subjective analyses to check sensitivity of results to prior specifications. The Schwartz criterion or BIC is appealing in that it can be applied even when the priors $p(\boldsymbol{\beta}_m|\mathcal{M}_m)$ are hard to specify precisely, provides a reference procedure for scientific reporting, and can be normalized to provide weights for BMA (Kass and Raftery 1995). While AIC is another commonly used default procedure for model selection, it has not been

used in BMA. In the next section, we describe objective prior distributions based on Jeffreys’ prior and discuss how to calibrate such priors based on standard model selection criteria, such as AIC and BIC. This provides a range of objective prior and posterior distributions for use in BMA and scientific reporting, requiring varying degrees of evidence.

4 Objective Prior Distributions

In estimation problems with vector valued parameters, Jeffreys (1961) suggested taking

$$p(\boldsymbol{\beta}_m|\mathcal{M}_m) = |\mathcal{I}(\boldsymbol{\beta}_m)|^{1/2} \quad (7)$$

where $|\mathcal{I}(\boldsymbol{\beta}_m)|$ is the determinant of the expected Fisher information matrix. In many problems, this leads to an improper distribution, which is determined only up to a multiplicative constant. We can always multiply an improper prior as in (7) by a constant a_m and still have a “valid” improper prior distribution. The constants do not affect the posterior distribution of $\boldsymbol{\beta}$ given \mathcal{M}_m , but are present in the marginal likelihood, and thus Bayes factors or posterior model probabilities contain the undefined constants.

We will show that specific choices of these constants yield what we will call Calibrated Information Criterion (CIC) prior distributions for generalized linear models, and that these CIC prior distributions can be used reconcile classical model selection and Bayesian model selection based on posterior model probabilities. Other approaches for specifying the constants based on imaginary training are in Spiegelhalter and Smith (1982). The CIC priors provide a general framework for model averaging, as opposed to model selection, when model uncertainty is an issue.

4.1 Calibrated Information Criterion Prior Distributions

For generalized linear models, we define the CIC prior distribution as

$$p(\boldsymbol{\beta}|\mathcal{M}_m)\pi(\mathcal{M}_m) = (2\pi)^{-d_m/2} \left| \frac{1}{c} \mathcal{I}(\hat{\boldsymbol{\beta}}_m) \right|^{1/2} \prod_{j=1}^p \delta_0(\beta_j)^{1-\gamma_j} \quad (8)$$

where d_m is the dimension of model \mathcal{M}_m and $\mathcal{I}(\hat{\boldsymbol{\beta}}_m)$ is the observed Fisher information for \mathcal{M}_m evaluated at the MLE's $\hat{\boldsymbol{\beta}}_m$ with j, k elements,

$$[\mathcal{I}(\boldsymbol{\beta}_m)]_{jk} = - \left[\frac{\partial^2}{\partial \beta_j \partial \beta_k} \mathcal{L}(\boldsymbol{\beta} | \mathcal{M}_m) \right]$$

and $\mathcal{L}(\boldsymbol{\beta} | \mathcal{M}_m)$ denotes the log likelihood under \mathcal{M}_m . This form automatically ensures that the prior distribution takes into account the link function. For the Poisson regression model with the log link the observed Fisher information is

$$\mathcal{I}(\hat{\boldsymbol{\beta}}_m) = \mathbf{X}'_m V(\hat{\boldsymbol{\beta}}_m) \mathbf{X}_m \quad (9)$$

where $V(\boldsymbol{\beta}_m)$ is the covariance matrix for \mathbf{Y} with elements $\exp(\mathbf{X}_m \boldsymbol{\beta}_m)$ on the diagonal and 0 elsewhere. For the canonical link, the observed and expected Fisher information are the same.

4.2 CIC Posterior Distributions

For the Poisson regression models under consideration we cannot obtain the marginal likelihood of the data (5) analytically. Laplace's method (Tierney and Kadane 1986) provides a useful approximation to the marginal likelihood as long as the likelihood is peaked near its maximum, which will be the case for large samples. Kass and Raftery (1995) have found that Laplace approximations for determining posterior model probabilities are accurate for sample sizes on the order of $20p$ or larger. As in classical approximations for obtaining the distribution of MLE's, we replace $\mathcal{L}(\boldsymbol{\beta} | \mathcal{M}_m)$ by a 2nd order Taylor series expansion about $\hat{\boldsymbol{\beta}}$, so that under \mathcal{M}_m

$$\tilde{\mathcal{L}}(\boldsymbol{\beta} | \mathcal{M}_m) = \mathcal{L}(\hat{\boldsymbol{\beta}} | \mathcal{M}_m) - \frac{1}{2} (\boldsymbol{\beta}_m - \hat{\boldsymbol{\beta}}_m)' \mathcal{I}(\hat{\boldsymbol{\beta}}_m) (\boldsymbol{\beta}_m - \hat{\boldsymbol{\beta}}_m).$$

Using $\exp(\tilde{\mathcal{L}}(\boldsymbol{\beta} | \mathcal{M}_m))$ in place of the actual likelihood, the (approximate) joint posterior for $\boldsymbol{\beta}_m$ and \mathcal{M}_m factors as

$$p(\boldsymbol{\beta}_m | \mathbf{Y}, \mathcal{M}_m) = N(\hat{\boldsymbol{\beta}}_m, \mathcal{I}(\hat{\boldsymbol{\beta}}_m)^{-1}) \quad (10)$$

$$\pi(\mathcal{M}_m | \mathbf{Y}) = \frac{\exp \left\{ \frac{1}{2} (D_m - d_m \log(c)) \right\}}{\sum_{m=1}^M \exp \left\{ \frac{1}{2} (D_m - d_m \log(c)) \right\}} \quad (11)$$

where D_m is the model deviance, which is the usual deviance (-2 times the log likelihood) under the null model minus the deviance under \mathcal{M}_m .

The log of the posterior model probability under the CIC prior distribution is

$$\log \pi(\mathcal{M}_m|\mathbf{Y}) = k + \frac{1}{2}(D_m - d_m \log(c))$$

which is proportional to the model deviance minus a model complexity penalty term depending on $\log(c)$. Posterior model probabilities under the CIC priors can be calibrated to classical model selection through the choice of $\log(c)$ (Table 2) for popular methods such as AIC (Akaike 1973, 1978), BIC (Schwarz 1978) and RIC (Foster and George 1994). For the values of $\log(c)$ in Table 2, the model with the highest posterior probability under that prior corresponds to the optimal CIC model using a penalized deviance criterion for model selection. Using the CIC prior distributions, posterior inference within a model is based on standard normal approximations in GLMs given by (10), but the CIC posterior model probabilities can also be used to incorporate model uncertainty using (6).

[Table 2 here]

While much of statistical practice has focused on selecting a single model, this may be unreasonable unless $\pi(\mathcal{M}_k|\mathbf{Y})$ is near one for one of the models under consideration or the best models provide similar inferences; even if the posterior probability of the best model is near one, there is little loss by using BMA as the best model dominates the mixture. Raftery (1996) used BMA based on BIC in GLMs, however, the usual justification of AIC provides no way of taking into account model uncertainty (Kass and Raftery 1995). The CIC prior distributions provide a justification for model averaging using AIC and other classical model selection criteria based on penalized deviance criteria of the form above.

The BIC and AIC priors provide a broad range for sensitivity analyses with BMA. AIC often favors models that are more complex than models selected by BIC, and tends to overestimate the number of parameters needed, even asymptotically (see discussion in Kass and Raftery 1995). BIC is a more conservative strategy and requires much stronger evidence to reject the null hypothesis, and hence often puts more weight on simpler models. For example in a one dimensional testing problem, t^2 statistics in favor of the alternative under AIC correspond to values greater than 2 while for BIC, t^2 values must be greater than

$\log(n)$. A high probability of variable inclusion under both AIC and BIC priors implies consistent strong support for an effect, while if both AIC and BIC lead to small probabilities of variable inclusion, this provides consistent support in favor of the null. Situations where the probability of variable inclusion is high under AIC, but low under BIC, require careful consideration of prior distributions, and may indicate that the sample is not large enough to detect an effect. Other choices of c may be more appropriate, and can be calibrated based on considerations of practical significance and costs/losses associated with decisions of accepting different hypotheses.

4.3 Implementing Model Averaging

For linear regression models with conjugate prior distributions and a small to moderate number of covariates (less than 20), posterior distributions for many quantities can be determined analytically (George and McCulloch 1997). For larger problems, we typically cannot enumerate all models so model averaging is approximated by using a sample of models from \mathcal{S} . Stochastic search using Markov Chain Monte Carlo (MCMC) methods or deterministic search methods such as leaps and bounds can be used to identify a sample of models that are used in BMA (see George and McCulloch 1997, Clyde 1999 and Hoeting *et al.* 1999 for discussion of approaches and methods for implementing BMA in the context of linear and generalized linear models). For large problems with highly correlated variables (such as the meteorological variables) using transformations based on factor analyses or principal components may lead to improved convergence with MCMC methods (Clyde 1999, Clyde and DeSimone-Sasinowska 1997).

For the application in the next section, we modified the `bic.glm` code (available on the BMA homepage, URL <http://att.research.com/~volinsky/bma.html>), written by C. Volinsky. This uses the leaps and bounds algorithm to provide a preliminary list of models for use with the objective CIC prior distributions; if a more detailed analysis including subjective information or other proper priors is warranted, then one may later implement Monte Carlo or MCMC methods to provide posterior samples for making inferences.

5 Results for Birmingham

We explore the use of the CIC prior distributions for assessing the effect of particulate matter on elderly mortality in Birmingham, AL and what impact model uncertainty may have on decisions. We use the Poisson model with canonical log link given in (1) and consider possible models based on the variables listed in Table 1.

As the combined design matrix for the thin-plate spline basis and meteorological and PM_{10} variables contains more than 30 variables, and the leaps and bounds code in SPlus is limited to 30 variables, one must either eliminate variables or use a multistage procedure. We implemented the model search in two stages: the first stage was used to estimate the non-parametric baseline trend using thin-plate smoothing splines; the second stage included the posterior mean of the baseline trend estimated in Stage 1 and all meteorological and PM_{10} variables listed in Table 1.

Figure 3 illustrates the degree of model uncertainty in the estimates of the baseline trend. The thin solid line is the GLM estimate under the full model with all 30 knots. This closely tracks a number of the high mortality episodes which may be a result of other (short-term) factors. The dashed lines correspond to the top 100 models under the CIC posterior using $c = n$, with the thick solid line corresponding to the predictive mean under BMA. This provides an objective estimate of the baseline trend, without subjective assessment of the number of knots, and appears to capture the necessary long-term variation without overfitting. We incorporate the BMA estimate of the baseline trend as a linear predictor in the model; this can be thought of as an independent underlying baseline estimate as in proportional hazards models. While this two stage approach ignores uncertainty in the baseline estimate, we can later account for the additional uncertainty by repeating Stage 2 for a number of the top models from Stage 1, or by refitting models identified in Stage 2 under the baseline models identified in Stage 1.

[Figure 3 here]

In Stage 2, 7860 models were selected by the leaps and bounds approach for use in BMA. We used all 860 of the observations that had complete records for all of the lagged pm variables, so that with all lags and the BMA baseline estimate the design matrix for the full

model included 29 candidate predictors.

Relative risks for each model were based on a simultaneous $100 \mu\text{g}/\text{m}^3$ increase in all PM_{10} variables (pm0 , pm1 , pm2 , pm3)

$$R = \exp(100(\beta_{\text{pm0}} + \beta_{\text{pm1}} + \beta_{\text{pm2}} + \beta_{\text{pm3}})). \quad (12)$$

The posterior distribution for the relative risk given \mathcal{M}_m was based on a standard normal approximation (using the Delta method) with mean based on (12) evaluated at the MLE (posterior mode) under \mathcal{M}_m , and variance,

$$\sigma_R^2 = \left(\frac{\partial R}{\partial \boldsymbol{\beta}} \right)' \Sigma_{\boldsymbol{\beta}|\mathcal{M}_m} \left(\frac{\partial R}{\partial \boldsymbol{\beta}} \right) \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}|\mathcal{M}_m}$$

where $\Sigma_{\boldsymbol{\beta}|\mathcal{M}_m}$ is the covariance matrix for $\boldsymbol{\beta}$ under model \mathcal{M}_m with a block for $\boldsymbol{\beta}_m$ that is the inverse of the observed Fisher information matrix under \mathcal{M}_m given in (9), and zero entries elsewhere for the degenerate components for variables not included under \mathcal{M}_m . More accurate approximations to the posterior distribution can be obtained by using Laplace approximations (Tierney and Kadane 1986) or by using MCMC sampling, however, under the Delta method approximations, Bayesian probability intervals correspond to classical confidence intervals used in other analyses.

Figure 4 links plots of the model space and corresponding MLE's of relative risks (posterior modes) under the top 25 models under the BIC and AIC priors. The model space is represented as a matrix, where rows correspond to models and columns to variables. In the model matrix, a black square in position jk indicates that the k th variable is not included in the j th model. The models are ordered from best (at the top) to worst (on the bottom) with the scale on the y-axis reflecting the log of ratio of the model probability of the best model to the worst model in the entire sample of models. A difference of 2 or less in this scale indicates that the top 25 models are more or less exchangeable.

[Figure 4 here]

The last 4 columns in the model space matrix correspond to the PM_{10} variables: pm0 - pm3 , the current day and 3 previous days. As the figure shows, the top models under the AIC and BIC priors are very different regarding the inclusion of PM_{10} : PM_{10} variables are included in all of the top AIC models, but in roughly one third of the top 25 BIC models.

The y-axis in the relative risk plots in Figure 4 also corresponds to models with the best model at the top, thus providing a static link between the model space, which identifies the confounding variables, and the corresponding relative risk estimates. Clearly models where all the coefficients for PM_{10} are zero have relative risks of 1. The points indicate the mean relative risk under each model, while the horizontal lines correspond to 95% probability (confidence) intervals. While the top 25 AIC models include PM_{10} , most of the probability (confidence) intervals for relative risk include one.

Figure 5 illustrates the distribution of relative risk over all sampled models that included PM_{10} in the model for both the BIC and AIC priors. The histograms at the top of the figure show the distributions of the MLE's of relative risks weighted by posterior model probabilities. This illustrates model uncertainty in the (point) estimates, but ignores parameter uncertainty. The range in the estimates (indicated by the dots) is almost as large as the length of individual probability intervals in Figure 4. The histograms at the bottom of Figure 5 are based on samples of relative risks from the posterior distribution under model averaging, which incorporates both model uncertainty and parameter uncertainty. While the range and values of the relative risks are the same under both the AIC and BIC posteriors (the same sampled models and realizations are used to reduce Monte Carlo variation), the weights used in constructing the histograms depend on the posterior model probabilities and priors.

[Figure 5 here]

Table 3 summarizes relative risks using the AIC and BIC prior distributions under model averaging. While posterior means for the relative risk given that PM_{10} (or some lag of PM_{10}) is included in the model are comparable under the BIC and AIC priors, (both approximately 1.05), the uncertainty over whether PM_{10} variables should be included in the model depends greatly on whether the AIC or BIC priors are used. Under AIC, the probability that the relative risk is 1 given the data is 0.03, while under the BIC prior the probability that the relative risk is 1 given the data is 0.72. This in turn impacts the overall estimate of the relative risk averaged over all sampled models (1.052 for AIC versus 1.015 for BIC). Because of the asymmetry of the mixture distribution for relative risk under BMA, symmetric confidence intervals are less appropriate; the highest posterior density interval under the BIC prior is

(0.99 - 1.11) while for the AIC prior the interval is (0.94 - 1.17). While the overall mean relative risk is higher under BMA with the AIC prior, the associated level of uncertainty is also greater, resulting in a wider probability interval.

[Table 3 here]

We also estimated relative risks using the area-wide average of PM₁₀ and lags, and found that the estimates of relative risk under BMA were 1.02 for AIC and 1.009 for BIC. For comparison, Schwartz (1993) obtained a relative risk of 1.11, while Davis *et al.* (1996) had relative risks in the range of 1.02 to 1.05. While the results of Schwartz are plausible (there are models that have relative risks in this range), results under both the AIC and BIC priors appear to support lower estimates of relative risks, with intervals that contain 1. While the AIC prior results in more support for the hypothesis that the relative risk is not 1, both approaches indicate that the increase in mortality is likely between 2 and 5 percent.

We may fail to reject the null hypothesis of “no effect” of PM₁₀ because either there is not enough data to detect an effect, or the data actually support the null hypothesis; one advantage of BMA over the use of traditional p-values is that posterior probabilities can distinguish between these two situations (Hoeting *et al.* 1999). The probability of an effect is equivalent to the posterior probability that the relative risk not 1. For the BIC prior the probability is 0.28, so that the data are indecisive, while under the AIC prior, the probability is 0.97. In both cases, the data do not provide evidence in favor of no PM₁₀ effect.

5.1 Validation Study

As the importance of the effect of PM₁₀ on mortality depends on the choice of prior, we carried out a small validation study to compare the predictive performance of Bayesian model averaging and model selection under the AIC and BIC prior distributions. For this we randomly selected 75 days with complete PM₁₀ data, and repeated the BMA analysis described in the previous section using the remaining data. For the 75 days in the validation set, we computed the predictive MSE,

$$\text{MSE} = \sum_{i=1}^{75} (Y_i - \hat{Y}_i)^2 / 75$$

where \hat{Y}_i is the predictive mean of daily elderly non-accidental mortality. We calculated predictive means and MSE's under BMA with the AIC and BIC priors, and for the best AIC and BIC models (Table 4). The predictive comparison favors model averaging under the BIC priors, with a gain in efficiency over the best AIC model of just over 5%. Similar results were obtained with other randomly selected validation sets. While this approach favors the simpler models under BIC and BMA, little research has been conducted on validating PM₁₀ mortality health effect models, and this is an open area for design of appropriate methods.

[Table 4 here]

6 Discussion

Using the methodology presented here, BMA can become a routine part of exploratory data analysis as most quantities of interest are based on standard output from GLM packages. The distributions of relative risk under BMA require generation of relative risks from the posterior distribution, but this is straightforward under most higher level statistical packages. The results are all conditional on the collection of models used, which does require special consideration. As we proceed, we may find it necessary to enlarge the class of models under consideration by allowing for overdispersion, nonlinear effects, or interactions. Bayesian model averaging can be applied sequentially as new information or variables become available.

The analyses presented illustrate how model averaging can be used in health effect studies using prior distributions based on AIC and BIC. In this example, AIC tends to favor models with more complexity while the top models under BIC are much simpler. This leads to wider probability intervals for relative risk under the AIC prior. While the predictive model validation suggests that BIC is better calibrated than AIC, model averaging using the AIC prior performs better than model selection using AIC. However, rather than trying to argue that BIC is better than AIC, one of our goals here is to provide an objective list of potential models ranked by AIC and BIC (or other CIC priors) as part of a sensitivity analysis. We hope that by exploring linked plots we can gain a better understanding of confounding variables. Given the collection of models, one can objectively discuss the relative merits

of these models and priors. In this example, it is clear that the relative risk estimates vary among the top models (particularly under AIC), and that the importance of the PM_{10} effect depends on whether the AIC or BIC prior is used. Ideally, we can also construct a calibrated prior, where c is based on determining which values of relative risks are practically significant for decision making based on the size of the population at risk for contrast with other subjective and objective prior distributions.

With the potential for reanalysis of existing studies due to legal challenges of the new standards, model averaging provides a coherent methodology for combining inferences from different models for the same data, as well as judging how well the data support competing models. By explicitly considering model uncertainty in analyses, more realistic measures of uncertainty for relative risks can be obtained, providing a more secure foundation for decision making.

References

- Akaike, H. (1973). “Information theory and an extension of the maximum likelihood”.
In *2nd International Symposium on Information Theory*, Budapest: Akademia Kaido.
267-281.
- Akaike, H. (1978). “A new look at the Bayes procedure”. *Biometrika* 65, 53–9.
- Berger, J.O. (1985). *Statistical Decision Theory and Bayesian Analysis*, Springer, New
York.
- Clyde, M. (1999). “Bayesian model averaging and model search strategies”. In *Bayesian
Statistics 6* J.M. Bernardo, A.P. Dawid, J.O. Berger, and A.F.M. Smith eds. Oxford
University Press. pages 157-185.
- Clyde, M. and DeSimone-Sasinowska, H. 1997. “Accounting for model uncertainty in Pois-
son regression models: Does particulate matter particularly matter?” ISDS Discussion
Paper 97-06.
- Clyde, M. and George, E. (1998). “Flexible empirical Bayes estimation for wavelets”.
Institute of Statistics and Decision Sciences, Duke University, Discussion Paper 98–21.
- Davis, J. Sacks, J. Saltzman, N. Smith, R., and Styer. P. (1996), “Airborne particular
matter and daily mortality in Birmingham, Alabama”. Tech Report 55. National
Institute of Statistical Sciences, Research Triangle Park, NC.
- Draper, D. (1995). “Assessment and propagation of model uncertainty (with Discussion)”.
Journal of the Royal Statistical Society, Series B, 56, 45-98..
- Foster, D.P. and George, E.I (1994). “The risk inflation criterion for multiple regression”.
Annals of Statistics, **22**, 1947–75.
- George, E.I. and McCulloch, R. (1997). “Approaches for Bayesian variable selection”.
Statistica Sinica **7**, 339-374.

- Hodges, J.S. (1987).” Uncertainty, policy analysis, and statistics”. *Statistical Science* 2, 259–291.
- Hoeting, J.A., Madigan, D., Raftery, A.E., and Volinsky, C.T. (1999). “Bayesian model averaging: A tutorial”. To appear in *Statistical Science*.
- Holmes, C.C. and Mallick, B.K. (1997). “Bayesian radial basis functions of unknown dimension”. Tech Report, Imperial College.
- Jeffreys, H. (1961). *Theory of Probability* 3rd Edition, Oxford University Press.
- Kass, R.E. and Raftery, A.E. (1995). “Bayes factors”. *Journal of the American Statistical Association*, 90, 773-795.
- Lamon, E.C. and Clyde, M.(1998). “Accounting for model uncertainty in prediction of chlorophyll *a* in Lake Okeechobee”. ISDS Discussion Paper 98-42
- National Research Council. (1998). “Research priorities for airborne particulate matter”. National Academy Press. Washington, DC.
- Raftery, A.E. (1996).” Approximate Bayes factors and accounting for model uncertainty in generalized linear models”. *Biometrika* 83, 251-266.
- Schwartz, J. (1993), “Air Pollution and daily mortality in Birmingham, Alabama”. *American Journal of Epidemiology*, **137**, 1136–1147.
- Schwarz, G. (1978). “Estimating the dimension of a model”. *Annals of Statistics* 6, 461-464.
- Smith, R.L., Davis, J.M. and Speckman, P. (1997), “Assessing the human health risk of atmospheric particles”. In *Proceedings of the ASA Section on Bayesian Statistical Science*.
- Spiegelhalter, D. J. and Smith, A. F. M. (1982). “ Bayes factors for linear and loglinear models with vague prior information”. *Journal of the Royal Statistical Society, Series B*, 44,377–387

Tierney, L. and Kadane, J.B. (1986). “Accurate approximations for posterior moments and marginal densities”. *Journal of the American Statistical Association*, 81, 82–86.

Viallefont, V. Raftery, A.E. and Richardson, S. (1998). “Variable selection and Bayesian model averaging in case-control studies”. Tech Report 343, Dept. of Statistics, University of Washington.

PM₁₀ current day, one, two, and three day lags.

Daily Monitor 0023 pm0, pm1, pm2, pm3

Area Wide Average pma0, pma1, pma2, pma3

Temperature (current day, one and two day lags)

daily minimum temperature (tmin, tmin1, tmin2)

daily maximum temperature (tmax, tmax1, tmax2)

average daily temperature (mntp, mntp1, mntp2)

Humidity (current day, one and two day lags)

average dew point temperature (dptp, dptp1, dptp2)

daily minimum relative humidity (mnrh, mnrh1, mnrh2)

daily maximum relative humidity (mxrh, mxrh1, mxrh2)

average daily specific humidity (mnsh, mnsh1, mnsh2)

Atmospheric Pressure (current day, one and two day lags)

average daily station pressure (pres, pres1, pres2)

Seasonal Trend

thin-plate spline basis with up to 30 potential knots

Posterior mean of trend **baseline**

Table 1: Explanatory variables used in the Birmingham, AL analysis.

$\log(c_m)$	CIC	Criterion
1	R^2	maximum R^2
2	AIC	Akaike Information Criterion
$\log(n)$	BIC	Bayes Information Criterion
$2\log(p)$	RIC	Risk Inflation Criterion

Table 2: Calibrated Information Criterion priors

SUMMARIES	AIC	BIC
P(Relative Risk = 1 data)	0.03	0.72
Posterior Mean of Relative Risk	1.052	1.015
Posterior Mean of Relative Risk for Models with PM ₁₀	1.054	1.053
Relative Risk for Best Model	1.025	1.00

Table 3: Summaries of the distribution of relative risk associated with a 100 unit increase in PM₁₀ under Bayesian model averaging (BMA) and model selection using the AIC and BIC prior distributions.

Predictive MSE	AIC	BIC
Model Selection	16.83	16.03
Bayesian Model Averaging	16.31	15.98

Table 4: MSE for predicting mortality for the validation set under Bayesian model averaging and model selection using the AIC and BIC prior distributions.

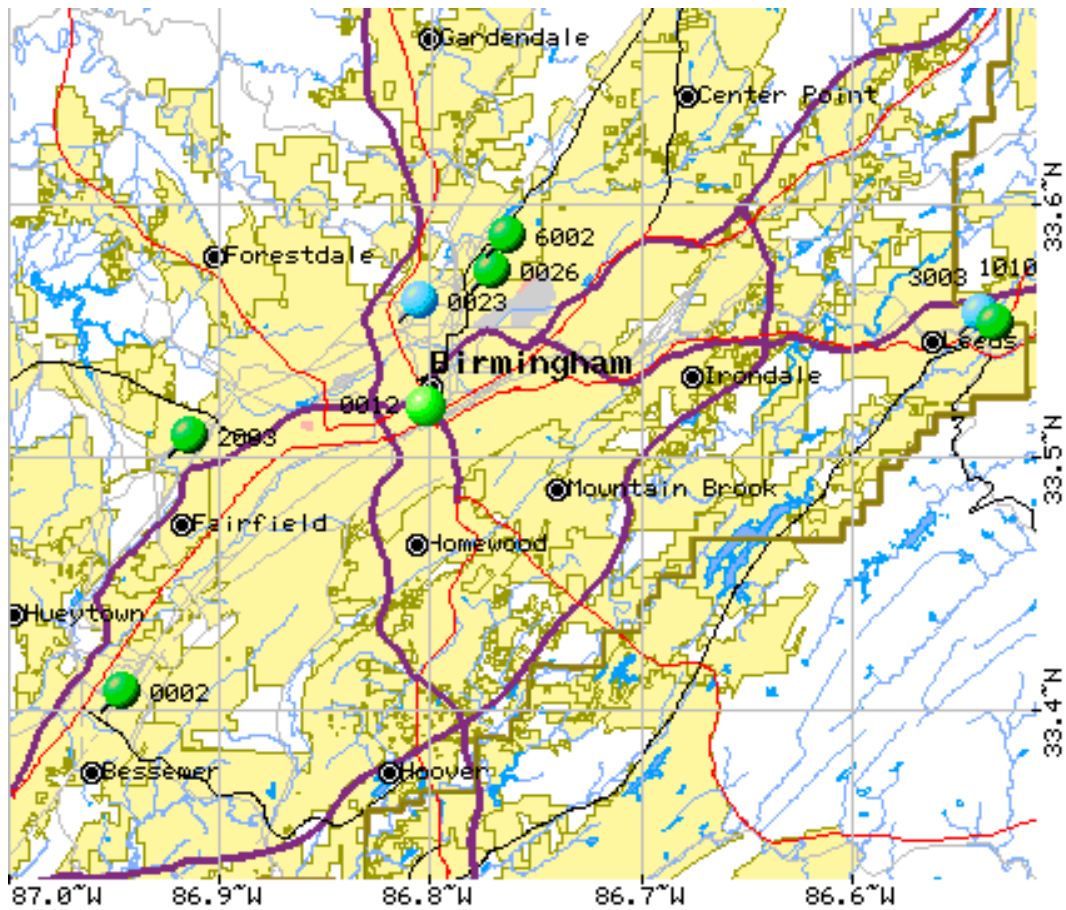


Figure 1: Location of PM₁₀ monitors within the Birmingham metropolitan area (shaded area).

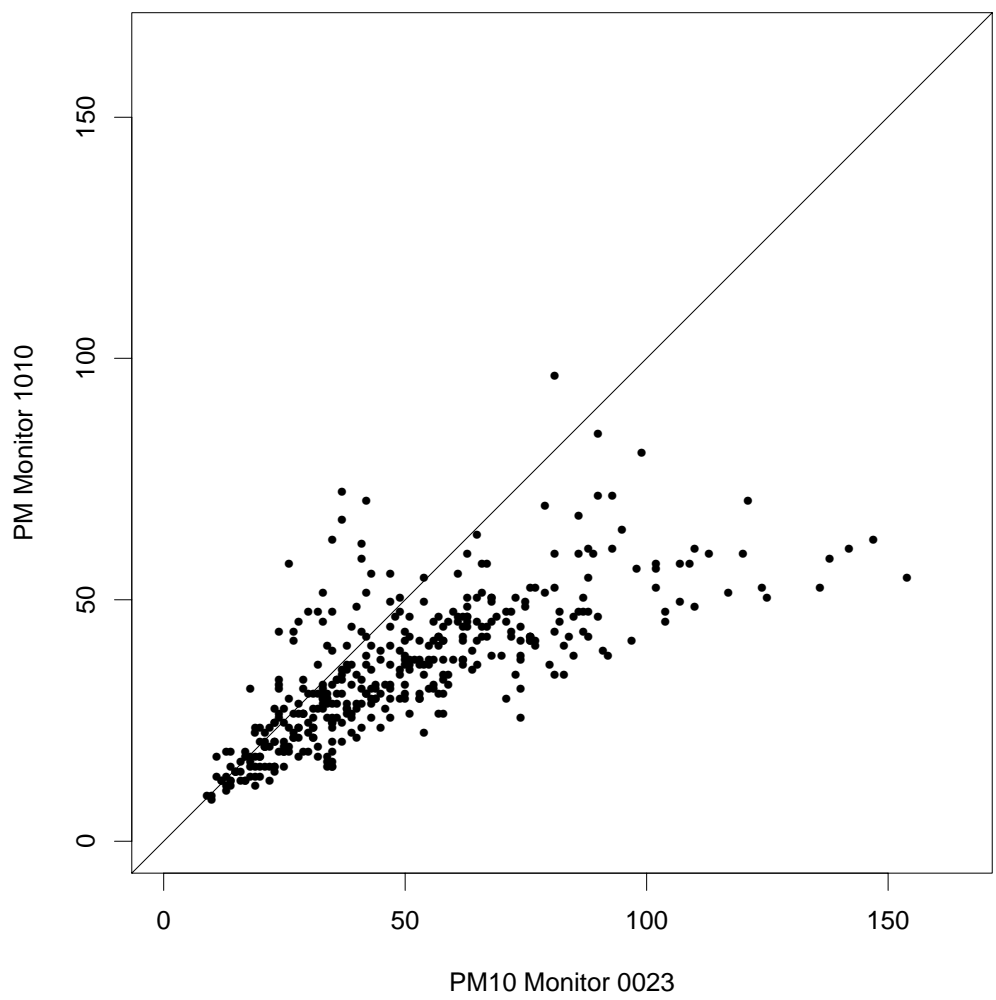


Figure 2: PM_{10} measurements for the daily monitor 1010 in Leeds versus PM_{10} measurements from the daily monitor 0023 in Birmingham.

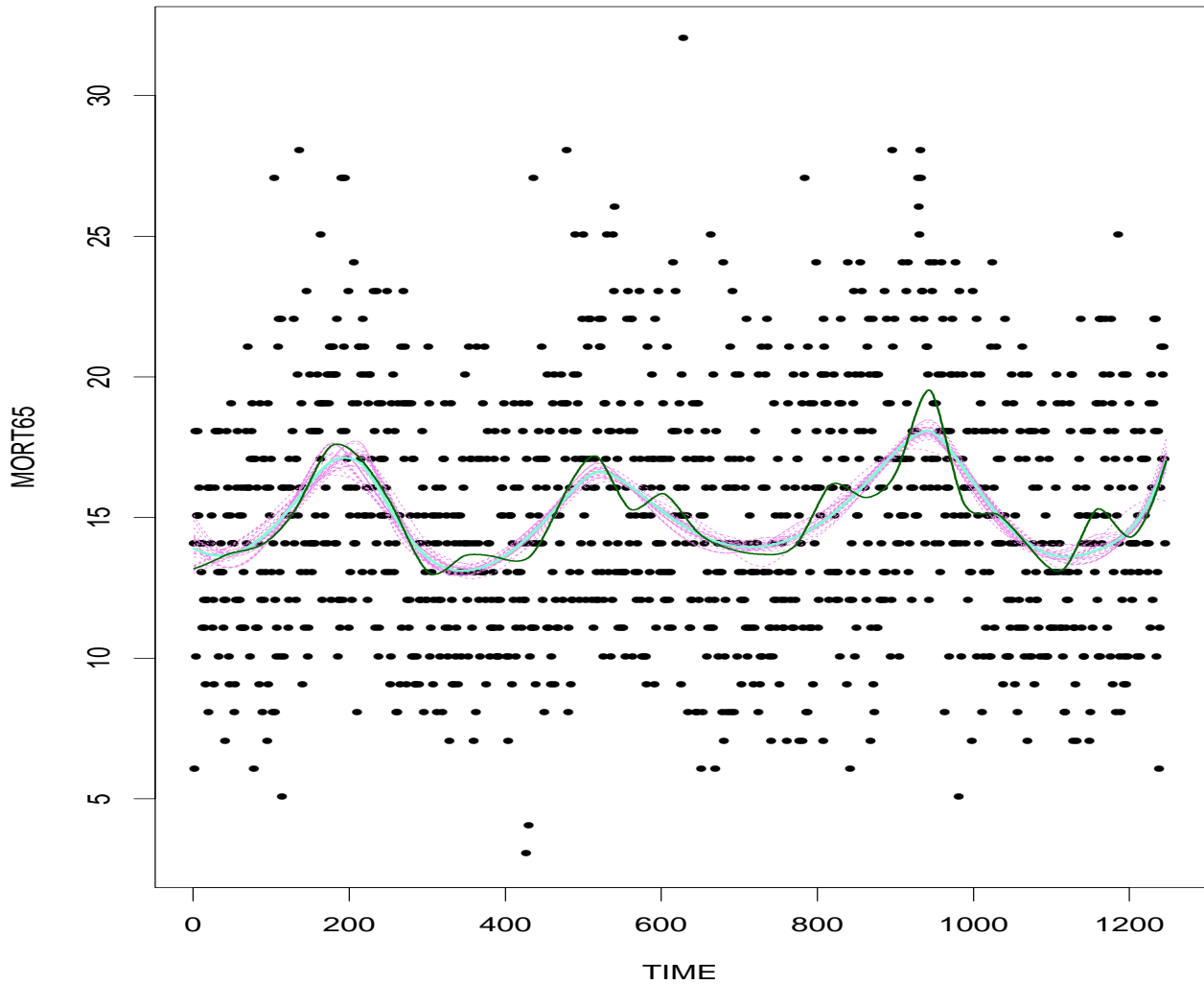


Figure 3: Plot of daily non-accidental mortality in the elderly (over 65) population for Birmingham, Alabama. The thick solid line corresponds to the baseline trend estimate under BMA with the BIC prior distribution; the thin solid line is the GLM estimate under the 30 knot model (roughly one knot for every 40 days); and the light dashed lines correspond to individual estimates from the top 100 models under the BIC prior.

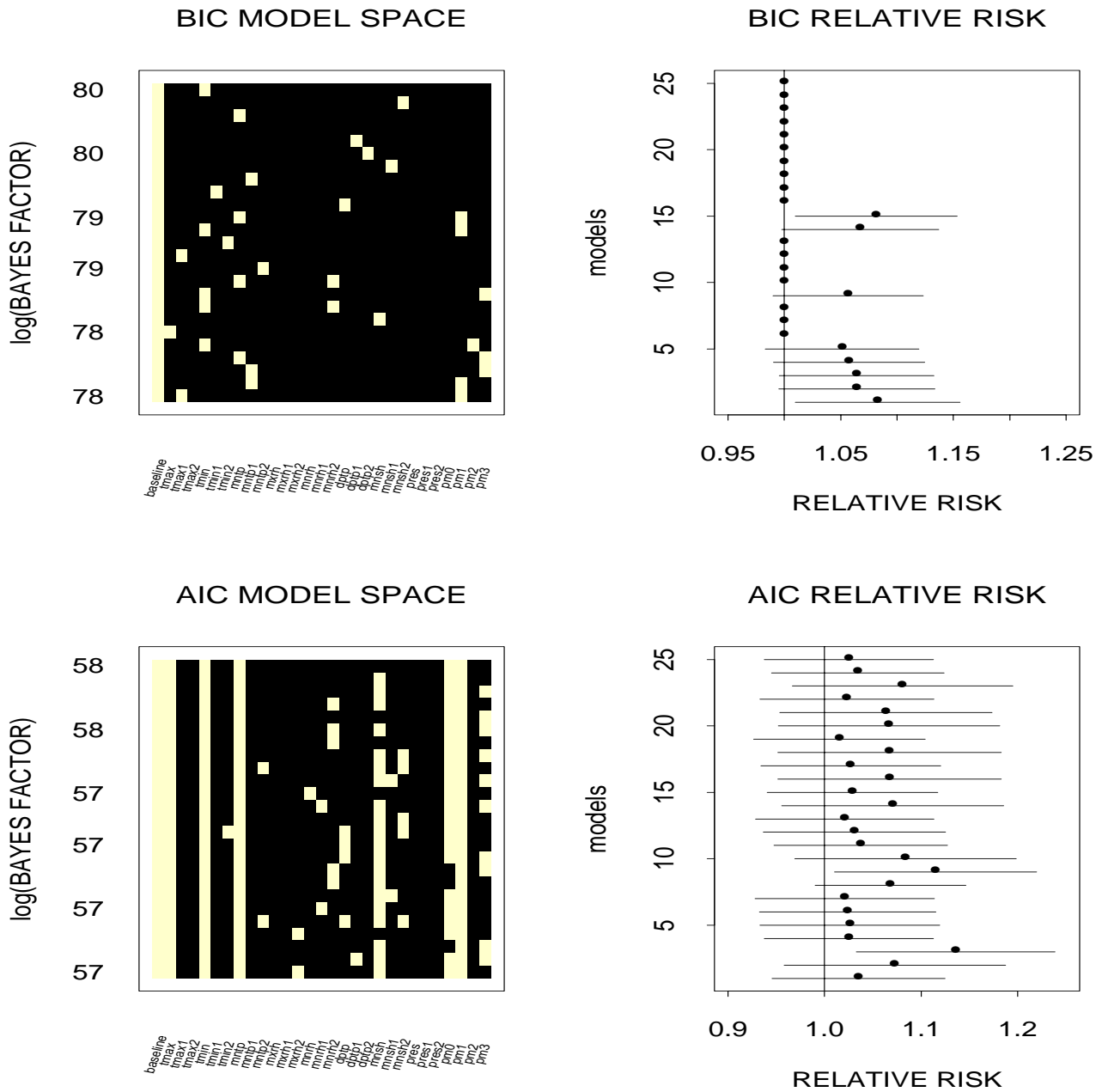


Figure 4: Top 25 models and the corresponding estimates of relative risks with 95% probability intervals under the CIC priors with $c = n$ (BIC) and $c = \exp(2)$ (AIC). Dark squares indicate that the variable in that column is not included in the model for that row.

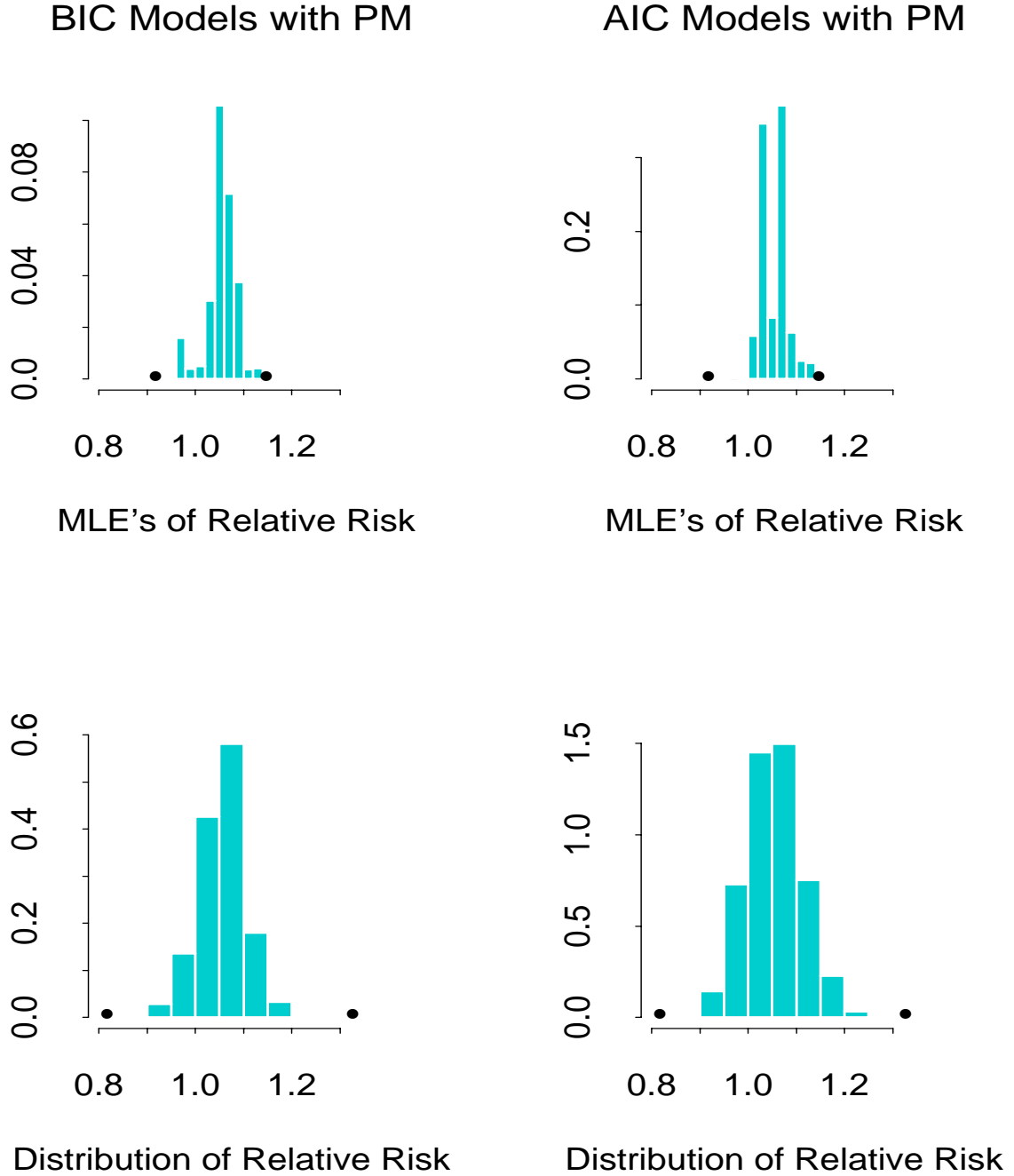


Figure 5: Distribution of relative risk associated with a $100 \mu\text{g}/\text{m}^3$ increase in PM_{10} under the BIC prior (left) and AIC prior (right) for all models that include PM_{10} . Top: histograms are posterior modes (MLE) for each model weighted by respective model probabilities. Lower: histograms of relative risk incorporating both model uncertainty and parameter uncertainty. The points indicate the range of the distribution. While the range and values of relative risk are the same under both posteriors, the weights depend on the prior.