

# Urban Form, Heart Disease, and Geography: A Case Study in Composite Index Formation and Bayesian Spatial Modeling

Gerald Shoultz · Jimmie Givens · J. Wanzer Drane

Published online: 27 September 2007  
© Springer Science+Business Media B.V. 2007

**Abstract** Recent studies indicate a relationship between measures of urban form as applied to urban and suburban areas, and obesity, a risk factor for heart disease. Measures of urban form for exurban and rural areas are considerably scarce; such measures could prove useful in measuring relationships between urban form and both mortality and morbidity in such areas. In modeling area-level mortality, geographic relationships between counties warrant consideration because geographically adjacent areas tend to have more in common than areas farther from each other. We modify county-level indices of urban form found in the literature so that they can be applied to exurban and rural counties. We then use these indices in a Bayesian spatial model that accounts for spatial autocorrelation to determine if there is a relationship between such measures and cardiovascular disease mortality for white males age 35 and older for the time period 1999–2001. Issues related to the formation and usefulness of the indices, and issues related to the spatial model, are discussed. Maps of observed and expected relative risk of mortality are presented.

**Keywords** Bayesian modeling · Mortality · Spatial modeling · Sprawl · Urban form

---

Jimmie Givens retired from his service.

G. Shoultz (✉)  
Department of Statistics, Grand Valley State University, 1133 Makinac Hall, 1 Campus Drive,  
Allendale, MI 49401, USA  
e-mail: shoultzg@gvsu.edu

J. Givens  
National Center for Health Statistics, Hyattsville, USA

J. W. Drane  
Department of Public Health, East Tennessee State University, Johnson City, USA

## Introduction

Over the past few decades there has been considerable increase in the number and percentage of adults and children considered obese and overweight (Flegal et al. 2002; Mokdad et al. 2001). The medical consequences of obesity, including increased risk of coronary heart disease, have been thoroughly documented (Rashid et al. 2003; Surgeon General 2001). Substantial evidence also indicates that lack of physical activity increases the risk of obesity and its consequences, while increased exercise reduces obesity and heart disease risk (Surgeon General 1996).

The purpose of the present study was to examine the relationship between urban form and heart disease. In this paper we use the term “urban form” to denote the organization and structure of buildings, roads, and other physical components. Urban form is not the same as urban sprawl. Urban sprawl denotes negative examples of urban form. For example, Ewing (1997, as cited in Ewing, Schieber et al. 2003, p. 1541) “consider[s] the term ‘sprawl’ to apply to any environment characterized by the following: a population widely dispersed in low-density residential development; rigid separation of homes, shops and workplaces; a lack of distinct, thriving activity centers, such as strong downtowns or suburban town centers; and a network of roads marked by very large block size and poor access from one place to another.”

Ewing, Schmid et al. (2003) in their study of 448 urban and suburban counties found a relationship between a six-variable county sprawl index and body mass index (BMI), minutes walked, and hypertension. Individual covariates in their study included age, education level, and smoking status. Their sprawl index is an abbreviated version of Smart Growth America’s sprawl index of metropolitan areas (Ewing et al. 2002). Data sources included the Behavioral Risk Factor Surveillance System surveys of the Center for Disease Control, U.S. Census data, and the Natural Resources Inventory of the U.S. Department of Agriculture. High BMI is an indicator of obesity, a risk factor for heart and cardiovascular disease (Rashid et al. 2003; Surgeon General 2001).

Applying urban-based concepts of urban form to nonurban areas is problematic. For example, Theobald (2001, pp. 544, 548) found that data aggregated at county and other coarse levels did not capture more finely-grained patterns in land use changes in more rural and exurban areas, and that “relying solely on urban-based land-cover classifications will likely underestimate large areas of low-density human settlement beyond the urban fringe.” He further notes that a poverty of data, clear definitions, and clear land-cover and land-use categories makes delineating land-use changes in exurbia difficult. One result of this is a potential underestimation of the amount of farm land lost to urbanization and development over time. Similarly, relationships between health or other issues and increased development may not be easily seen if measuring devices and techniques for determining land use and other changes in less urban areas do not account for the unique characteristics of such areas.

Geography is essential in understanding both urban form and disease mortality. Longley et al. (2001, p. 99), summarizing Tobler’s (1970) so-called First Law of Geography, said “Everything is related to everything else, but near things are more related than distant things.” Within a geographic area (e.g., a county) measures of

spatial autocorrelation, both alone or with other characteristics, can be important measures of urban form that in turn can illuminate other issues (housing, poverty, etc.).

Spatial autocorrelation between areas also warrants consideration. Most statistical models assume that the response variables (county-level heart disease mortality for our context) are independent. If spatial autocorrelation exists between nearby geographic areas for the response variable, then independence does not hold and standard statistical variables are not valid. Today methods exist to perform both normal and Poisson regression with a spatial autocorrelation component (Lawson et al. 2003; Anselin et al. 2006).

The focus of the present study is more pedagogical and demonstrative, with primary emphasis on index formation and the presentation of spatial modeling issues. Substantive results, while not a primary emphasis of the paper, are briefly presented and discussed. Using only 1990 U.S. Census data, we adapt components of the Smart Growth America index of Ewing et al. (2002) that is designed for suburban and urban counties to create measures of urban form that could be applied to exurban and rural areas. These indices are then used in a Bayesian spatial model to determine if they are predictors of county-level heart disease mortality counts for white males age 35 and over for 1999–2001. In our spatial model areal socioeconomic deprivation, age, lung cancer mortality (used here as a proxy for air pollution), and spatial autocorrelation are all controlled for. The study includes a discussion of issues of model fit and parameter estimation peculiar to spatial modeling.

## Study Population

The study population was comprised of white males age 35 years of age and older who resided in the United States during the 1999–2001 period. About 35 years of age was the cut-off age because the prevalence of cardiovascular disease is considerably lower for younger people (Benjamin et al. 2003, p. 1065) while the risk of heart attack and stroke is markedly higher for those aged 35 and over (Mayo Clinic 2007). The study population has been restricted to white men because of differences between men and women and between people of different races. There are differences between males and females with regard to diagnosis and treatment of heart attacks (AHRQ 2006; Riccotti 2003) as well as risk (NWHIC 2007). African-Americans and other minority groups face a higher risk of heart disease than whites (American Heart Association 2007a) and risk probably differs between nonwhite groups. Future research on nonwhite groups will necessarily include developing models involving counties with small or nonexistent nonwhite populations.

## Data Sets

Counts of deaths due to heart disease, and age-adjusted mortality rates for lung cancer mortality, were obtained using Type II Multiple Cause of Death Files for

1999–2001 (NCHS 1999–2001) from the National Center for Health Statistics and 1990 and 2000 Population SF1 and SF3 Census Files (U.S. Census Bureau 1990, 2002a, b). All deaths, and only such deaths, whose underlying cause as per the death certificate can be classified as I00–I09, I11, I13, or I20–I51 of the International Classification of Diseases, 10th Revision (World Health Organization 1992) are counted as due to heart disease. Measures of mortality, deprivation indices, and 1990 urban form indices (indices explained below) were calculated for 3,137 counties, combinations of counties, or county equivalents in the United States. Counties and their equivalents in parts of Alaska, Virginia, and Montana were combined (see Appendix A).

## Urban Form Measurement

In a study of 101 metropolitan areas, Ewing et al. (2002, p. 10) note that “Sprawl, and its antithesis, compact development, are constructs...they must be *operationalized* to be investigated empirically.” In operationalizing these constructs, the authors created four factors using data from U.S. government and private databases: density, land use mix, degree of centering, and street accessibility. The density factor focused on residential population density and included measures of residential population density for both the entire county and for more urban and rural subsets, and average lot sizes of single-family homes. The land use mix factor focused on diversity of land uses within a county: are homes, jobs, and businesses geographically mixed or separated, and how much? A centering factor measured to what degree residences and businesses cluster in limited areas within a metropolitan area. Variables included percentages of the population within certain distances from the central business district and the “coefficient of variation of population density across census tracts” (2002, p. 23). The street accessibility factor was formed with measures of block size. They reasoned that smaller blocks correspond to more compact development while the larger blocks formed from curvilinear and discontinuous road networks indicated less dense development.

After the variables were sorted into their four factors, principal components analysis (PCA) (Kim and Mueller 1978a, b) was performed separately on the variables comprising each factor. The principal component (PC) used to represent the factor was the one that accounted for the greatest amount of variation among the variables comprising that factor. For example, the component for street accessibility that captured the largest amount of variance among its variables was considered to be a reliable and valid measure of street accessibility. The authors then computed one composite sprawl index from the four individual factors.

## Formation of Urban Form Indices

The present study builds on the sprawl index developed by Ewing et al. (2002), discussed above, which included over 100 metropolitan areas. We wanted to measure every county in the United States. We found, however, that many of Ewing et al.’s

variables were not easily (if at all) available for counties or nonurban areas. Therefore, only 1990 U.S. Census data were used to derive the indices. The variables used by Ewing et al. (2002) focused on urban characteristics, but we wanted both urban and rural characteristics to be included. Therefore their factors were modified so they could be used with both urban and rural areas. Given these restrictions, only measures of density and street accessibility were generated for this study. Formation of land use mix and degree of centering measures is left for later research. Finally, while Ewing et al. used population counts, we use housing counts in our study. Theobald (2001) recommends using housing density as opposed to population density to gauge changes in land use, and our concern with the impact of the built environment points to use of an indicator of built environment.

The density index of Ewing et al. was modified to add measures of exurban and rural counties obtainable with U.S. Census data. As census definitions of urban and rural have changed from 1990 to 2000, and since our eventual dependent variable is 1999–2001 mortality, we used Theobald's (2001) definitions of urban and exurban/rural housing density, listed below.

1. Urban blockgroups have a housing density of at least 1 unit/acre or 640 units/square mile, slightly higher than U.S. Census.
2. Suburban blockgroups have a housing density of between 1 unit/acre and 1 unit/10 acres (64–640 units/square mile).
3. Exurban/rural blockgroups have a housing density of at most 1 unit/10 acres or 64 units/square mile.

Five variables comprise the density index. They are listed below.

1. Overall housing density (units/square mile).
2. Percentage of the area of the county covered by urban blockgroups (100 times the total area of the urban blockgroups divided by the area of the county).
3. Percentage of the area of the county covered by exurban/rural blockgroups (100 times the total area of the exurban and rural blockgroups divided by the area of the county).
4. Density (units/square mile) of urban blockgroups.
5. Density (units/square mile) of exurban/rural blockgroups. If the county has no such blockgroups, then the variable is set at 64 units/square mile, the upper bound for the density of such blockgroups.

These five variables were obtained or derived for each county in the United States. Principal components were extracted on the five variables comprising the density factor using PROC FACTOR (SAS Institute 1999a, pp. 1121–1193) in SAS 8.2 (SAS Institute 1999–2001). The factor that accounted for the greatest amount of variation among the five variables was used to form the index. Factor loadings were obtained for each of the five variables based on the PCA. The density index value for a county is the sum of the five variables for that county weighted by the factor loadings. For later analysis these indices were then standardized to have mean of zero and standard deviation of one.

The density factor accounts for 58% of the variation in our five-variable data set. The reliability coefficient  $\alpha = 0.81$ , indicating strong internal consistency among the five variables. A higher index value indicates a more condensed area and less sprawl.

The road accessibility index modifies Ewing et al.'s streets factor to include larger blocks more likely to be found in exurban and rural areas. The variables comprising our factor are:

1. The percentage of blocks in the county with area less than 0.01 square miles (also used in Ewing et al.; 100 times the count of such blocks divided by the total count of blocks).
2. The percentage of blocks in the county with area more than 1 square mile (intended as a measure of proportion of larger blocks found in exurban and rural areas, it is calculated as 100 times the count of such blocks divided by the total count of blocks).
3. The mean area of all blocks with area less than 1 square mile (intended as mean area of small and medium size blocks, it is calculated by finding the total area of such blocks divided by the total count of such blocks).

These three variables were obtained or derived for each county in the United States. Principal components were extracted on the three variables comprising the road accessibility factor, and indices calculated for each county obtained, in the same manner as for the density factor. The road accessibility factor accounted for 72% of the variation in our three-variable data set. The reliability coefficient,  $\alpha = 0.81$ , indicates strong internal consistency among the three variables. A higher index value indicates a more condensed area and less sprawl.

### Usefulness of Indices

An important question to consider is whether the indices discriminate between urban and exurban/rural counties. We compared our indices against the 1993 Rural–Urban Continuum Codes (RUCC) (USDA 2003). These codes classify counties that are part of a U.S. Census-defined Metropolitan Statistical Area (MSA) by population, and remaining counties by their level of urbanization and whether or not they are adjacent to an MSA. While we do not claim that the RUCC codes are necessarily the best standard of comparison or that our method is the best or only method to use, we do submit that it is a reasonable first attempt.

The density and road accessibility indices were separately compared to the 1993 RUCC codes as follows. The RUCC classify all United States counties into 10 groups based on population and adjacency to a metropolitan county. We condensed their 10 groups into three and classified each county into one of the three groups below.

1. RUCC-1 consists of counties in MSAs with population of at least 250,000 (813 counties). This category covers most urban and suburban areas.

2. RUCC-2 consists of non-MSA counties with urban population of 20,000 or more (244 counties). Some of these are adjacent to an MSA and hence could be smaller suburbs. Many of the counties not adjacent to an MSA may be stand-alone small towns. We use this category as a buffer zone between clearly urban/suburban counties and the exurban/rural counties.
3. RUCC-3 consists of non-MSA counties with urban population of less than 20,000 (2,029 counties).

The density and road accessibility index values were put in descending order. Higher indices imply a more condensed environment and less sprawl. The 813 counties with the highest index values were assigned to Density-1 (analogous to the RUCC-1), the counties with the 244 next highest indices were assigned to Density-2, and the remaining counties were assigned to Density-3. The road accessibility indices were handled the same way. Nine counties with 1993 RUCC codes that did not have index values because of boundary changes since 1990 were removed from this analysis. Crosstabulation tables and inferential statistics for Chi-Square tests of independence were then obtained on the three-level density variable versus the three-level RUCC codes, and on the three-level road accessibility variable versus the three-level RUCC codes.

For the density index 2,351 (76%) of the 3,076 counties were in the same numbered Density Level and RUCC Level (e.g., in Density-1 and RUCC-1, etc.). For these counties both the index and the RUCC group indicated that they were in the same group on the three-group urban-rural continuum. On the other hand, 303 (10%) of the 3,076 counties were at opposite levels of the continuum (e.g., in RUCC-1 and Density-3). These counties were either in an MSA with a low density index or a county with a small or nonexistent metro area with a high density index. The remaining counties differed in their level numbers by one unit (e.g., in Density-1 and RUCC-2).

We found the road accessibility index to be slightly less satisfactory than the density index. For the road accessibility index 2,203 (76%) of the 3,076 counties were in the same numbered Road Accessibility Level and RUCC Level. On the other hand, 455 (15%) of the 3,076 counties were at opposite levels of the three-level continuum. The remaining counties differed in their level numbers by one level.

Based on the comparison with the three-level RUCC continuum we conclude that both the Density and Road Accessibility indices work reasonably well for discriminating between urban and exurban/rural counties. With further study and adjustments they may prove to be useful in making distinctions among exurban and rural counties. The inconsistencies between our indices and the RUCC codes indicate that further study is required to see how our (or their) measures can be adjusted to account for these inconsistencies. So, while we cannot say that our indices are a major improvement in urban form measurement for exurban and rural counties, we do hold that they are a respectable step in that direction.

## Spatial Modeling

A Bayesian spatial model was used to determine whether our indices of urban form are predictors of heart disease mortality for our study population, controlling for certain covariates (listed below). Before proceeding, we state two sets of disclaimers

on what this section will do and what it will not do. First, we present and explain our Bayesian model and its components *only* as applied to our study. We do not argue that a Bayesian hierarchical approach is the best or only proper approach. Second, with regard to spatial autocorrelation we *only* present and apply a model that could be used in a Bayesian context to account for spatial autocorrelation. A statistical test for spatial autocorrelation is presented later in the paper. The model used here, the Intrinsic Conditional Autoregression (ICAR) model, is a specific form of the Conditional Autoregression (CAR) model (Wakefield et al. 2001; Clayton and Kaldor 1987). We do not argue for the ICAR or CAR as the best or only proper approach to model spatial autocorrelation.

### Statement of the Model

The outcome variable is the observed number of deaths due to heart disease given the expected number of deaths and other covariates. It is reasonable to assess observed and expected mortality or incidence counts in an area and then compare the observed and expected counts (Lawson et al. 2003, p. 4). Since the dependent variable is mortality or disease counts, a distribution appropriate for count data should be used (Lawson et al. 2003, p. 38). The Poisson distribution is indeed such a distribution. Therefore, Poisson regression with a log-linear link function was used.

Our model (in matrix form below) is the Poisson version of the generalized linear model:

$$\begin{cases} \mathbf{Y} \sim \text{Poisson}(\lambda) \\ \ln(\lambda) = \boldsymbol{\alpha} + \ln(\mathbf{e}) + \mathbf{X}\boldsymbol{\beta} + \varepsilon \end{cases} \quad (1)$$

where:

- $i$  = county (1–3,137),
- $j$  = age group ( $j = 1, 2, \dots, 6$  for age groups 35–44, 45–54, ..., 85+),
- for each  $(i, j)$  set a unique index  $k = 3137(j - 1) + i, k = 1, 2, \dots, 18822$ ,
- $\mathbf{Y} = [Y_1 \dots Y_{18822}]^T$  is the observed number of deaths of white males age 35+ for county-ages  $k = 1, 2, \dots, 18822$ ; the  $Y_k$ 's are assumed to be independent.
- $\boldsymbol{\alpha} = 18822 \times 1$  matrix with all elements equal to the logarithm of overall relative risk  $\alpha$ ; in our Bayesian context  $\alpha$  has an improper uniform distribution,
- $\ln(\mathbf{e}) = [\ln(e_1) \dots \ln(e_{18822})]^T$  is the matrix of the natural logarithms of the expected number of deaths  $e_k$  of white males age 35+ (defined below) for county-ages  $k$ ,
- $\mathbf{X}$  = covariance matrix with  $k$ th row  $\mathbf{x}_k^T = [x_{k1} \ x_{k2} \dots x_{k6}]$ ,
- $x_{k1} = j - 3.5$  for age group  $j$ , set so that the mean of the  $x_{k1}$ 's = 0,
- $x_{k2} = (j - 3.5)^2 - 35/12 =$  square of age terms, set so that the mean of the  $x_{k2}$ 's = 0,
- $x_{k3} = (j - 3.5)^3$  for age group  $j$ ,
- $x_{k4} = 1990$  Deprivation Index (standardized) for county  $i$ ,
- $x_{k5} = 1990$  Sprawl Index (either Density or Road Accessibility, standardized) for county  $i$ ,



- $x_{k6}$  = County-level Lung Cancer Mortality for white males age 35 and over (standardized) for county  $i$ ,
- $\beta$  = The  $6 \times 1$  matrix of covariate parameters,
- $\varepsilon = [\varepsilon_1 \dots \varepsilon_{18822}]^T$  = the matrix of regression error terms, and
- $\varepsilon_k = u_i + v_k$ , where  $u_i$  = structured (spatial) variation for county  $i$  based on ICAR model (defined below) with overall variance  $\sigma_u^2$ ,  $v_k$  = unstructured (random) variation having normal distribution with mean 0 and variance  $\sigma_v^2$ , and  $u_i$  and  $v_k$  are assumed to be independent of each other. This particular set-up of the error terms is sometimes called the convolution model (Besag et al. 1991).

Parameter estimates were calculated using a Bayesian approach. In the classical (frequentist) statistical approach the parameters are estimated via sample data, while in the Bayesian approach the parameters are treated as random variables. Prior distributions are set on the parameters before data analysis, samples are taken for the random variables and then based on the data the prior distributions (called priors) are adjusted using Bayes' Rule. The adjusted priors are called posterior distributions (see Casella and Berger 2002, p. 324; NIST 2005). Here priors are set on  $\sigma_u^2$ ,  $\sigma_v^2$ ,  $\alpha$ , and the individual  $\beta$ 's, a Poisson distribution is the sampling distribution, and Bayes' rule is used to obtain posterior distributions. See Appendix B for a listing of the priors for our parameters.

Prior to the advent of powerful computing technology statisticians avoided Bayesian methods because these posterior distributions usually cannot be written in a mathematically closed form. Hence, finding characteristics of the distributions (mean, variance, etc.) was difficult or impossible. Today computer simulations based on Markov Chain Monte Carlo (MCMC) methods "allow posterior sampling from models of considerable complexity" (Lawson 2001, p. 239). After an initial "burn-in" phase of some number of iterations the MCMC method converges to the true posterior distribution. One MCMC method, the Gibbs sampler, generates random variables from a marginal distribution without explicitly knowing the distribution. Large samples are simulated from the posterior distribution so that desired characteristics of the distribution can be obtained accurately (Casella and George 1992). A summary of MCMC methods is in Lawson (2001, pp. 239–243).

Much of the procedure for obtaining parameter estimates parallels that of Johnson (2004) and closely parallels that of Shoultz and Givens (in progress). Details are found there and in Appendix C. For our model 10,000 iterations of our simulation were obtained to make sure that parameter estimates converged, and an additional 30,000 iterations to obtain the parameter estimates.

### Calculation of Expected Counts and Relative Risk (Rate Ratio)

The standardized mortality ratio (SMR) compares the actual number of deaths with the number of deaths one would expect from an external reference population. An SMR greater than one indicates that the actual number of deaths is more than what would be expected of the reference population. Here the reference population is the same as the study population, white males age 35+. The SMR for our model (Szklo and Nieto 2000, p. 272) is

$$\text{SMR} = \frac{\text{Observed Number of Deaths}}{\text{Expected Number of Deaths}} = \frac{Y_k}{e_k}, \quad e_k = P_k \bullet \frac{\sum_{k=1}^{18822} Y_k}{\sum_{k=1}^{18822} P_k}, \quad (2)$$

where  $P_k$  is the population of white males for county-age  $k$ . The SMR is an estimate of relative risk for each county-age  $k$  (Lawson et al. 2003, p. 4). Setting  $e_k$  in (1) as defined in (2) leads to estimates of the true relative risk  $\theta_k$  that account for desired covariates (i.e., smoothed SMRs) are obtained. The true relative risk for county-age  $k$  with county  $i$  is

$$\theta_k = \exp\{\alpha + \mathbf{x}_k^T \boldsymbol{\beta} + u_i + v_k\} \quad (3)$$

To obtain smoothed mortality rates instead of smoothed SMR's set  $e_k = P_k$  in (1). The raw mortality rate for index  $k$  is  $Y_k/P_k$ , and the smoothed mortality rate is  $\theta_k$ .

#### Further Definition of Variables

Age as a controlling variable is modeled via a cubic polynomial. Pickle et al. (1996) used a cubic spline with a knot to fit chronic disease mortality. Since our age range (35 and over) is smaller than Pickle et al. the knot is discarded.

Lung Cancer Mortality serves as a proxy for several items in combination, including smoking and the environment. "Cigarette and tobacco smoke...are [two of] the six major independent risk factors for coronary heart disease that you can modify or control...[Cigarette smoke] increases the risk of coronary heart disease by itself. When it acts with other factors, it greatly increases risk...[Cigarette smoke is] the most important risk factor for young men and women" (American Heart Association 2007b). Exposure to air pollution has been found to be related to both Lung Cancer and Cardiopulmonary mortality (Pope et al. 2002). Many epidemiological studies find greater risk of lung cancer among those in urban areas (Nielsen et al. 1996). Finally Lung Cancer Mortality has been previously used as a proxy for smoking in modeling Chronic Obstructive Pulmonary Disease mortality (Nandram et al. 2000).

Greater socioeconomic deprivation is also associated with higher levels of cardiovascular disease mortality (Hayes et al. 2005). County-level socioeconomic deprivation is controlled for with a 16 variable composite deprivation index based on 1990 census data. The index is a modification of Singh (2003) and Singh and Siahpush (2002). The 16 variables included housing, education level, and income data. The index accounted for 54% of the variation in the 16 variables controlling the index. The reliability coefficient  $\alpha = 0.94$  for the standardized county indices, indicating a high level of internal consistency among the variables (Carmines and Zeller 1979). Higher index values indicate higher levels of deprivation. Indices were calculated for 3,137 counties or combinations, with boundaries adjusted as per Appendix A.

Variation is separated into two terms  $u_i$  and  $v_k$ . The SMR  $\hat{\theta}_k = Y_k/e_k$  is the maximum likelihood estimate (MLE) of  $\theta_k$ ; for small areas, the  $\hat{\theta}_k$ s will be unstable because of the sparseness of the data (Wakefield et al. 2001). To handle this problem Wakefield et al. recommend using a multivariate probability distribution for  $\theta = (\theta_{11}, \dots, \theta_{18822})^T$ , where each estimate of  $\theta_k$  is formed by “borrowing strength” from the other estimates. This is done by adding the component  $v_k$  to the log-link function, with  $v_{ij}$  being normally distributed with mean zero and common variance  $\sigma_v^2$ .

**Intrinsic Conditional Autoregression (ICAR) Model for Spatial Autocorrelation**

Numerous methods for modeling spatial and other autocorrelation exist in the literature (e.g., Anselin 1993, 1988; Cressie 1993; Whittle 1954). A concise overview of such methods is found in Bao (no date). The ICAR model (Besag et al. 1991), a special case of the CAR model used to map disease and cancer rates (Wakefield et al. 2001, pp. 116–122; Clayton and Kaldor 1987) is used here. In the CAR approach the distribution of each mortality count  $Y_k$  (or its regression error) is conditioned only upon the remaining  $Y_1, Y_2, \dots, Y_{k-1}, Y_{k+1}, \dots$ . A CAR model can be formed using geographic adjacencies or distances between points. The ICAR estimates the distribution of each mortality count  $Y_k$  (or its regression error) using only its geographic neighbors. The ICAR model is given below. The CAR model, and a proof of the ICAR’s equivalence to it, is found in Appendix C.

The ICAR model for our 3,137 counties is

$$U_i | (U_d = u_d, d \neq i) \sim N \left( \frac{\sum_{d=1}^{3137} u_d w_{id}}{a_i}, \frac{\sigma_u^2}{a_i} \right), \tag{4}$$

where:

- $U_i | (U_d = u_d, d \neq i)$  = the distribution of the spatial autocorrelation term  $u_i$  for county  $i$  conditioned upon the spatial autocorrelation terms for all the other 3,136 counties,  $i, d = 1, 2, \dots, 3, 137$ ,
- $N \left( \frac{\sum_{d=1}^{3137} u_d w_{id}}{a_i}, \frac{\sigma_u^2}{a_i} \right)$  is a normal distribution with mean  $\left( \sum_{d=1}^{3137} u_d w_{id} \right) / a_i$  and variance  $\sigma_u^2 / a_i$ ,
- $w_{id} = \begin{cases} 1 & \text{if } i \text{ and } d \text{ are adjacent, } i \neq d \\ 0 & \text{otherwise} \end{cases}$ ,
- $a_i = \sum_{d=1}^{3137} w_{id}$  is the number of counties adjacent to county  $i$ , and
- $\sigma_u^2$  is a measure of overall variance of the  $u_i$ ’s.

The distribution of each  $u_i$  has mean equal to the average of the  $u_d$ ,  $d \neq i$  that are its neighbors. The overall variance  $\sigma_u^2$  is a random variable with its own prior distribution (see Appendix B).

CAR and ICAR models have at least two benefits. The explicit conditional structure of the CAR and ICAR models is readily modeled with a hierarchical Bayesian approach. The spatial weightings  $w_{id}/a_i$  have a symmetry common to CAR models:  $w_{id}/a_i = w_{di}/a_d$ . Unfortunately, CAR models often have computational and theoretical difficulties. Covariance matrices for the joint distributions of the  $u_i$ 's that do not have inverses lead to joint distributions whose means and/or variances are infinite (Arab et al. 2007; Bannerjee et al. 2004). While methods to handle these problems have been proposed (Bannerjee et al. 2004; Cressie 1993), means to handle such issues are a source of ongoing research (Arab et al. 2007).

## Results

### Parameter Estimates

Posterior parameter estimates for the density and road accessibility indices were obtained using WINBUGS 1.4 (The BUGS Project 2004; Spiegelhalter et al. 2003) and are in Tables 1 and 2 respectively. Signs of the indices for the density and road accessibility indices are consistent with Ewing et al. (2002): Higher levels of sprawl indicate higher risk of heart disease mortality. Credible intervals are the Bayesian equivalent of confidence intervals. A 95% credible interval is formed by the endpoints of the middle 95% of the posterior estimates of the parameters. The data indicate that our indices of urban form are indicators of heart disease mortality.

For the 3,137 counties the relative risk estimates (rate ratios) of mortality due to urban form as measured by our indices were at most 1.02 and 1.09, respectively for the density and road accessibility indices (with a rate ratio of 1 indicating no impact). So while the indices were statistically significant in the model the “real world” impact of urban form as measured here was relatively small.

**Table 1** Parameter estimates with 1990 density index

	Mean	Standard deviation	95% Credible interval
Intercept	0.0264	0.00253	(0.0214, 0.0314)
Age	0.8940	0.00284	(0.8883, 0.8995)
Age <sup>2</sup>	-0.0143	0.00089	(-0.0161, -0.0126)
Age <sup>3</sup>	0.0179	0.00058	(0.0168, 0.0190)
1990 Deprivation Index (Mean 0, SD 1)	0.1086	0.00412	(0.1004, 0.1166)
Lung Cancer Mortality (Mean 0, SD 1)	0.0534	0.00419	(0.0453, 0.0616)
1990 Density Index (Mean 0, SD 1)	-0.0288	0.00401	(-0.0366, -0.0209)

Higher Deprivation Index implies greater socioeconomic deprivation. Higher Density Index indicates more condensed area (less sprawl)

**Table 2** Parameter estimates with 1990 road accessibility index

	Mean	Standard deviation	95% Credible interval
Intercept	0.0288	0.00249	(0.0239, 0.0337)
Age	0.8937	0.00292	(0.8879, 0.8984)
Age <sup>2</sup>	-0.0144	0.00088	(-0.0161, -0.0127)
Age <sup>3</sup>	0.0179	0.00059	(0.0168, 0.0191)
1990 Deprivation Index (Mean 0, SD 1)	0.1030	0.00378	(0.0955, 0.1105)
Lung Cancer Mortality (Mean 0, SD 1)	0.0586	0.00406	(0.0507, 0.0665)
1990 Road Accessibility Index (Mean 0, SD 1)	-0.0180	-0.00330	(-0.0117, -0.0247)

Higher Deprivation Index implies greater socioeconomic deprivation. Higher Road Accessibility Index indicates more condensed area (less sprawl)

### Evidence of Spatial Autocorrelation: Model Comparison

Consider two models  $H_1$  and  $H_0$ , identical in every way except  $H_1$  has a term to account for spatial autocorrelation and  $H_0$  does not. Then  $H_1$  and  $H_0$  are nested models, with  $H_0$  nested within  $H_1$ . To determine if there is evidence of spatial autocorrelation “one could calculate their goodness-of-fit statistics for the two models (call them  $G_1$  and  $G_0$ , respectively), find the difference (or ratio) in their goodness-of-fit statistics  $G_1 - G_0$  (or  $G_1/G_0$ ), and compare them to some sampling distribution of  $G_1 - G_0$  (or  $G_1/G_0$ )” (Dobson 2002, p. 69).

In classical statistics the deviance or some form of the log likelihood is used to compare nested models. Bayesian statisticians tend to use other measures (Congdon 2001, pp. 465–494). The Deviance Information Criterion (DIC) (Spiegelhalter et al. 2002, pp. 583–584, 2003), one such Bayesian measure for model comparison, was developed to handle model comparison when it is not clear how many parameters there are in the model. The DIC essentially adds a factor for the level of complexity of the model. A lower DIC implies improved fit of the model. A more complete explanation of the DIC (including definition and calculation) for our model is in Appendix F.

The DIC was used to determine whether there is evidence of spatial autocorrelation in and of itself, without controlling for any other variable. The log-link functions are  $F_1 = \ln(\mu_{ij}) = \ln(e_{ij}) + u_i + v_{ij}$  (spatial autocorrelation included and no covariates) and  $F_0 = \ln(\mu_{ij}) = \ln(e_{ij}) + v_{ij}$  (spatial autocorrelation not included and no covariates) respectively, with variables defined as in (1). For  $F_1$  and  $F_0$  the DIC's were 117,477 and 117,519, respectively. The addition of the spatial autocorrelation component drops the DIC 42 units—a statistically significant amount in the eyes of many (David Spiegelhalter personal correspondence; Holsinger 2006), but as it is a less than 0.1% drop one could reasonably conclude that the impact of spatial autocorrelation on risk of mortality due to heart disease is small.

### Chloropleth Maps

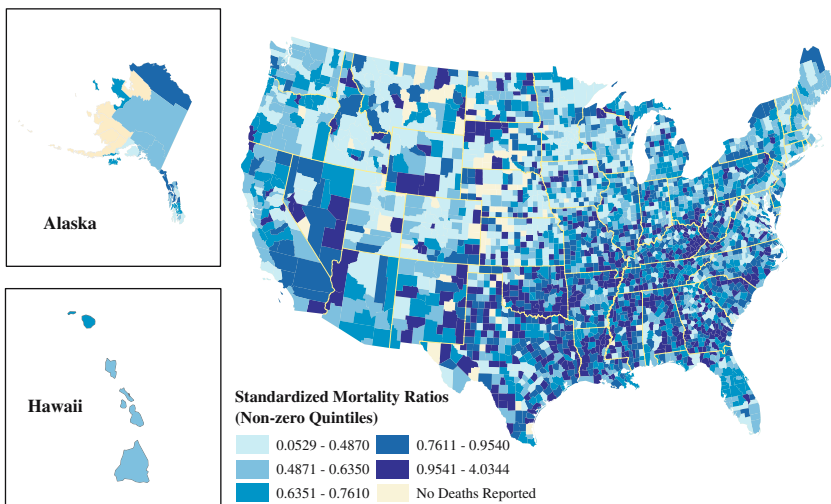
Maps of observed standardized mortality ratios (also found in Shoultz and Givens in progress) and smoothed relative risk estimates for those 55–64 years of age are

found in Figs. 1 and 2 at the end of the paper. Both maps clearly indicate higher relative risk estimates of heart disease mortality found in the southeastern part of the country. The smoothing impact of the spatial model is clearly seen in the second map.

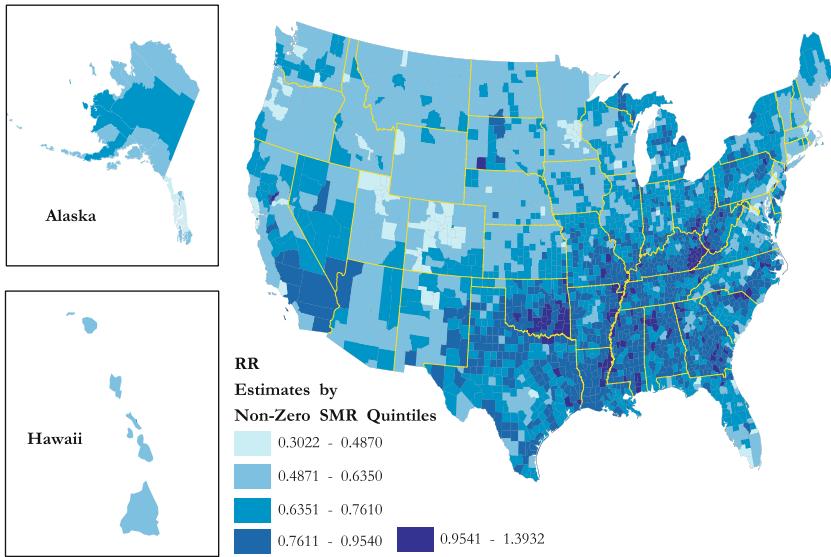
A related issue is the relative risk (RR) of heart disease mortality due to geographically based variables while controlling for remaining nongeographic variables. Rejection of the null hypothesis  $H_0 : RR = 1$  in favor of the alternative  $H_1 : RR > 1$  implies that overall the geographic variables for that county indicate risk of mortality due to heart disease for white males age 35 and over living in the county even after controlling for age and unstructured variation. In a Bayesian context one can take posterior estimates of relative risk and use those posterior estimates to obtain empirical estimates of the probability of that alternative hypothesis. A map of these estimated probabilities with non-geographic variables (age and unstructured variation) excluded is found in Fig. 3. The darker areas have higher empirical probabilities that  $RR > 1$ . For these areas the probability that the population relative risk due to mortality in these geographic areas is higher than one. The darker areas are concentrated in the southeast, southern California and parts of the northeastern United States.

## Discussion

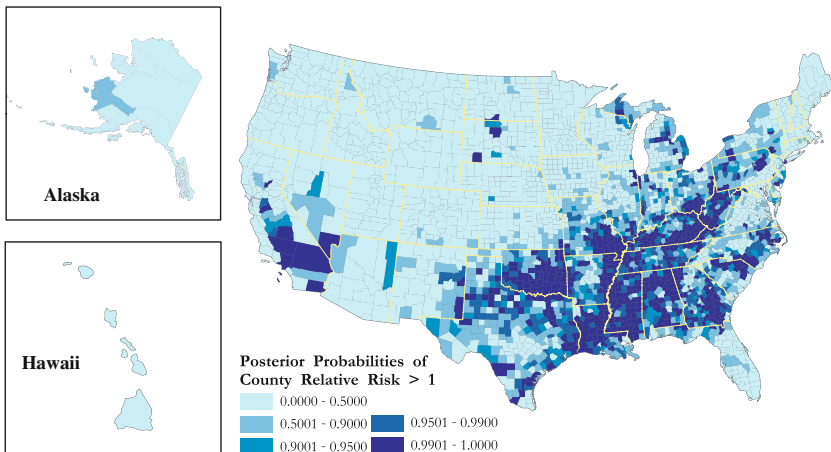
In this paper we modified two county-level measures of urban form found in the literature for use with exurban and rural areas. We then used the indices in a Bayesian hierarchical model that accounted for spatial autocorrelation to determine whether these measures were predictors of heart disease mortality. As substantive



**Fig. 1** Standardized mortality ratios (raw estimates of relative risk), age 55–64. *Note:* SMR = observed count of deaths/expected counts of deaths



**Fig. 2** Overall estimates of relative risk (covariates: age, lung cancer mortality, 1990 density index, 1990 socioeconomic deprivation index, spatial autocorrelation and unstructured variation, age 55–64. *Note:* Mean of 30,000 Posterior Estimates of  $\theta_{ij} = \exp \left\{ \begin{aligned} &\alpha + \beta_{age}Age_j + \beta_{age^2} \left[ (Age_j)^2 - 35/12 \right] + \beta_{age^3} (Age_j)^3 + \\ &\beta_{Dep90}Dep90_i + \beta_{Sprawl90}Sprawl90_i + \beta_{LCM}LCM_i + u_i + v_{ij} \end{aligned} \right\}$



**Fig. 3** Posterior probabilities of county-level relative risk being greater than one, controlling for age, estimates of relative risk (covariates: lung cancer mortality, 1990 density index, 1990 socioeconomic deprivation index, structured variation. *Note:* Empirical estimates of  $P(\beta_{Dep90}Dep90_i + \beta_{Sprawl90}Sprawl90_i + \beta_{LCM}LCM_i + u_i > 1)$  based on 30,000 iterations

results are not a primary emphasis of the paper we briefly summarize them here before discussing the index formation and spatial modeling issues that form the primary emphasis of this paper.

We found the density and road accessibility indices to be predictors of heart disease mortality in the expected direction: less condensed counties have higher rates of heart disease mortality. The small rate ratios of our indices (1.09 at most, with rate ratio of 1 implying no effect) imply that the substantive impact of urban form (as measured) on heart disease mortality was small. These results are consistent with those of Ewing, Schmid et al. (2003). Spatial autocorrelation was found to be a statistically significant predictor of heart disease mortality but the substantive impact on such mortality is small. Maps of raw and smoothed estimates of heart disease mortality indicate that the southeastern and far western areas of the country are of higher risk of heart disease mortality even after accounting for measures of urban form, lung cancer mortality, spatial autocorrelation, and socioeconomic deprivation.

Serious problems exist with the urban form indices. One reviewer found the lack of substantive importance of our results to be a sign that (1) the indices themselves are questionable measures of urban form and (2) a different spatial scale other than the county level is needed. Both may well be the case. The two indices cover related yet separate aspects of urban form. Further research could include combinations of the variables forming the indices, as per Ewing, Schmid et al. (2003). Generally, sound urban form measurement may require a smaller land scale than the county level, and any county measure will likely lose the potentially rich variation of land uses within a county.

Data limitations are also a problem with our indices. Desiring to develop a scale usable for all U.S. counties and not just urban and suburban ones, we restricted ourselves to U.S. census data. Ewing et al. (2002) used other data sources in addition to the U.S. census and examined only areas where such data was available. Our modifications of Ewing et al.'s scales to account for urban characteristics gave indices with questionable usefulness. More detailed data would potentially improve the accuracy and usefulness of our indices. The lack of detailed land-cover and land-use data, and other pertinent data sets, for smaller communities makes urban form measurement for more rural areas difficult (Theobald 2001). We conclude that studies similar to this one are better done with smaller geographic areas where more precise data in addition to census data could be obtained. This is especially important for analysis of exurban and rural areas where data are likely to be considerably sparse in comparison to cities.

Our testing of the indices for whether or not they made distinctions between urban and rural areas is another concern. More extreme inconsistencies between the RUCC designations and our indices occurred in 10–15% of urban and exurban/rural counties. The RUCC codes used urban population within a county and living in (or county-level adjacency to) a metro area to set its scale. For our density scale we used measures of housing density, and for our road accessibility scale we used block sizes and block areas. While all of these are related in some manner they measure very different characteristics. These differences likely account for some of the inconsistencies. Further research could include examining counties with such inconsistencies to better understand the reasons for inconsistencies. Such information could result in more useful and accurate indices that account for these inconsistencies.



A Bayesian hierarchical version of Poisson regression was used here. With diseases having relatively rare mortality, a Zero-Inflated Model or Negative Binomial (Lee et al. 2002; Durham et al. 2004) may be more appropriate. The availability of such software as WinBUGS that capably handles the computational issues that come with Bayesian methodology allows modeling possibilities that were previously difficult or impossible.

Spatial autocorrelation was controlled for in the model. Given the geographic nature of both mortality and of urban form we argue that one should continue to consider some form of spatial autocorrelation in models of mortality and morbidity. Other models exist to measure spatial autocorrelation. The Simultaneous Autoregression (SAR) model (Whittle 1954) is popular in the econometrics literature. While the CAR approach models the distribution of each  $Y_k$  (or its regression error) conditional only upon the remaining  $Y_1, Y_2, \dots, Y_{k-1}, Y_{k+1}, \dots$  the SAR approach models the distribution of the regression errors (or the  $Y_k$ 's) simultaneously. A presentation of a SAR model, and a comparison of the CAR and the SAR models, is in Appendix F. Other software packages, especially GeoDa (Anselin et al. 1998–2004), allow modeling to account for spatial autocorrelation by other means (e.g., Spatial Lag and Spatial Error models).

By mapping measures of relative risk and related posterior probabilities we saw geographic patterns of mortality. Such maps have interpretation issues. Moulton et al. (1994, p. 297) observe that “[while] mapping the actual rates of standardized mortality ratios (SMRs) can be informative, such an approach can be misleading when the denominators vary across geographic units. Those regions with the smallest populations will tend to have both the highest and lowest rates merely because they have the greatest variability.” Therefore, Bayesian methods, in which a proposed prior distribution of rates is combined with the observed rates to obtain posterior rates, are applied. These posterior rates are called stabilized rates because the variability of the original estimated rates is reduced. In examining mortality rates for cancer Gelman and Price (1999) found a pattern similar to Moulton et al. when mapping both posterior estimates and measures of statistical significance. Explanatory variables, according to Gelman and Price, will not abolish such artifacts but can sometimes mitigate them by reducing uncertainty in the parameters.

**Acknowledgments** *Disclaimer* This publication was made possible through a fellowship sponsored by the Center for Disease Control (CDC), National Center for Health Statistics (NCHS) and the Association of Schools of Public Health (ASPH). The findings and conclusions contained in this paper represent the views of the authors. No official support or endorsement by either the Grand Valley State University Department of Statistics or the Centers for Disease Control and Prevention, Department of Health and Human Services is intended, nor should be inferred.

## Appendix A

This study used 1990 Census data but 1999–2001 mortality data. Therefore changes in county boundaries 1990 and 2000 (U.S. Census Bureau 2001–2005) are accounted for as follows: In Montana the portion of Yellowstone National Park inside of Montana was divided between Gallatin and Park Counties; therefore,

Gallatin and Park Counties are pooled for this study. In Virginia South Boston town (a county equivalent) was merged into Halifax County; for this study South Boston and Halifax County are pooled. In Alaska Denali borough was formed from parts of Yukon-Koyukak and Southeast Fairbanks boroughs; for this study the areas for Yukon-Koyukak, Southeast Fairbanks, and Denali boroughs are merged. Also in Alaska, Skagway-Yakutat-Angoon borough was divided into Skagway-Hoonan-Angoon borough and Yakutat Borough; for this study the boroughs are merged. Annexations of portions of a county into another county that did not dissolve a county were not accounted for.

## Appendix B

For our model the priors are:  $1/\sigma_u^2, 1/\sigma_v^2 \sim \Gamma(0.5, 1/0.0005)$ , and all  $\beta'_s \sim N(0, 1/0.00001)$ , where  $\Gamma(\alpha, \varepsilon)$  is a Gamma function with shape parameters  $\alpha$  and  $\varepsilon$  and  $N(\mu, \sigma^2)$  is a Normal distribution with mean  $\mu$  and variance  $\sigma^2$ .

## Appendix C

For each sprawl index WINBUGS 1.4 (The BUGS Project 2004; Spiegelhalter et al. 2003) was used for running three independent Markov Chains. First, for each of our two urban form indices (density and road accessibility), model (1) without spatial autocorrelation and unstructured variation components (e.g., without the terms  $u_i$  and  $v_k$ ) was run in PROC GENMOD in SAS (SAS Institute 1999b, pp. 1365–1464) to obtain maximum likelihood estimates of the coefficients. Initial values for the three chains were those estimates plus 4, 0 and  $-4$  standard deviations. Initial values for  $1/\sigma_u^2$  and  $1/\sigma_v^2$  were taken as 0.001, 1,000 (mean of the prior distribution) and 7,000, respectively. Initial values for  $u_i$  and  $v_k$  were all set at 0.

Time series Gelman–Rubin diagnostic graphs and trace graphs were used to check convergence of relative risk estimates and parameter estimates for all three chains (Spiegelhalter et al. 2003). After 10,000 iterations the traces of the parameter estimates had good mixing around a common value, with varying degrees of white noise around those values. Gelman–Rubin graphs were convergent and stable. Chains were examined for  $1/\sigma_u^2, 1/\sigma_v^2, \sigma_u$  and  $\sigma_v$ ; these chains also mixed well. Hence, after 10,000 iterations convergence of all parameters was concluded. An additional 30,000 iterations were then taken for each chain, but to reduce autocorrelation every third iteration was kept for later calculations. As a result, a total of 30,000 iterations (10,000 from each chain) were used to obtain parameter estimates.

Convergence of parameter estimates was checked according to the “Checking convergence” section of the WINBUGS 1.4 (Spiegelhalter et al. 2003) manual. The chains clearly appeared to be overlapping one another, and parameter estimates look stable. Posterior graphs had the desired bell shape.

The Brooks and Gelman (1998) version of the Gelman–Rubin convergence statistic (GR) as given in WINBUGS 1.4 was used for iterations 5,001–10,000. Let

**Table C1** Summary of Gelman–Rubin statistics for convergence

Variable	Iteration	Middle 80% of estimates				GR = P/W
		Raw		Normalized <sup>a</sup>		
		P <sup>b</sup>	W <sup>c</sup>	P <sup>b</sup>	W <sup>c</sup>	
Density	9,851	0.00849	0.00844	0.9318	0.9264	1.006
	9,901	0.00846	0.00842	0.9291	0.9243	1.005
	9,951	0.00845	0.00841	0.9283	0.9233	1.005
	10,001	0.00850	0.00844	0.9332	0.9269	1.007
Road accessibility	9,851	0.00785	0.00784	0.9524	0.9510	1.001
	9,901	0.00784	0.00782	0.9508	0.9489	1.002
	9,951	0.00784	0.00782	0.9510	0.9498	1.001
	10,001	0.00784	0.00783	0.9511	0.9502	1.001

<sup>a</sup> Reset to have maximum value of 1

<sup>b</sup> Parameter Estimates for all three chains pooled together and width determined from the pooled set

<sup>c</sup> Widths of the three chains determined separately and the three widths then averaged

X be the width of the middle 80% of the parameter estimates of all three chains pooled together, and let Y be the average of the widths of the middle 80% of the parameter estimates for each of the three chains individually. Then GR = X/Y. Here, we want GR to converge close to 1 and both X and Y to converge to some number. Numerical values for X, Y and GR for later iterations are found in Table C1. The X and Y columns are numerically close to each other and the GR values are nearly equal to one.

Finally, we calculated parameter estimates for a variety of prior distributions to determine whether said estimates were sensitive to the choice of the prior. We found the parameter estimates to be very similar for all our choices.

### Appendix D

This discussion of CAR and ICAR follows much of Wakefield et al. (2001, p. 110ff). We first define the CAR model. Define  $N_n(\mathbf{0}_n, \sigma^2\Sigma)$  as an n-dimensional normal distribution with  $n \times n$  positive definite (i.e., the matrix has an inverse) correlation matrix  $\Sigma$  and parameter  $\sigma^2$ , and let  $\mathbf{Q} = \Sigma^{-1}$  have elements  $Q_{id}$ ,  $i, d = 1, \dots, n$ . The general CAR model can be written as (Besag and Kooperburg 1995)

$$U_i | (U_d = u_d, d \neq i) \sim N \left( \sum_{d=1}^n M_{id} u_d, \sigma_u^2 V_{ii} \right), \tag{5}$$

where

- $M_{id} = \begin{cases} -Q_{id}/Q_{ii} & \text{if } i \neq d \\ 0 & \text{if } i = d \end{cases}$
- $\sigma_u^2$  is a measure of overall variance of the  $u_i$ 's and
- $V_{ii} = 1/Q_{ii}$ .

For a complete derivation of the above relationships see Wakefield et al. (2001, pp. 124–125). Since  $\mathbf{Q}$  is symmetric  $M_{id}V_{dd} = M_{di}V_{ii}$ . In matrix form the correlation matrix  $\Sigma$  is:

$$\Sigma = \mathbf{Q}^{-1} = \mathbf{V}^{-1}(\mathbf{I} - \mathbf{M}), \tag{6}$$

where:

- $\mathbf{V}$  is a matrix with elements  $V_{ii}, i = 1, \dots, n$  and 0 otherwise, and
- $\mathbf{M}$  is the matrix of spatial weights  $M_{id}$ .

If the matrix  $\mathbf{Q}$  has an inverse than the matrix  $\mathbf{U} = [u_1, u_2, \dots, u_n]^T$  has distribution  $\mathbf{U} \sim N_n(0_n, \sigma_u^2(\mathbf{I} - \mathbf{M})^{-1}\mathbf{V})$ .

To obtain the ICAR model (4) from the general CAR model of (5) set  $V_{ii} = 1/a_i$  and  $M_{id} = w_{id}/a_i$ . Here  $\mathbf{Q}$  does not have an inverse. Proof: Each row  $i$  of the matrix  $\mathbf{I} - \mathbf{M}$  has a solitary 1 on the diagonal,  $a_i$  elements with value  $-1/a_i$ , and the remaining elements equal to zero. Since the sum of the elements on each row all equal zero the matrix  $\mathbf{Q}$  has rank  $n - 1 < n$ , and so is not full rank and is not invertible.

### Appendix E

We first define the DIC and then derive the DIC for the model  $F_1 = \ln(\mu_k) = \ln(e_k) + u_i + v_k$ , where  $i$  and  $k$  are defined as in (1). Recall that after 10,000 “burn-in” iterations obtain convergence we obtained 30,000 more iterations for the parameter estimates. The Log Likelihood  $LL$  is

$$LL = LL(\mathbf{Y}, \{\hat{\mathbf{Y}}|\varpi\}) = \sum (Y_k \log \hat{Y}_k - \hat{Y}_k - Y_k!), \tag{7}$$

where  $\mathbf{Y} = \{Y_k\}$  are the observed counts and  $\{\hat{\mathbf{Y}}|\varpi\} = \{\hat{Y}_k\}$  is the set of predicted counts derived from the set of parameter estimates  $\varpi$  (Dobson 2002, p. 76). Using (7), the DIC (Spiegelhalter et al. 2002, 2003) is:

$$DIC = \overline{-2LL(\mathbf{Y}, \{\hat{\mathbf{Y}}|\varpi\})} - (-2LL(\mathbf{Y}, \{\hat{\mathbf{Y}}|\varpi\})), \tag{8}$$

where:

- $\overline{-2LL(\mathbf{Y}, \{\hat{\mathbf{Y}}|\varpi\})}$  is the mean of the 30,000 individual iterations of  $-2LL(\mathbf{Y}, \{\hat{\mathbf{Y}}|\varpi\})$  and
- $\bar{\varpi}$  is the posterior mean of the parameter estimates from the 30,000 iterations.

To calculate the DIC for model  $F_1$ , follow these three steps. First, calculate:

$$\overline{-2LL(\mathbf{Y}, \{\hat{\mathbf{Y}}|\varpi\})} = -2 \sum_{l=1}^{30000} \sum_{k=1}^{18822} \left( (Y_k \log \hat{Y}_k^{(l)} - \hat{Y}_k^{(l)} - Y_k!) \right) / 30000, \tag{9}$$

where

- $l =$  iteration number ( $l = 1, 2, \dots, 30,000$ . taken after the 10,000 iterations for convergence),
- $\hat{Y}_k^{(l)}$  = the expected number of deaths predicted via the  $l$ th iteration from model  $F_1 : \hat{Y}_k^{(l)} = e_k \bullet \exp(\hat{u}_i^{(l)} + \hat{v}_k^{(l)})$ ,
- $\hat{u}_i^{(l)}$  = the spatial autocorrelation component predicted via the  $l$ th iteration,
- $\hat{v}_k^{(l)}$  = the unstructured variation component predicted via the  $l$ th iteration,
- $\varpi = \{\hat{u}_i^{(l)}, \hat{v}_k^{(l)}\}$ , the set of parameter estimates for iteration  $k$ , and
- $\{\hat{Y}|\varpi\}$  = the set of individual  $\hat{Y}_k^{(l)}$ s.

Second, calculate

$$-2LL(Y, \{\hat{Y}|\bar{\varpi}\}) = -2 \sum_{k=1}^{18822} (Y_k \log \hat{Y}_k^\bullet - \hat{Y}_k^\bullet - Y_k!), \tag{10}$$

where

- $\hat{Y}_k^\bullet$  = the estimated number of deaths using  $\bar{\varpi} : \hat{Y}_k^\bullet = e_k \bullet \exp(\bar{u}_i + \bar{v}_k)$ ,
- $\bar{u}_i$  = the mean of the  $\hat{u}_i^{(l)}$ s,  $\bar{u}_i = \sum_{l=1}^{30000} \hat{u}_i^{(l)} / 30000$ ,
- $\bar{v}_k$  = the mean of the  $\hat{v}_k^{(l)}$ s,  $\bar{v}_k = \sum_{l=1}^{30000} \hat{v}_k^{(l)} / 30000$ ,
- $\bar{\varpi} = \{\bar{u}_i, \bar{v}_k\}$ , the set of means of the 30,000 parameter estimates, and
- $\{\hat{Y}|\bar{\varpi}\}$  = the set of individual  $\hat{Y}_k^\bullet$ s.

Finally, substitute the results of Eqs. 9 and 10 into Eq. 8 to find the DIC.

### Appendix F

SAR models are similar to the AR(1) model in time series. Note that in our model (1) we set  $\varepsilon_k = u_i + v_k$  with structured variation  $u_i$  independent of unstructured variation  $v_k$ . As this is not possible with the SAR model we default to the error term  $\varepsilon_k$ .

For our conditions (1) we can write a SAR model as:

$$\begin{cases} \mathbf{Y} \sim \text{Poisson}(\lambda) \\ \ln(\lambda) = \boldsymbol{\alpha} + \ln(\mathbf{e}) + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \\ \boldsymbol{\varepsilon} = \rho\mathbf{S}\boldsymbol{\varepsilon} + \boldsymbol{\eta} \end{cases} \tag{11}$$

where

- $\mathbf{Y}$  is the matrix of observed mortality counts as per (1),
- $\boldsymbol{\lambda} = [\lambda_1 \dots \lambda_{18822}]^T$ ,  $\lambda_k = e_k \bullet \exp(\alpha + \mathbf{x}_k^T \boldsymbol{\beta})$  for each  $k$  with remaining variables defined as per (1), and  $\alpha$  is a constant,
- $\boldsymbol{\varepsilon} = [\varepsilon_1 \dots \varepsilon_{18822}]^T =$  regression error terms,
- $\rho$  is a measure of spatial correlation ( $-1 \leq \rho \leq 1$ ),
- $\mathbf{S}$  is an  $18,822 \times 18,822$  neighborhood (spatial weighting) matrix with zeroes on the diagonals (*not necessarily symmetric*), standardized so the row sums add to one, and
- $\boldsymbol{\eta} = [\eta_1 \dots \eta_{18822}]^T$ ,  $\eta_k \sim N(0, \sigma^2)$ .

It follows from (10) that  $\boldsymbol{\varepsilon} = (\mathbf{I} - \rho \mathbf{S})^{-1} \boldsymbol{\eta}$  and so the covariance matrix  $\Sigma$  for the SAR model is  $\Sigma = \sigma^2 (\mathbf{I} - \rho \mathbf{S})^{-1} (\mathbf{I} - \rho \mathbf{S}^T)^{-1}$ .

If desired SAR can also use nearest neighbor weighting: Set  $s_{id} = w_{id}/a_i$  with  $w_{id}$ ,  $a_i$  defined as in (4)

Briefly comparing and contrasting the CAR and SAR models (Arab et al. 2007; Cressie 1993; Whittle 1954; Bao no date):

- CAR sets specifications for the  $Y_k$ s conditionally, while SAR does so simultaneously.
- Spatial weighting matrices do not have to be symmetric in the SAR model but do in the CAR model.
- A SAR model can always be restated as a CAR model, but not vice versa.
- The SAR model and the CAR model are the same if and only if the covariance matrices are the same.
- The CAR model is more computationally efficient than the SAR model because the matrix  $\mathbf{I} - \mathbf{M}$  in the CAR model is symmetric while the matrix  $\mathbf{I} - \rho \mathbf{S}$  in the SAR model is not.
- In some cases the spatial weights in the SAR model may not be identifiable.
- Parameter estimates for the SAR model are statistically not consistent. That is, for increasing sample sizes the parameter estimates may not converge with high probability to the actual parameter.
- CAR gives the best (that is, the minimum mean squared prediction error) estimates of  $Y_k$  based on all the other  $Y_l$ s,  $l \neq k$ .
- If it makes more sense to specify the model conditionally, or if there is a symmetric structure in the correlation matrix, use a CAR model.

## References

- AHRQ (2006). *Research on cardiovascular disease in women*. Program Brief. AHRQ Publication No. 06-P016. Rockville: Agency for Healthcare Research and Quality. Retrieved March 8, 2007 from <http://www.ahrq.gov/research/womheart.htm>
- American Heart Association (2007a). *Risk factors and coronary heart disease*. Retrieved March 6, 2007 from <http://www.americanheart.org/presenter.jhtml?identifier=4726>
- American Heart Association (2007b). *Cigarette smoking and cardiovascular disease*. Retrieved March 9, 2007 from <http://www.americanheart.org/presenter.jhtml?identifier=454>
- Anselin, L. (1988). The maximum likelihood approach to spatial process models. In L. Anselin (Ed.), *Spatial econometrics: Methods and models* (Chapt. 6, pp. 57–80). Dordrecht: Kluwer Academic.

- Anselin, L. (1993). Discrete space autoregressive models. In M. Goodchild, B. Parks, & T. Steyaert (Eds.), *Environmental modeling with GIS* (pp. 454–469). Oxford: Oxford University Press.
- Anselin, L., & The Regents of the University of Illinois (1998–2004). *GeoDa 0.95-i*. Available at the University of Illinois GeoDa web site: <https://www.geoda.uiuc.edu/default.php>
- Anselin, L., Syabri, I., & Kho, Y. (2006). GeoDa: An introduction to spatial data analysis. *Geographical Analysis*, 38, 5–22.
- Arab, A., Hooten, M. B., & Wikle, C. K. (2007). Hierarchical spatial models. In *Encyclopedia of geographical information science*. New York: Springer. Retrieved March 9, 2007 from Utah State University Department of Mathematics and Statistics web site: <http://www.math.usu.edu/~hooten/papers/HSMv4.pdf> (in press).
- Banerjee, S., Carlin, B. P., & Gelfand, A. E. (2004). *Hierarchical modeling and analysis for spatial data*. Boca Raton: Chapman and Hall/CRC Press.
- Bao, S. (n.d.). *An overview of spatial econometric models*. Retrieved June 23, 2005 from China Data Center, University of Michigan website: [http://www.umich.edu/~inet/chinadata/docs/topic\\_3.pdf](http://www.umich.edu/~inet/chinadata/docs/topic_3.pdf)
- Benjamin, S. M., Geiss, L. S., Pan, L., Engelgau, M. M., & Greenlund, K. J. (2003). Self-reported heart disease and stroke among adults with and without diabetes—United States, 1999–2001. *Morbidity and Mortality Weekly Report*, 52(44), 1065–1070. Retrieved March 6, 2007 from <http://www.cdc.gov/mmwr/preview/mmwrhtml/mm5244a2.htm>
- Besag, J. E., & Kooperburg, C. (1995). On conditional and intrinsic autoregressions. *Biometrika*, 82, 733–746.
- Besag, J., York, J., & Mollie, A. (1991). Bayesian image restoration with two applications in spatial statistics. *Annals of the Institute of Statistical Mathematics*, 43, 1–59.
- Brooks, S. P., & Gelman, A. (1998). Alternative methods for monitoring convergence of iterative simulations. *Journal of Computational and Graphical Statistics*, 7, 434–455.
- Carmines, E. G., & Zeller, R. A. (1979). *Reliability and validity assessment: Sage University paper series on quantitative applications in the social sciences* (Vol. 7). London: Sage Publications.
- Casella, G., & Berger, R. L. (2002). *Statistical inference* (2nd ed.). Pacific Grove: Duxbury Publishing.
- Casella, G., & George, E. I. (1992). Explaining the Gibbs sampler. *The American Statistician*, 46(3), 167–174.
- Clayton, D. G., & Kaldor, J. (1987). Empirical Bays estimates of age-standardized relative risks for use in disease mapping. *Biometrics*, 43, 671–691.
- Congdon, P. (2001). *Applied Bayesian modeling*. Chichester: John Wiley and Sons.
- Cressie, N. (1993). *Statistics for spatial data* (revised ed.). New York: John Wiley and Sons.
- Dobson, A. J. (2002). *An introduction to generalized linear models* (2nd ed.). Boca Raton: Chapman and Hall.
- Durham, C. A., Pardoe, I., & Vega, E. (2004). A methodology for evaluating how product characteristics impact choice in retail settings with many zero observations: An application to restaurant wine purchase. *Journal of Agricultural and Resource Economics*, 29(1), 112–131.
- Ewing, R. (1997). Is Los Angeles style sprawl desirable? *Journal of the American Planning Association*, 63(1), 107–126.
- Ewing, R., Pendall, R., & Chen, D. (2002). *Measuring sprawl and its impact* (Vol. 1). Retrieved March 9, 2007 from the Smart Growth America website: <http://www.smartgrowthamerica.org/sprawindex/MeasuringSprawlTechnical.pdf>
- Ewing, R., Schieber, R. A., & Zegeer, C. V. (2003). Urban sprawl as a risk factor in motor vehicle occupant and pedestrian fatalities. *American Journal of Public Health*, 93, 1541–1545.
- Ewing, R., Schmid, S., Killingsworth, R., Zlot, A., & Raudenbush, S. (2003). Relationship between urban sprawl and physical activity, obesity, and morbidity. *American Journal of Health Promotion*, 18(1), 47–57.
- Flegal, K. M., Carroll, M. D., Ogden, C. L., & Johnson, C. L. (2002). Prevalence and trends in obesity among US adults, 1999–2000. *Journal of the American Medical Association*, 288, 1723–1727.
- Gelman, A., & Price, P. (1999). All maps of parameter estimates are misleading. *Statistics in Medicine*, 18, 3221–3234.
- Hayes, D. K., Greenlund, K. J., Denny, C. H., Croft, J. B., & Keenan, N. L. (2005). Racial/ethnic and socioeconomic disparities in multiple factors for heart disease and stroke—United States, 2003. *Morbidity and Mortality Weekly Report*, 54(05), 113–117. Retrieved March 9, 2007 from: <http://www.cdc.gov/mmwr/preview/mmwrhtml/mm5405a1.htm>

- Holsinger, K. (2006). The deviance information criterion. Retrieved March 9, 2007 from University of Connecticut Department of Ecology and Evolutionary Biology faculty member Kent Holsinger's website: <http://www.darwin.eeb.uconn.edu/eeb348/lecture-notes/testing-hardy-weinberg/node5.html>
- Johnson, G. D. (2004). Smoothing small area maps of prostate cancer incidence in New York state using fully Bayesian hierarchical modeling. *International Journal of Health Geographics*, 3, 29. Retrieved March 9, 2007 from: <http://www.ij-healthgeographics.com/content/3/1/29>
- Kim, J., & Mueller, C. (1978a). *Factor analysis: Statistical methods and practical issues*. Sage University Paper Series on Quantitative Applications in the Social Sciences, 14. Beverly Hills: Sage Publications.
- Kim, J., & Mueller, C. (1978b). *Introduction to factor analysis: What it is and how to do it*. Sage University Paper Series on Quantitative Applications in the Social Sciences, 14. Beverly Hills: Sage Publications.
- Lawson, A. (2001). *Statistical methods in spatial epidemiology*. Chichester: John Wiley and Sons.
- Lawson, A., Browne, W., & Rodeiro, C. (2003). *Disease mapping with WinBUGS and MLwiN*. Chichester: John Wiley and Sons.
- Lee, A. H., Stevenson, M. R., Wang, K., & Yau, K. K. W. (2002). Modeling young driver motor vehicle crashes: Data with extra zeroes. *Accident Analysis and Prevention*, 34, 515–521.
- Longley, P. A., Goodchild, M. F., Maguire, D. J., & Rhind, D. W. (2001). *Geographic information systems and science*. Chichester: John Wiley and Sons.
- Mayo Clinic (2007). *Heart disease prevention: 5 strategies keep your heart healthy*. January 15. Retrieved March 6, 2007 from Mayo Clinic website: <http://www.mayoclinic.com/health/heart-disease-prevention/WO0004>
- Mokdad, A. H., Bowman, B. A., Ford, E. S., Vinicor, F., Marks, J. S., & Koplan, J. P. (2001). The continuing epidemics of obesity and diabetes in the United States. *Journal of the American Medical Association*, 286, 1195–1200.
- Moulton, L. H., Foxman, B., Wolfe, R. A., & Port, F. K. (1994). Potential pitfalls in interpreting maps of stabilized rates. *Epidemiology*, 5(3), 297–301.
- Nandram, B., Sedransk, J., & Pickle, L. W. (2000). Bayesian analysis and mapping of mortality rates for chronic obstructive pulmonary disease. *Journal of the American Statistical Association*, 95, 1110–1118.
- NCHS (1999–2001). *Type II multiple cause of death files 1999–2001*. Hyattsville: National Center for Health Statistics.
- NWHIC (2007). *Heart disease*. Retrieved March 6, 2007 from National Women's Health Information Center website: <http://www.4woman.gov/faq/heartdis.htm>
- Nielsen, P. S., Okkels, H., Sigsgaard, T., Kyrtopoulos, S., & Autrup, H. (1996). Exposure to urban and rural air pollution: DNA and protein adducts and effect of glutathione-S-transferase genotype on adduct levels. *International Archives of Occupational and Environmental Health*, 68(3), 170–176.
- NIST (2005). 8.1.10: How can Bayesian methodology be used for reliability evaluation? In *NIST/SEMATECH e-Handbook of statistical methods*. Retrieved March 7, 2007 from the National Institute of Standards and Technology website: <http://www.itl.nist.gov/div898/handbook/apr/section1/apr1a.htm#What%20is%20Bayesian%20Methodology%20and%20why%20is%20it>
- Pickle, L. W., Mungiole, M., Jones, G. K., & White, A. A. (1996). *Atlas of United States mortality*. Hyattsville: U.S. Department of Health and Human Services.
- Pope, C. A., Burnett, R. T., Thun, M. J., Calle, E. E., Krewski, D., Ito, K., & Thurston, G. D. (2002). Lung cancer, cardiopulmonary mortality, and long-term exposure to fine particulate air pollution. *Journal of the American Medical Association*, 287(9), 1132–1141.
- Rashid, M. N., Fuentes, F., Touchon, R. C., & Wehner, P. S. (2003). Obesity and the risk for cardiovascular disease. *Preventive Cardiology*, 6(1), 42–47.
- Riccotti, H. (2003). *Heart disease: Differences between men and women*. Retrieved March 6, 2007 from Beth Israel Deaconess Medical Center website: [http://www.bidmc.harvard.edu/display.asp?node\\_id=4952](http://www.bidmc.harvard.edu/display.asp?node_id=4952)
- SAS Institute (1999a). *SAS/STAT user's guide, version 8* (Vol. 1). Cary: SAS Institute.
- SAS Institute (1999b). *SAS/STAT user's guide, version 8* (Vol. 2). Cary: SAS Institute.
- SAS Institute (1999–2001). *SAS 8.2*. Cary: SAS Institute.
- Shoultz, G., & Givens, J. (in progress). Heart disease, change in economic deprivation over time, and Bayesian spatial modeling: A case study.



- Singh, G. K., & Siahpush, M. (2002). Increasing inequalities in all-cause and cardiovascular mortality among US adults aged 25–64 years by area socioeconomic status, 1969–1998. *International Journal of Epidemiology*, *31*, 600–613.
- Singh, G. K. (2003). Area deprivation and widening inequalities in U.S. mortality, 1969–1998. *American Journal of Public Health*, *93*(7), 1137–1143.
- Spiegelhalter, D. J., Best, N., Carlin, B. P., & Van der Linde, A. (2002). Bayesian deviance, the effective number of parameters and the comparison of arbitrarily complex models. *Journal of the Royal Statistical Society B*, *64*, 583–640.
- Spiegelhalter, D. J., Thomas, A., Best, N., & Lunn, D. (2003). *WinBUGS user manual, version 1.4, January 2003*. Retrieved March 7, 2007 from the WinBUGS website: <http://www.mrc-bsu.cam.ac.uk/bugs/winbugs/manual14.pdf>
- Surgeon General (1996). *Physical activity and health: A report of the Surgeon General*. Atlanta: U.S. Department of Health and Human Services, Centers for Disease Control and Prevention, National Center for Chronic Disease Prevention and Health Promotion. Retrieved June 24, 2005 from <http://www.cdc.gov/nccdphp/sgr/pdf/sgrfull.pdf>
- Surgeon General (2001). *The Surgeon General's call to action to prevent and decrease overweight and obesity*. Rockville: U.S. Department of Health and Human Services, Public Health Service Office of the Surgeon General. Retrieved June 24, 2005 from <http://www.surgeongeneral.gov/topics/obesity/calltoaction/CalltoAction.pdf>
- Szklo, M., & Nieto, F. J. (2000). *Epidemiology: Beyond the basics*. Sudbury: Jones and Bartlett.
- The BUGS Project (2004). *WINBUGS 1.4*. Updated version 1.4.3 retrieved September 14, 2007 at: <http://www.mrc-bsu.cam.ac.uk/bugs/winbugs/contents.shtml>
- Theobald, D. M. (2001). Land-use dynamics beyond the American urban fringe. *Geographical Review*, *91*(3), 544–564.
- Tobler, W. R. (1970). A computer movie simulating urban growth in the Detroit region. *Economic Geography*, *46*, 234–240.
- U.S. Census Bureau (1990). *1990 census of population and housing block statistics, CD90-1B-(1-10)*. Washington: U.S. Census Bureau.
- U.S. Census Bureau. (2002a). *Census 2000 summary file 1*. Washington: U.S. Census Bureau. Available at <http://www.census.gov>
- U.S. Census Bureau. (2002b). *Census 2000 summary file 3, CD-ROM CS-D00-S3ST-08-US1*. Washington: U.S. Census Bureau.
- U.S. Census Bureau (2001–2005). *Geographic changes for Census 2000 + Glossary*. Retrieved March 7, 2007 from: <http://www.census.gov/geo/www/tiger/glossary.html#states>
- U.S. Department of Agriculture [USDA] (2003). *Measuring rurality: Rural-urban continuum codes*. Retrieved April 7, 2007 from the Economic Research Service, United States Department of Agriculture, at: <http://www.ers.usda.gov/Briefing/Rurality/RuralUrbCon/>
- Wakefield, J. C., Best, N. G., & Waller, L. (2001). Bayesian approaches to disease mapping. In P. Elliott, J. Wakefield, N. Best, & D. Briggs (Eds.), *Spatial epidemiology, methods and applications* (pp. 104–127). Oxford: Oxford University Press.
- Whittle, P. (1954). On stationary process in the plane. *Biometrika*, *41*, 434–449.
- World Health Organization (1992). *International statistical classification of diseases and related health problems, 10th revision*. Geneva: World Health Organization.