



Effect of Sample Size on the Performance of Ordinary Least Squares and Geographically Weighted Regression

Mitra L. Devkota^{1*}, Gary Hatfield¹ and Rajesh Chintala²

¹Department of Mathematics and Statistics, South Dakota State University, Brookings, SD 57007, USA.

²Department of Plant Science, South Dakota State University, Brookings, SD 57006, USA.

Research Article

Received: 20 July 2013

Accepted: 20 September 2013

Published: 05 October 2013

Abstract

A recently developed spatial analytical tool, Geographically Weighted Regression (GWR) was used to deal with spatial nonstationarity in modeling the crop residue yield potential for North Central region of the USA. Average of daily mean temperature and total precipitation of crop growing season were the explanatory variables. In this study, the model performance of Ordinary Least Squares (OLS) and GWR were compared in terms of coefficient of determination (R^2) and corrected Akaike Information Criterion (AICc). Moran's I and Geary's C were used to test the spatial autocorrelation of OLS and GWR residuals. The explanatory power of the models was assessed by approximate likelihood ratio test. Furthermore, the test of spatial heterogeneity of the GWR parameters was conducted by using data sets with small and large samples. The comparative study of R^2 and AICc between the models showed that all the GWR models performed better than the analogous OLS models. Test of spatial autocorrelation of residuals revealed that the OLS residuals had higher degrees of spatial autocorrelation than the GWR residuals indicating that GWR mitigated the spatial autocorrelation of residuals. Results of the approximate likelihood ratio test showed that GWR models performed better than the OLS models suggesting that the OLS relationship was not constant across the space of interest. More importantly, it was demonstrated that the data set would have to be large enough for the individual parameters of GWR models to be spatially heterogeneous.

Keywords: Geographically weighted regression, spatial autocorrelation, spatial heterogeneity, residuals, crop residue yield potential.

1 Introduction

A traditional modeling technique used in geographical analysis is based on Ordinary Least Squares (OLS). In this technique, it is assumed that patterns in the data are spatially constant, and

*Corresponding author: mitra.devkota@sdstate.edu;

therefore, the parameter estimates are the same for the whole study area. This technique does not account for location in the analysis of the relationship between the variables. Therefore, such parameter estimates can be considered as global statistics [1]. A simple form of the OLS model (in matrix form) is given by

$$Y = X\beta + \varepsilon \tag{1}$$

where Y is the $n \times 1$ vector of the dependent variable, X is the design matrix of independent (explanatory) variables, which includes a column of 1s for the intercept, β is the vector of regression coefficients (global parameter estimates for the independent variables) and ε is a random vector, where $\varepsilon \sim N(0, \sigma^2 I)$. The maximum likelihood estimate of β under such model is given by

$$\hat{\beta} = (X^T X)^{-1} X^T Y \tag{2}$$

However, in practice, there may be situations when the pattern of the data is not spatially constant. This is referred to spatial nonstationarity [2]. Under spatial nonstationarity, the measurement of a relationship depends partly on the location where the measurement is taken. Such a relationship which is nonstationary over space will not be well addressed by OLS regression. The violation of the assumptions of independence of residuals in spatial data creates additional problems, such as spatial autocorrelation. Spatial autocorrelation is a phenomenon in which the value of a given variable at a location is related to the values of the same variable in a nearby location [3].

In an attempt to explore and explain the spatially varying relationships by allowing the model parameters to vary over space and thus to attempt to overcome some restrictive assumptions of OLS regression, the concept of Geographically Weighted Regression (GWR) was developed [1]. GWR is an extension of OLS regression in which the parameters are allowed to vary spatially. As such, spatial relationships between the variables can be examined and patterns can be identified. GWR is often referred to be a local model because it provides estimates to local statistics and hence it is more appropriate when the relationships vary spatially [1]. GWR model is defined [2] as

$$Y_i = \sum_{j=1} X_{ij} \beta_j(p_i) + \varepsilon_i \tag{3}$$

where Y is again a $n \times 1$ vector of dependent variables, p_i is the coordinates for the observation i , and ε_i is the random error term for the i th observation. These $\beta_j(p_i)$ contain the local parameters to be estimated. The concept here is that these expressions would be substituted into the original model and the resulting expanded model (which may now be nonlinear) is calibrated [2].

Previous studies have shown that GWR performs better than OLS regression model [4, 5]. GWR is not recommended in situations with small sample sizes ($n \approx 160$ in their experiments) [6]. To the best of our knowledge, no work has been done in the literature to test the spatial heterogeneity of individual parameters of GWR model when applied to data of small and large samples. Therefore,

the primary purpose of this research work was to test the spatial heterogeneity of individual parameters of GWR model when applied to spatially varying data of small and large samples. Similarly, the secondary purpose of this research work was to test the performance of GWR and OLS regression models when applied to these data sets. In addition to comparing the performances of these two models, our third objective was to assess the influence of temperature and precipitation on crop residue yield potential and to assess the degree of spatial autocorrelation of OLS and GWR residuals using the outputs generated.

In this study, the spatial dependence of GWR and OLS residuals were compared. The coefficient of determination (Adjusted R^2) and the Akaike Information Criterion (AIC and AICc) were used to compare the models. The widely accepted rule of thumb [7] that a decrease of 3 in AIC considered as an improvement was also used to evaluate the performance of the models. In addition to comparing the models with different independent variables, AIC_C can also be used to compare OLS and GWR models (as long as the response variable is the same in the models) and also in the software that determines the optimal value for the bandwidth; the bandwidth with lowest AIC_C value will be used to estimate the parameters [7]. An approximate likelihood ratio test based on the F test was used to check for significant improvement of GWR over OLS. A test of spatial heterogeneity of individual parameters was also used.

In Section 2, the weighting functions and kernel bandwidths for GWR were discussed. Section 3 had a comparative study of GWR and OLS models. In addition, the test of spatial heterogeneity of individual parameters of GWR and the outputs from GWR model were discussed. In section 4, OLS and GWR regressions were used to model the crop residue yield potential data of South Dakota and the ten states of the North-Central region (Illinois, Indiana, Iowa, Minnesota, Montana, Nebraska, North Dakota, South Dakota, Wisconsin, and Wyoming) of the USA for the years 1970, 1980, and 2008 and evaluated their corresponding model performances. Summary of this study was given in Section 5.

2 Weighting Functions and Kernel Bandwidths

The basic methodology behind GWR is that a weighted distance decay function w_{ij} is used for model calibration [1, 2]. One of the methods for such calibration is to take a point p_i and to fit a weighted OLS regression with observations within the circle weighted as 1 and zero otherwise. That is, the weight w_{ij} for a given p_i assigned to observation j would be

$$w_{ij} = \begin{cases} 1 & \text{if } d_{ij} < r \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

where d_{ij} is the distance between locations of observations i and j and r is the radius of the circle.

A variety of weighting functions are available [2]. In all of these, it is assumed that the weight decreases with an increase in distance between the observations. One of the weighting functions could be an exponential function given by

$$w_{ij} = \exp\left(-\frac{d_{ij}}{r}\right) \tag{5}$$

where d_{ij} is, as before, the distance between locations of observations i and j , and r is the kernel bandwidth.

Another weighting function commonly used in practice is to calculate the weight by using a Gaussian distance decay function (also referred to as the Continuous Weighting Function) given by

$$w_{ij} = \exp\left(-\frac{d_{ij}^2}{2r^2}\right) \tag{6}$$

where the values of w_{ij} decay gradually with distance to the extent that when $d_{ij} = 3r$, the weight virtually decreases to zero [8]. Results from GWR are not very sensitive to the choice of the weighting function, provided that it is smooth and follows a distance decay property [2]. In this paper, the weights w_{ij} are calculated at each calibration location using the Gaussian distance decay function given by (6).

The kernel bandwidth parameter r is first estimated from the data to fit a GWR model. At present, there are three methods for this: direct assignment of the bandwidth of number of nearest neighbors [9], cross validation process (an iterative process that finds the kernel bandwidth that minimizes the prediction error of all the response variables using a subset of data for prediction) [10, 11], and a corrected Akaike Information Criterion (AIC_C) [1]. As of now, the most commonly used approach is cross validation, and this method was used in our analysis. The vector of estimated regression coefficients is given by

$$\hat{\beta}_i = (X^T W_i X)^{-1} X^T W_i Y \tag{7}$$

where W_i is a square matrix of weights relative to the position i and is given by

$$W_i = \begin{bmatrix} w_1 & 0 & 0 & 0 \\ 0 & w_2 & 0 & 0 \\ 0 & 0 & \dots & \dots \\ 0 & 0 & \dots & w_N \end{bmatrix}$$

and contains the geographical weights in its main diagonal and 0 in its off diagonal elements, N is the number of observations, $X^T W_i X$ is the geographically weighted variance covariance matrix [7], and Y is the vector of values of the response variables.

At present, GWR can be run using GWR3.0 or R. In both of these, fixed and adaptive bandwidth kernels can be chosen. The former computes a bandwidth which is held constant over space while the latter adapts the bandwidth distance according to the density of the data (the denser the data, the smaller the bandwidth and vice versa) [1, 5, 12]. In this study, an adaptive kernel bandwidth is chosen since sample densities varied over the study area. The data analyses were done using the *spgwr* package [13] for GWR in the statistical software package R, version 3.0.1 [14].

3 Model Comparisons

The comparison between OLS and GWR models was performed by comparing their R^2 and AIC/AIC_c using a small data set (crop residue data of 66 counties of South Dakota, U.S. for the years 1970, 1980, and 2008) to those with a larger data set (crop residue yield potential data of 743 counties of ten states of the North-Central region (Illinois, Indiana, Iowa, Minnesota, Montana, Nebraska, North Dakota, South Dakota, Wisconsin, and Wyoming) of the USA for the years 1970, 1980, and 2008). The conventional interpretation of R^2 is that a larger value has greater explanatory power. While comparing two models using AIC_c, the model with lower AIC_c value is considered to be a better model. More precisely, if the AIC_c values for two models differ by at least 3, then the model with the smaller AIC_c will be considered to be the better model. Furthermore, the results of an approximate likelihood ratio test (also called goodness of fit test) based on the F test [1] was used to assess the model improvement of GWR over OLS models. A statistically significant result from the F test (at a given level of significance) indicates that the GWR model performs better than the analogous OLS model.

Residual Analyses were conducted for both of the models, which included histograms of the OLS and GWR residuals. Global Moran's I and Geary's C statistics (which are special cases of the general cross product statistic applied to continuous data) were calculated to test the presence of spatial autocorrelation of the residuals. These statistics are particularly useful with spatial data whose covariance structures are defined by neighborhoods, a common situation for lattice data. Similar to the coefficient of correlation, the values of global Moran's I also range from -1 to 1 where the value of 1 and -1 indicate perfect positive spatial autocorrelation and perfect negative spatial autocorrelation, respectively. A global Moran's I with a value of 0 indicates no spatial autocorrelation. On the other hand, the values of Geary's C range from 0 to 2. It should be noted that Geary's C is interpreted much differently than Moran's I statistic. First, Geary's C has a mean of 1 under the null hypothesis of no spatial autocorrelation, and can never be negative. Second, values of Geary's C between 0 and 1 indicate positive autocorrelation, whereas values of C greater than 1 indicate negative spatial autocorrelation. The calculations of Moran's I and Geary's C were performed using the *spdep* package [15] in the statistical software package R, version 3.0.1 [14]. Spatial plots were used to detect the pattern of the OLS and GWR residuals.

3.1 Test of Spatial Heterogeneity of Individual Parameters of GWR

The individual parameter estimates of the GWR model are considered spatially heterogeneous if the inter-quartile range (difference between the third quartile and the first quartile) of GWR coefficients is greater than the range of values at ± 1 standard deviations of the respective global estimates estimates i.e., $Q_3 - Q_1$ of GWR coefficients is greater than the range of

$(b_1 - 1 S.D., b_1 + 1 S.D.)$ of OLS estimates which is simply twice the standard deviation ($S.D.$) of the OLS estimates [1,4,16].

3.2 Outputs from GWR

The spatial variation in parameter estimates, also considered as the main output from GWR, were used to see the variations in relationships revealed by these coefficients. In addition, histograms of the parameter estimates were used to detect the spatial variation revealed by these coefficients. Scatter plots of estimated GWR coefficients for the pairs of predictor variables of the regression model were used to visualize the nature of dependence in the estimated coefficients [17].

4 Applications

4.1 Data

The data set was obtained from the USDA National Agricultural Statistics Service [18]. North-Central region of the USA was selected for this observational study which included the states of Illinois, Indiana, Iowa, Minnesota, Montana, Nebraska, North Dakota, South Dakota, Wisconsin, and Wyoming. This selected region is located between 37.0 and 49.0 degrees north latitude and 84.8 and 116.0 degrees west longitude [19]. The data set consisted of the variables: County names, FIPS ID of the county, crop residue yield potential, temperature, and precipitation for the years 1970-2008. The county level dry crop residue yield potential ($Mg\ ha^{-1}$) was calculated using crop yield data collected from the USDA-NASS, 2009 [19]. The average daily mean temperature ($^{\circ}C$) and total precipitation (mm) of crop growing season (April – October) were the climate parameters used as independent (predictor) variables to model the county level crop residue potential (response / dependent variable) in this study. This climate data was obtained from the monthly gridded Parameter-elevation Regressions on Independent Slopes Model (PRISM) weather data [20]. For this analysis, the subset of the data for the three years 2008, 1980, and 1970 was considered. The response variable crop residue yield potential was closely normal and several transformations were attempted to normalize it. However, none of the transformations perfectly normalized it, and so it was kept as is for further analysis. Transformations were not attempted for the predictor variables. The data set had missing values of these variables for 56, 3, and 5 counties for 2008, 1980, and 1970 data respectively. Some examples of the counties with the missing data were Corson, Custer, Lawrence, Shannon, and Stanley. An imputation technique was used to estimate the missing values of crop residue yield potential, temperature, and precipitation. Thus, we had a complete data set for all 743 counties of the ten states of the North-Central region of the USA.

The original version of the data set with crop residue yield potential, temperature and precipitation for the years 1970 to 2008, has already been analyzed by averaging over time and then Conditional Autoregressive and Simultaneous Autoregressive models [21] were fit [22]. To the best of our knowledge, this dataset has not been analyzed using GWR yet.

4.2 Results and Discussion

OLS and GWR regression models were fit for the crop residue yield potential as a function of two climate variables (temperature and precipitation) for small as well as large data sets (66 counties of South Dakota and 743 counties of ten states of the North-Central region of the USA) for the years 1970, 1980, and 2008. In all the cases, the value of Moran's I and Geary's C for residuals (Table 1) indicated that there was a positive autocorrelation of OLS residuals. For instance, for the year 1970, the value of Moran's I for OLS residuals was 0.33 while that of Geary's C was 0.63 (both agreeing on positive autocorrelation of residuals). On the other hand, the values of Moran's I and Geary's C (Table 1) indicated that there was a negative autocorrelation of GWR residuals, for example, the value of Moran's I for GWR residuals was -0.12 and Geary's C was 1.11 (both agreeing on small negative autocorrelation of residuals). It should also be noted that in all the years, the degree of autocorrelation of OLS residuals was higher than that of GWR residuals indicating that GWR mitigated the autocorrelation of residuals. This was also supported by the spatial plot of OLS and GWR residuals (Figs. 1 and 2). In addition, there was more variability of OLS residuals than GWR residuals, for instance, in 2008 data, the OLS residuals ranged from -5.286 to 7.1 while the GWR residuals ranged from -4.23 to 3.35. The variances of OLS and GWR were 3.92 and 0.84 respectively. These numbers helped us quantify the visible relationships between the two map patterns illustrated in Figs. 1 and 2. The residuals under all the OLS models were non normal, while the residuals under all the GWR models were almost normal (based on histograms, Fig. 3). Significant improvements in the fit of both the OLS and GWR models were observed for all models on comparing R^2 and AIC_C . In all cases, the values of R^2 under GWR models were higher than those of the corresponding OLS models (Table 1) indicating the higher explanatory power of the GWR models over the analogous OLS models. In addition, all the GWR models had smaller AIC_C values than the analogous OLS models (Table 1).

Table 1. Comparison of OLS and GWR

Data	Year	Moran's I		Geary's C		Adjusted R^2		AIC_C	
		OLS	GWR	OLS	GWR	OLS	GWR	OLS	GWR
66 counties	1970	<i>0.33</i>	-0.12	<i>0.63</i>	1.11	0.08	0.74	57.00	37.63
	1980	<i>0.54</i>	0.05	<i>0.45</i>	0.99	0.45	0.94	132.53	111.41
	2008	<i>0.58</i>	-0.01	<i>0.40</i>	0.95	0.01	0.74	287.48	237.50
743 counties	1970	<i>0.62</i>	-0.07	<i>0.39</i>	1.06	0.22	0.84	2645.44	1955.68
	1980	<i>0.47</i>	-0.05	<i>0.53</i>	1.06	0.46	0.87	2606.92	2047.67
	2008	<i>0.43</i>	0.21	<i>0.56</i>	0.96	0.39	0.87	3131.34	2412.15

The adjusted R^2 for GWR models are the mean values of the corresponding R^2 values. Moran's I and Geary's C for OLS residuals (italicized) had statistically significant p -values at the 0.05 level. The p -values for the Moran's I and Geary's C for GWR residuals were not statistically significant.

The GWR models in all cases performed better than the corresponding OLS models, according to the approximate likelihood ratio test (Table 2) with extremely small p -values in each case suggesting that the ordinary least squares relationship between crop residue yield potential, temperature, and precipitation was not constant across the study area. Table 2 also showed the significant improvements in explained variance (reduced sum of squared errors) by the GWR models over the analogous OLS models.

Table 2. Goodness of fit test for improvement in model fit of GWR over OLS

Data	Year	Source			F-statistic	P-value
		OLS residuals	GWR residuals	GWR improvement		
66 counties	1970	8.12	2.40	5.73	3.54	6.7e-06
	1980	25.50	3.09	22.40	5.11	7.1e-08
	2008	266.76	72.65	194.11	7.14	1.0e-11
743 counties	1970	1514.04	317.75	1196.29	7.53	< 2.2e-16
	1980	1437.55	336.52	1101.03	5.88	< 2.2e-16
	2008	2911.72	626.50	2285.22	8.27	< 2.2e-16

Note: $F = \frac{22852/22645}{626.5/51355} = 8.27$ for 743 counties of 2008 data as given by GWR output.

It was observed that the inter-quartile ranges (IQR) of GWR coefficients were greater than twice the standard deviations of their respective global estimates for small data set with 66 counties of South Dakota for the years 1970 and 1980 (Table 3). But the IQRs of the intercept and precipitation of GWR coefficients were less than twice the standard deviations of their respective global estimates in the year 2008 model for the small data set. On the other hand, for the large data set with 743 counties of the North-Central region of the USA in all the years, the IQRs of all the GWR coefficients were greater than twice the standard deviations of their respective global estimates. Thus, the parameters for intercept and precipitation for the small data set were not spatially heterogeneous in the 2008 model. However, all the parameters of models for the large data set were spatially heterogeneous. We again repeated the experiment by fitting the GWR model for a data set with 194 counties of the two states Illinois and Indiana of the USA (results not shown). We observed that the IQRs of GWR coefficients were greater than twice the standard deviations of their respective global coefficients. Hence, a caution needs to be exercised when using GWR for the purpose of assessing spatial heterogeneity of parameters when analyzing small data sets.

Table 3. Test of spatial heterogeneity of individual parameters

Data	Year	Variable	S.D.(Global)	2S.D.(Global)	IQR (GWR)	Results
66 Counties	1970	Intercept	0.74	0.148	2.134	$Q_3-Q_1 > 2$ S.D.
		Temp	0.041	0.082	0.136	$Q_3-Q_1 > 2$ S.D.
		Precpt	0.0011	0.0022	0.004	$Q_3-Q_1 > 2$ S.D.
	1980	Intercept	1.44	2.88	9.394	$Q_3-Q_1 > 2$ S.D.
		Temp	0.08	0.16	0.649	$Q_3-Q_1 > 2$ S.D.
		Precpt	0.0013	0.0026	0.006	$Q_3-Q_1 > 2$ S.D.
	2008	Intercept	5.6	11.2	8.31	$Q_3-Q_1 < 2$ S.D.
		Temp	0.335	0.670	0.691	$Q_3-Q_1 > 2$ S.D.
		Precpt	0.007	0.014	0.011	$Q_3-Q_1 < 2$ S.D.
743 Counties	1970	Intercept	0.34	0.68	10.59	$Q_3-Q_1 > 2$ S.D.
		Temp	0.027	0.054	0.67	$Q_3-Q_1 > 2$ S.D.
		Precpt	0.0004	0.0008	0.006	$Q_3-Q_1 > 2$ S.D.
	1980	Intercept	0.36	0.72	11.58	$Q_3-Q_1 > 2$ S.D.
		Temp	0.022	0.044	0.75	$Q_3-Q_1 > 2$ S.D.
		Precpt	0.0003	0.0006	0.008	$Q_3-Q_1 > 2$ S.D.
	2008	Intercept	0.52	1.04	20.2	$Q_3-Q_1 > 2$ S.D.
		Temp	0.04	0.08	1.43	$Q_3-Q_1 > 2$ S.D.
		Precpt	0.0005	0.001	0.012	$Q_3-Q_1 > 2$ S.D.

Note: Temp and Precpt stand for temperature and precipitation respectively.

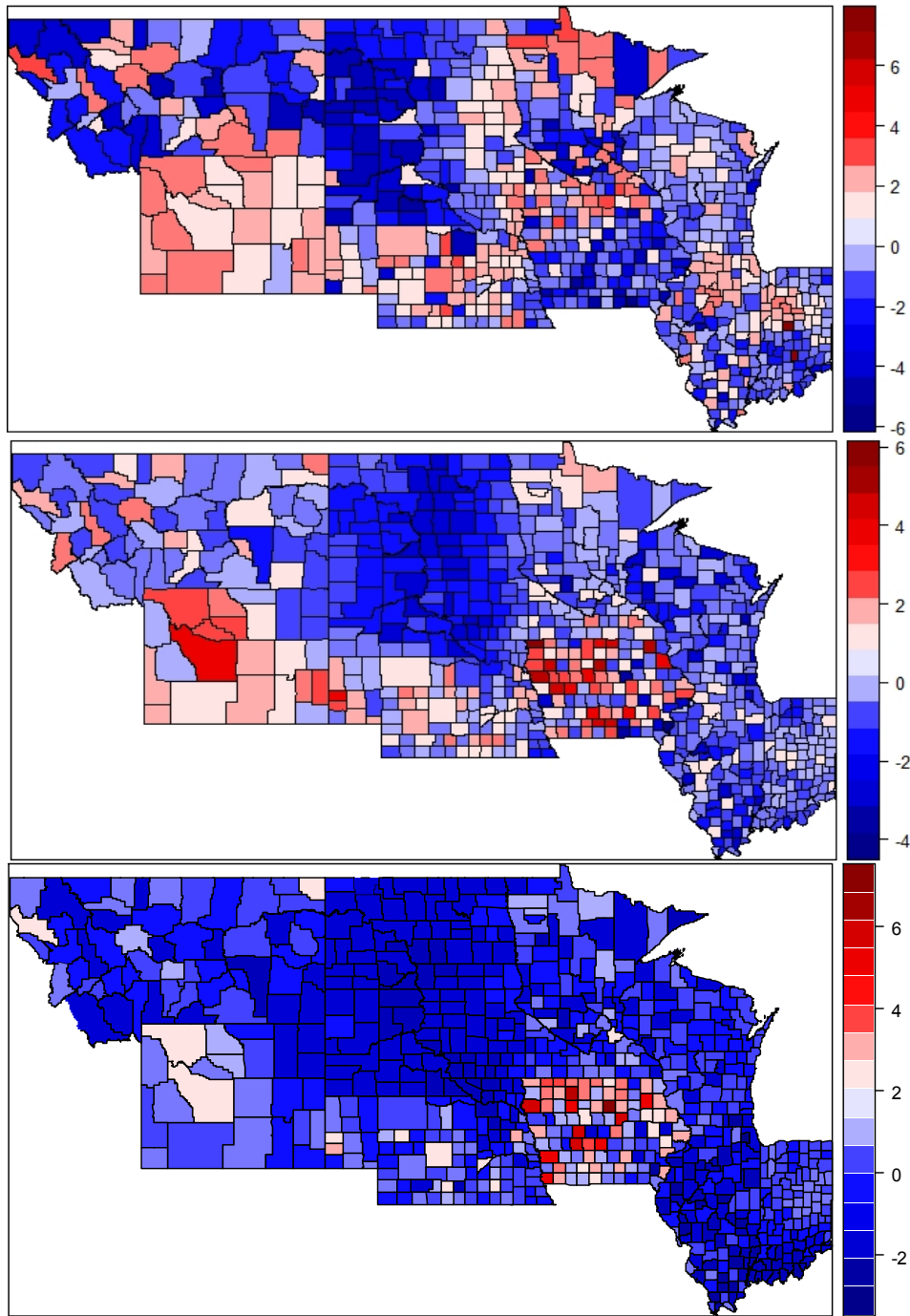


Fig. 1. Spatial distribution of OLS residuals for the ten north central states of the US. Top to bottom: 2008, 1980, and 1970

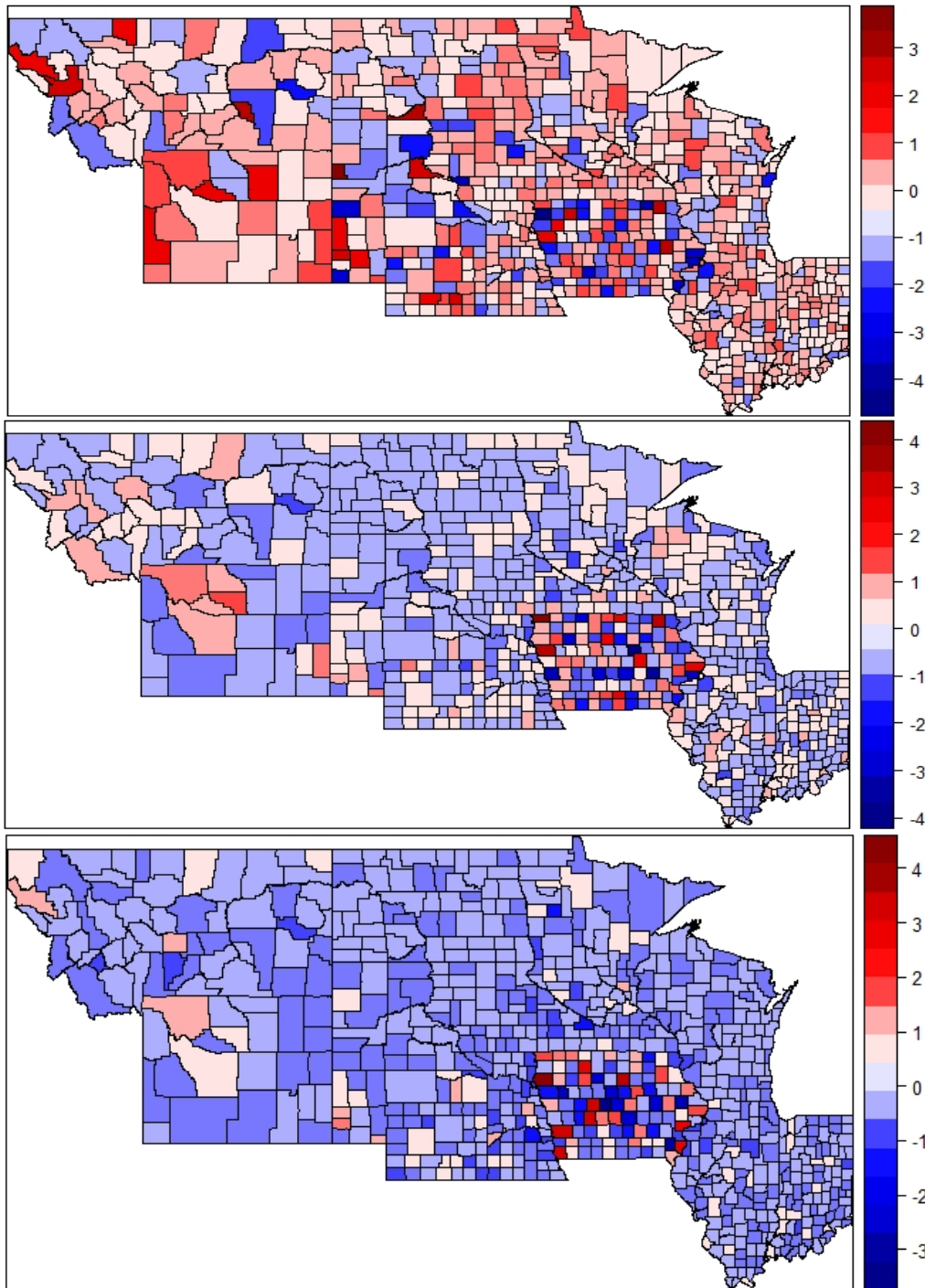


Fig. 2. Spatial distribution of GWR residuals for the ten north central states of the US. Top to bottom: 2008, 1980, and 1970

The normality of OLS and GWR residuals for all the models was compared using histograms (Fig. 3). It was found that the GWR residuals for all the models were more closely normally distributed than OLS residuals. It supported our previous conclusion that GWR had the ability to better model spatially varying data. In addition, the non-normality of OLS residuals implied that the assumptions of OLS regression models were not satisfied, raising a question about the validity of these models. When the assumptions of OLS models are violated, the results of such models could be misleading, and statistical inferences provided by such models could be spurious. In addition, the strong autocorrelation of OLS residuals suggests that the OLS relationship between crop residue yield potential, temperature, and precipitation was not stationary. Moreover, the lower degree of spatial autocorrelation of GWR residuals implied that the GWR models were more trustworthy. Hence, the non-normality and spatial autocorrelation of the OLS residuals and lower values of the coefficients of determination of OLS models supported the assertion that GWR modelled the spatially varying data more efficiently than the analogous OLS models.

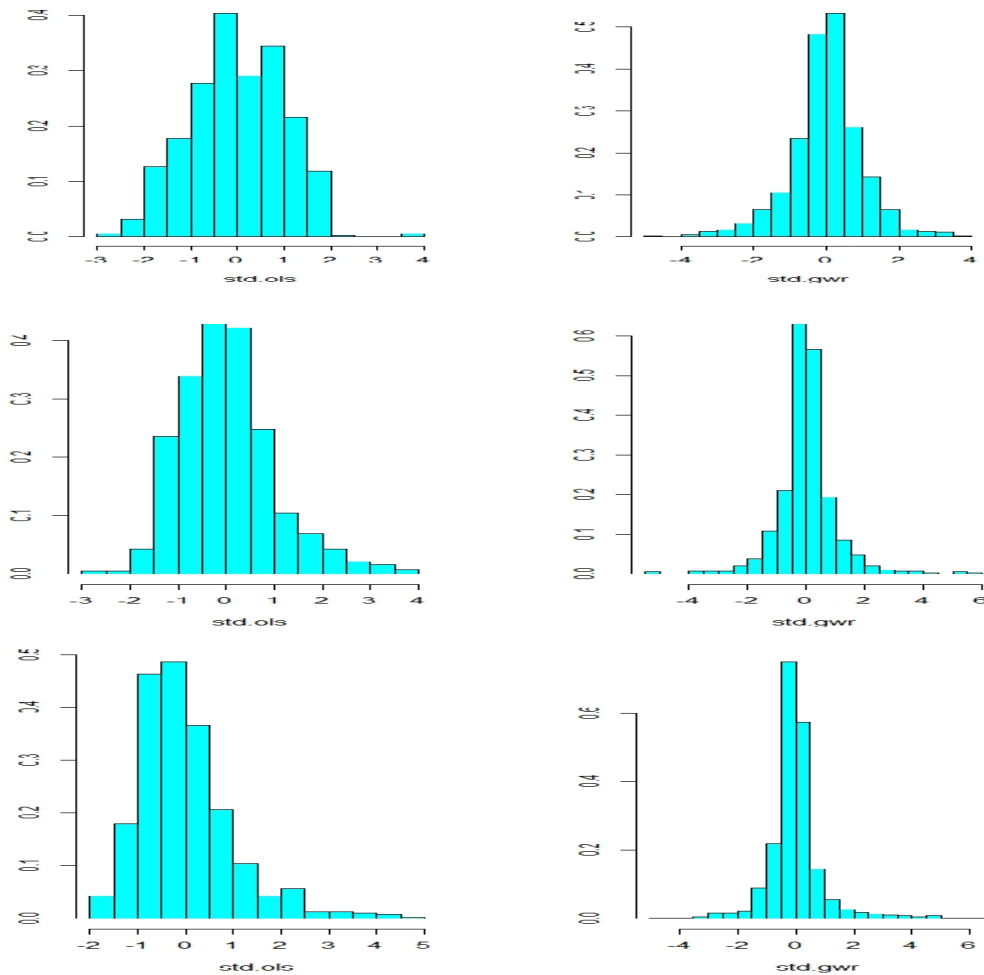


Fig. 3. Histogram of residuals. Top to bottom: 2008, 1980, and 1970; left to right: OLS, and GWR

The histograms of estimated GWR coefficients for the years 2008, 1980, and 1970 are given in Fig. 4. These histograms allowed the distributions of the coefficients to be easily compared for the three years. The distributions for intercept were right skewed, which indicated that a few counties

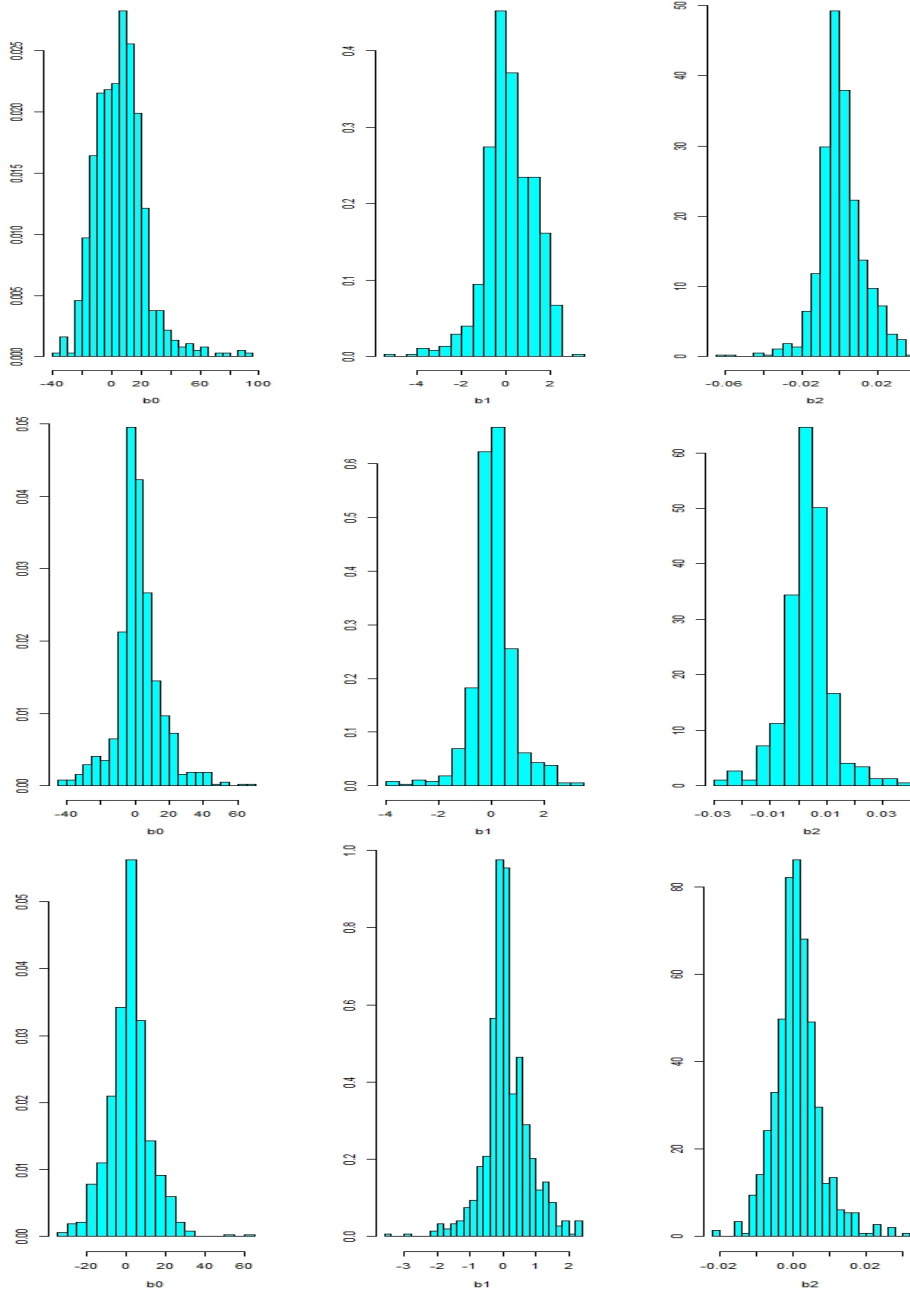


Fig. 4. Histogram of estimated GWR coefficients. Top to bottom: 2008, 1980, and 1970; left to right: coefficients of intercept, temperature, and precipitation

had much larger intercepts than other counties. The distributions for temperature coefficients were left skewed, which indicated that a few counties had larger negative coefficients for temperature than other counties. The distributions for precipitation coefficients were left skewed for 2008, somewhat symmetrical for 1980, and right skewed for 1970. This indicates that the shape of the distribution of coefficients for precipitation has changed over time.

As recommended by [17], scatter plots of estimated GWR coefficients for the pairs of predictor variables of the regression model were used to visualize the nature of dependence in the estimated coefficients. Fig. 5 showed the scatter plots for the pairs of the three regression terms: intercept, temperature, and precipitation. The vertical and the horizontal dashed reference lines denote the levels of analogous global parameter estimates (Table 4 has the parameter estimates for the OLS regression models for 2008, 1980, and 1970 models). The Pearson product moment correlation coefficients for the pairs of the regression terms were given in Table 5. These results strongly agreed with the scatter plots from Fig. 5. There was negative correlation between the coefficients of correlations (with the exception of the coefficients between intercept and precipitation for 1970). The degree of correlation between intercept and temperature was very strong. The correlations between the other pairs of regression terms were comparatively weaker. The scatter plots showed that there was large variation around the OLS estimates for intercept; however, there was small variation around the OLS estimates for temperature and precipitation. Furthermore, the estimated GWR coefficients do not center on their analogous OLS estimates. The strong correlation between the coefficients of intercept and temperature could signal the presence of local collinearity in the GWR model which could lead to complications on the statistical inference drawn from these GWR coefficients.

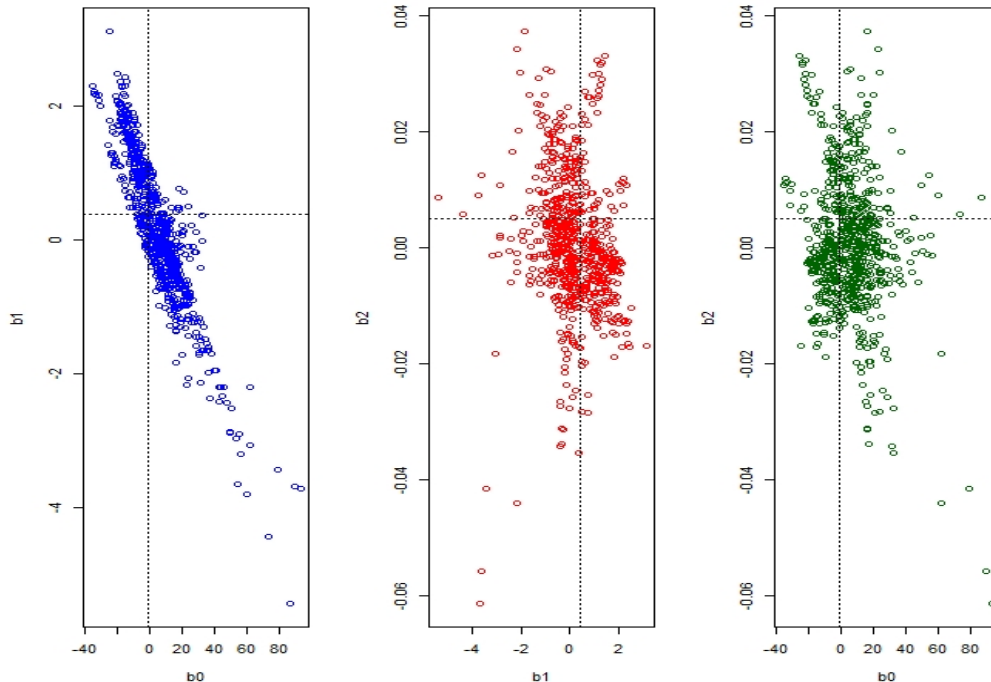


Fig. 5 (a). Scatter plots of estimated GWR coefficients for 2008 model

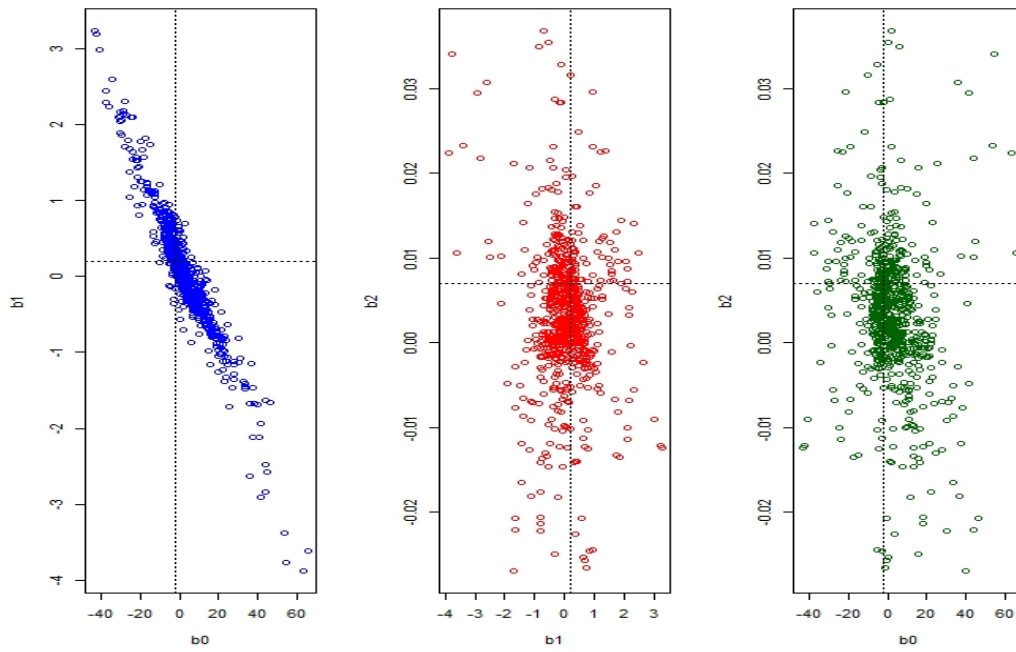


Fig. 5 (b). Scatter plots of estimated GWR coefficients for 1980 model

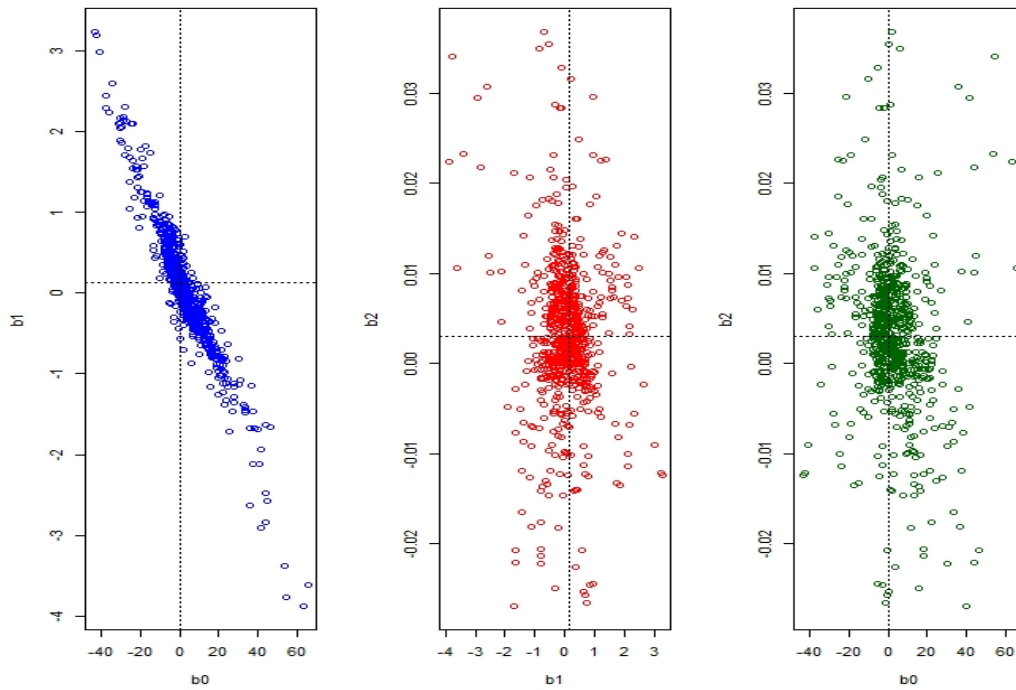


Fig. 5 (c). Scatter plots of estimated GWR coefficients for 1970 model

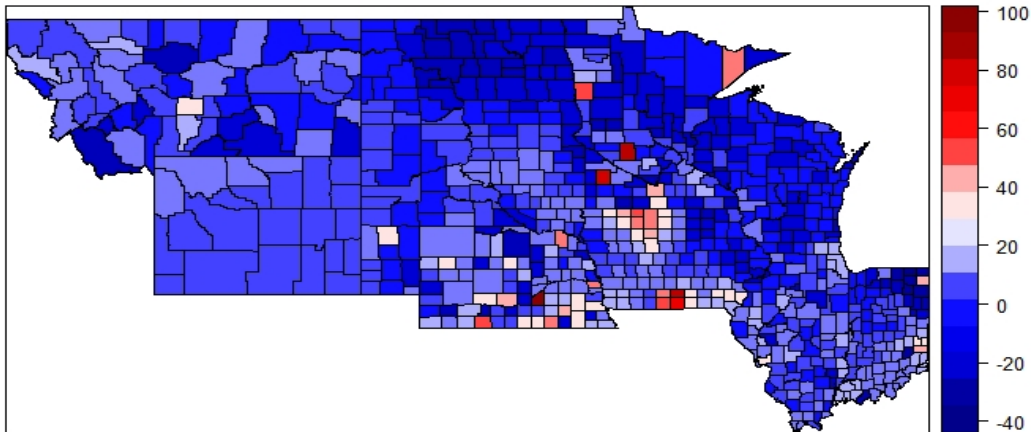
Table 4. Parameter estimates for the OLS regression models

Year	Intercept	Temperature	Precipitation
2008	-1.16	0.393	0.005
1980	-1.888	0.195	0.007
1970	0.145	0.126	0.003

Table 5. Correlation coefficients between the estimated GWR coefficients

Year	cor(b_0, b_1)	cor(b_1, b_2)	cor(b_0, b_2)
2008	-0.911	-0.144	-0.234
1980	-0.955	-0.097	-0.152
1970	-0.951	-0.333	0.0734

The spatial variation of the estimated GWR coefficients for the years 2008, 1980 and 1970 for intercept, temperature, and precipitation were shown in Figs. 6, 7, and 8 respectively. The spatial variations in relationship shown by these distributions were interesting. A comparative study of Figs. 6 and 7 showed a clear inverse map pattern between intercept and temperature (as one would expect, since the correlation coefficients between intercept and temperature from Table 5 were large and negative in all the cases): in general, when the intercept parameter was high, the temperature parameter was low and vice versa. On the other hand, there was no clear pattern between the maps of intercept and precipitation and temperature and precipitation (Figs. 7 and 8, as the coefficients of correlations between both the pairs are small in both the cases). Thus, the correlation coefficients given in Table 5 helped us quantify the three map patterns illustrated in Figs. 6, 7, and 8. For instance, the coefficients of correlations between temperature and precipitation in the years 2008 and 1980 were small. Due to this, there was not much difference in the patterns of the plots of temperature and precipitations for these years.



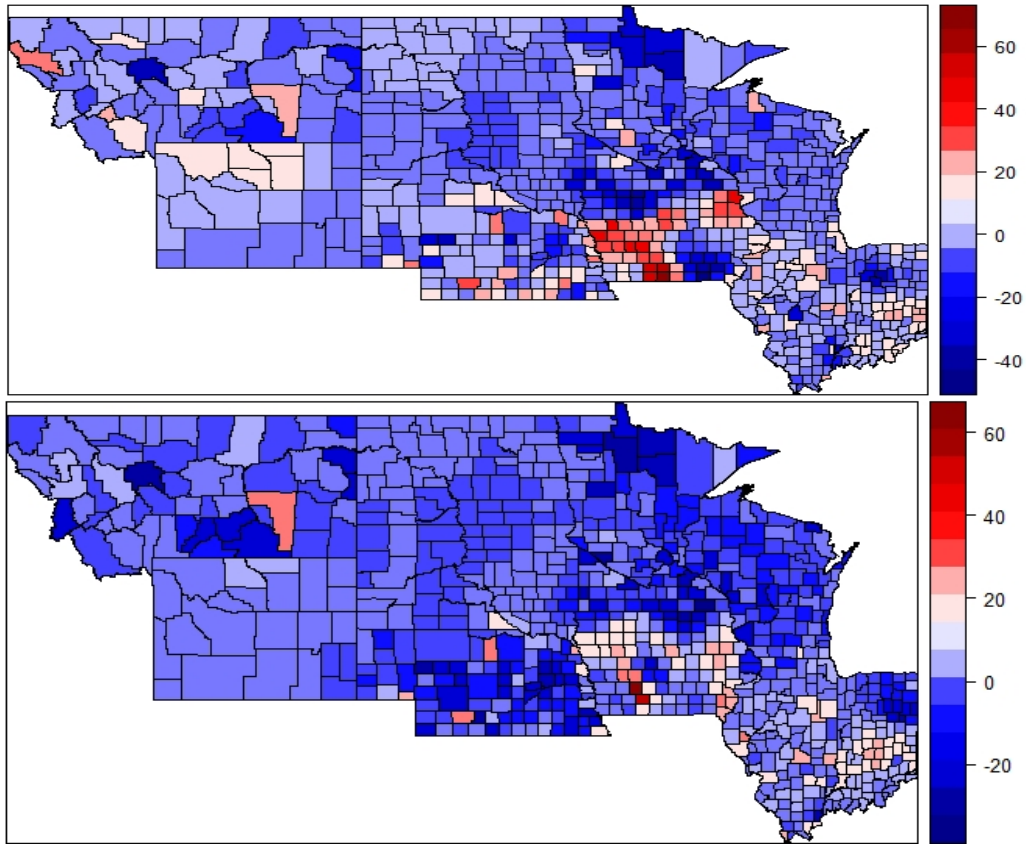
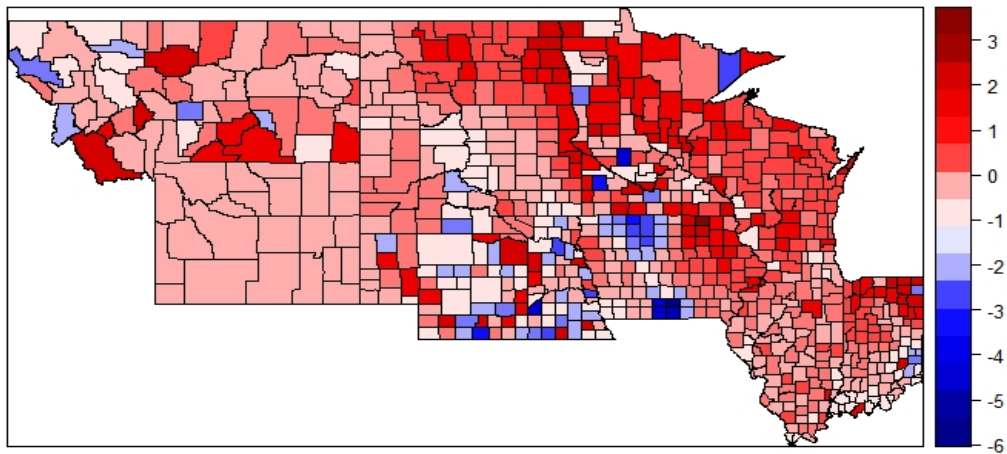


Fig. 6. Spatial distribution of estimated GWR coefficients for intercept for the ten north central states of the US. Top to bottom: 2008, 1980, and 1970



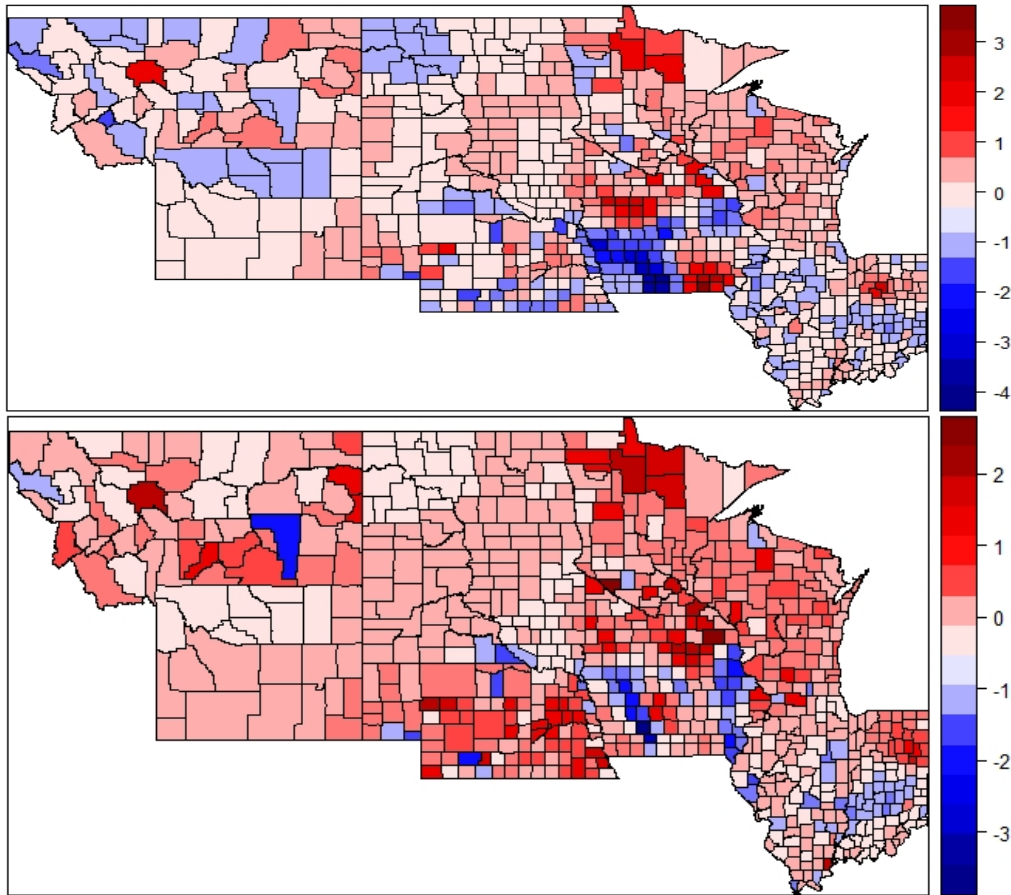
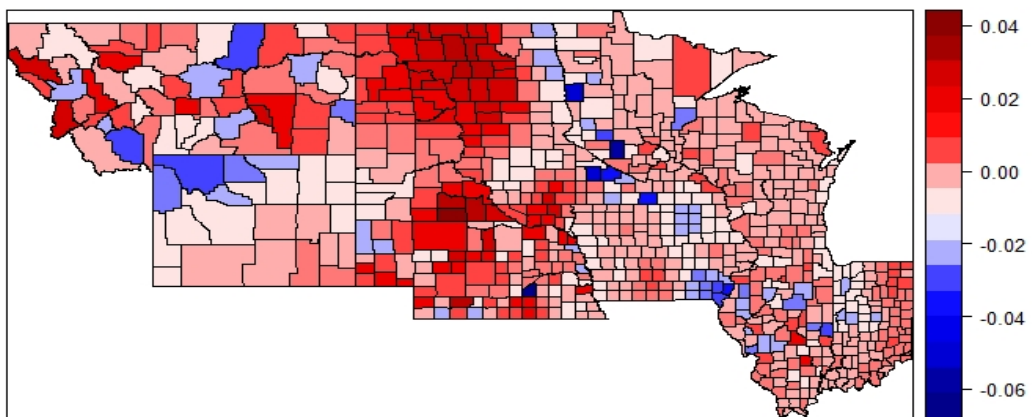


Fig. 7. Spatial distribution of estimated GWR coefficients for temperature for the ten north central states of the US. Top to bottom: 2008, 1980, and 1970



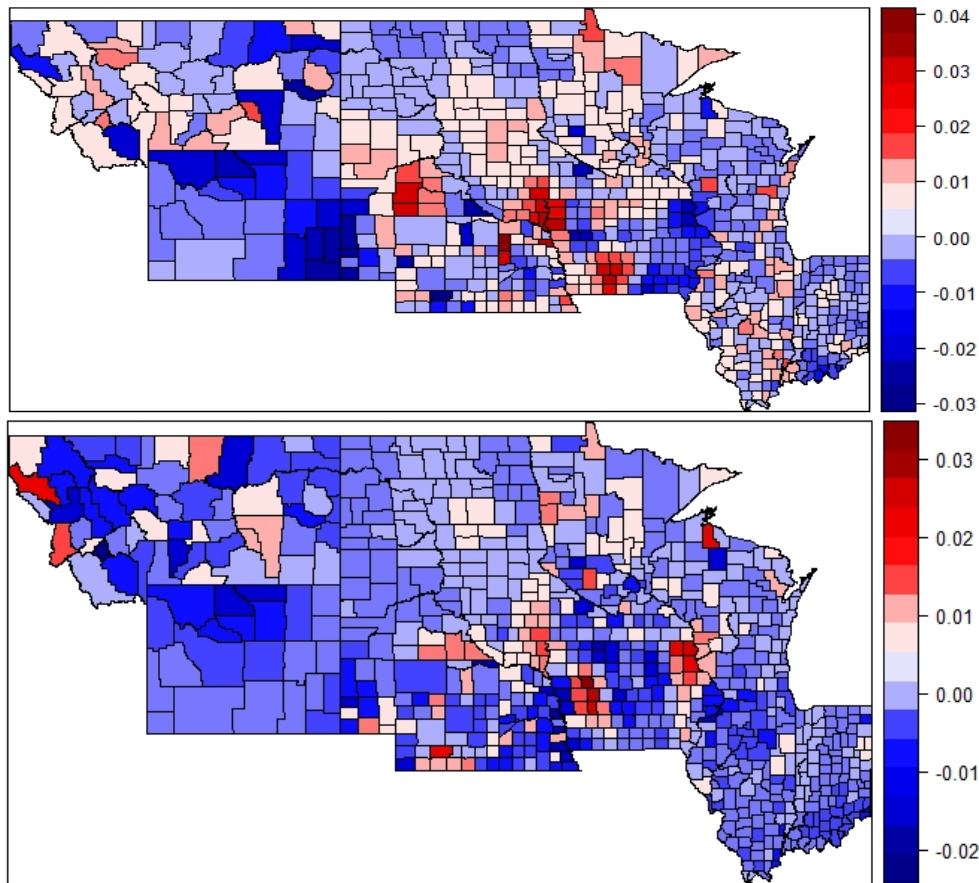


Fig. 8. Spatial distribution of estimated GWR coefficients for precipitation for the ten north central states of the USA. Top to bottom: 2008, 1980, and 1970

5 Summary

In this study, OLS and GWR models were fit for crop residue yield potential of South Dakota, United States for the years 2008, 1980, and 1970 and the performance of these models were compared. The same analyses were repeated with a larger data set (crop residue yield potential data of ten states of the North-Central region of the USA). The approximate likelihood ratio test (also called the goodness of fit test) suggested that all the GWR models performed better than the analogous OLS regression models indicating that OLS relationship between crop residue yield potential, temperature, and precipitation was not constant across the study area. The coefficients of determination in all the GWR models were higher than in the analogous OLS models, indicating that GWR models had higher power at explaining spatially varying relationships than the analogous OLS models. All the GWR models had smaller AIC_C values than those of the analogous OLS models by at least 3 indicating that GWR models performed better and provided more parsimonious results than the analogous OLS models. In this study, the residuals of all the GWR models were approximately normally distributed while those of OLS models were not

normally distributed. Moran's I and Geary's C statistics showed that the degree of autocorrelation of OLS residuals were higher than that of GWR residuals indicating that GWR mitigated the autocorrelation of residuals.

Out of the three cases of our analysis, two of the three parameters of GWR model (using small data set for 2008) were not spatially heterogeneous. This could be due to the small sample size, which increases the possibility of spurious correlations between local coefficients [6]. In the other two cases, i.e. 1980 and 1970; all the parameters of the GWR model were spatially heterogeneous for small and large data sets. GWR should not be practiced for data with small sample size; and serious caution should be exercised when interpreting the results from such outputs [6].

In our case, there was strong (negative) correlation between the estimated GWR coefficients of intercept and temperature. However, the correlations between the other pairs of coefficients were weaker. The strong correlation between the GWR coefficients of intercept and temperature could be the result of local collinearity in the model. This could be investigated using maps of approximate local regression coefficient correlations [17]; local variance inflation factors (VIFs), variance-decomposition proportions and the associated condition indices [23, 24] as diagnostic tools for collinearity when estimating GWR coefficients.

Acknowledgements

We thank the anonymous reviewers for their helpful comments on an earlier draft.

Competing Interests

Authors have declared that no competing interests exist.

References

- [1] Fotheringham AS, Brunson C, Charlton. Geographically weighted regression. M. 1st ed. New York: Wiley; 2002.
- [2] Brunson C, Fotheringham S, Charlton M. Geographically weighted regression-modeling spatial non-stationarity. *The Statistician*. 1998;47:431-43.
- [3] Zhang L, Bi H, Cheng P, Davis, CJ. Modeling spatial variation in tree diameter–height relationships. *Forest Ecology and Management*. 2004;189(1):317-329.
- [4] Tulbure, MG, Wimberly MC, Roy DP, Henebry GM. Spatial and temporal heterogeneity of agricultural fires in the central United States in relation to land cover and land use. *Landscape Ecology*. 2011;26(2):211-224.
- [5] Brown S, Versace VL, Laurenson L, Ierodiaconou D, Fawcett J, Salzman S. Assessment of spatiotemporal varying relationships between rainfall, land cover and surface water area using geographically weighted regression. *Environmental Modeling & Assessment*. 2012;17(3):241-54.

- [6] Páez A, Farber S, Wheeler D. A simulation-based study of geographically weighted regression as a method for investigating spatially varying relationships. *Environment and Planning-Part A*. 2011;43(12):2992.
- [7] Charlton M, Fotheringham S, Brunson C. Geographically weighted regression. White paper. National Centre for Geocomputation. National University of Ireland Maynooth; 2009.
- [8] Brunson C, McClatchey J, Unwin DJ. Spatial variations in the average rainfall–altitude relationship in Great Britain: an approach using geographically weighted regression. *International Journal of Climatology*. 2001;21(4):455-66.
- [9] McMillen DP. One hundred fifty years of land values in Chicago: a nonparametric approach. *Journal of Urban Economics*. 1996;40(1):100-124.
- [10] Brunson C, Fotheringham AS, Charlton ME. Geographically weighted regression: a method for exploring spatial nonstationarity. *Geographical analysis*. 1996;28(4):281-98.
- [11] Farber S, Páez A. A systematic investigation of cross-validation in GWR model estimation: empirical analysis and Monte Carlo simulations. *Journal of Geographical Systems*. 2007;9(4):371-396.
- [12] Tu J, Xia ZG. Examining spatially varying relationships between land use and water quality using geographically weighted regression I: model design and evaluation. *Science of the total environment*. 2008;407(1):358-378.
- [13] Bivand R, Yu D. spgwr: Geographically weighted regression. R package version 0.6-22; 2013.
Available: <http://CRAN.R-project.org/package=spgwr>.
- [14] R Development Core Team. R: a language and environment for statistical computing. R Foundation for Statistical Computing. Vienna, Austria; 2013.
- [15] Bivand R. spdep: Spatial dependence: weighting schemes, statistics and models. R package version 0.5-60; 2013.
Available: <http://CRAN.R-project.org/package=spdep>.
- [16] Charlton M, Fotheringham S, Brunson C. Geographically Weighted Regression, ESRC National Center for Research Methods; 2006.
- [17] Wheeler D, Tiefelsdorf M. Multicollinearity and correlation among local regression coefficients in geographically weighted regression. *Journal of Geographical Systems*. 2005;7(2):161-187.
- [18] NASS [Database on Internet] USDA; National Agricultural Statistical Service [cited 2012 Jan 1] Annual Data. Available from: <http://www.nass.usda.gov>. Files updated annually.

- [19] Chintala R, Wimberly MC, Djira GD, Tulbure MG. Interannual variability of crop residue potential in the north central region of the United States. *Biomass and Bioenergy*. 2013;49:231-38.
- [20] Daly C, Gibson WP, Dogget M, Smith J, Taylor G. Up-to-date monthly climate maps for the conterminous United states. Proceedings of 14th American Meteorological Society Conference on Applied Climatology, 84th American Meteorological Society Annual Meeting Combined Preprints; 13e16 Jan 2004; Seattle, USA. 2004. Paper P5.1, CD-ROM. Available: <http://ams.confex.com/ams/pdfpapers/71444.pdf>.
- [21] Schabenberger O, Gotway CA. Statistical methods for spatial data analysis (Vol. 64). Chapman and Hall/CRC; 2004.
- [22] Chintala R, Djira G, Devkota M. Modeling the Effect of Climate Variables on Crop Residue Potential for the North Central Region of the United States. *Agronomy Journal* (Submitted).
- [23] Belsley D. Conditioning diagnostics: collinearity and weak data in regression. Wiley, New York; 1991.
- [24] Wheeler DC. Diagnostic tools and a remedial method for collinearity in geographically weighted regression. *Environment and Planning A*. 2007;39(10):2464-2481.

© 2014 Devkota et al.; This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/3.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Peer-review history:

The peer review history for this paper can be accessed here (Please copy paste the total link in your browser address bar)

www.sciencedomain.org/review-history.php?iid=276&id=6&aid=2174