

SPATIAL STATISTICAL MODELING OF EXPLICIT AND IMPLICIT  
DETERMINANTS OF REGIONAL FERTILITY RATES  
- A CASE STUDY OF HE-NAN PROVINCE, CHINA

H. Michael Feng  
Department of Geography, Syracuse University  
144 Eggers Hall, Syracuse, NY 13244-1090  
hfeng@rodan.acs.syr.edu

I. Preliminary Considerations of the Spatial Statistical Application

This paper attempts to use spatial statistical techniques to simultaneously model effects of explicit and implicit fertility determinant variables. Explicit variables refer to those whose values may be measured and indeed are computed at various levels of surveys, and implicit variables refer to those whose values may not be collected in such a manner. I contend that implicit variables may be analyzed by spatial statistical methods when they possess three characteristics: [1] The causal relationship between such implicit variables and their response variable is established through theory; [2] the objects these implicit variables act on may be collectively measured as relatively homogeneous areal units, and their effects are generally measurable, although it may not be immediately separable from the effects of explicit variables; and, [3] the mechanism by which these variables take effect consists of either human groups responding simultaneously to a common source, or interacting with each other. In other words, the effects of such variables may be found either as a spatial trend, or as spatial correlation and autocorrelation.

II Model Specification

The technique is demonstrated through a case study of He-Nan Province of China (see map 1 for location of the province). The dataset is based on the Fourth National Population Census publications (He-Nan Province Census Bureau, 1992). Selection of the province has to do with the fact that He-Nan is one of the most representative of the entire country in population density, living standards, economic structure, and cultural characteristics. The census took place at 0:00 am, July 1, 1990 for all provinces of China. Indices are enumerated at county level and may be aggregated into levels of prefecture and province, which represent the three basic areal units of Chinese administrative hierarchy. The raw data employed in this study may be categorized under three general groups, namely, (1) demographic structure, (2) economic structure, and (3) educational attainment. In later spatial analysis a fourth group of variables is added, which mainly pertains to the effects of government population policies.

The specified model is composed of two sub-models. The first is the following general linear regression model, which ignores relative geographical locations:  $[\text{Fertility Rate}] = [\text{GroupI}]b_1 + [\text{GroupII}]b_2 + [\text{GroupIII}]b_3 + e$ , where  $[\text{Fertility Rate}]$  is the vector whose elements are the fertility rates for each county, and  $[\text{GroupI}]$  through  $[\text{GroupIII}]$  are variable matrices where elements are the three groups of variables described above.  $b_1$ ,  $b_2$  and  $b_3$  are parameter vectors, and  $e$  is an error vector.

## SPATIAL STATISTICAL MODELING

The second, or spatial sub-model, starts with an inspection of this error vector,  $\mathbf{e}$ . The values of  $\mathbf{e}$  may be considered as the portion of variation in [Fertility Rate] that is *not* associated with any of the explicit variables in the first model. In other words, it is the variation associated with implicit variables as well as random noise. It is at this juncture that the relative location and interaction of the regions is introduced.

The first question addressed in building the spatial sub-model asks whether or not any spatial pattern is latent in the error terms. If the answer is no, then a case can be made for collectively denying the relevance of effects of implicit variables. But if the answer is yes, then further investigation is in order. Detection of a spatial pattern relies on two indices, namely the Moran Coefficient (I) and Geary Ratio (c) (see Griffith, 1987, Griffith and Amrhein, 1991, pp115-43). Their computational formulae are

;

where  $i$  and  $j$  denote areal units,  $y_i$  and  $y_j$  are the attribute values for areal units  $i$  and  $j$ , and  $c_{ij} = 1$  if units  $i$  and  $j$  are adjacent and 0 otherwise.

When a spatial pattern is detected, we need to identify its form and to provide an explanation in light of it. This starts with a basic assumption about the shape of the autocorrelation model. Suppose it is specified as  $\mathbf{B}(\mathbf{Y}-\boldsymbol{\mu})=\mathbf{D}\mathbf{e}$  ... (1) or  $\mathbf{Y}=\boldsymbol{\mu}+\mathbf{B}^{-1}\mathbf{D}\mathbf{e}$  ... (2), where  $\mathbf{Y}$  is  $n$ -by-1 vector of spatially correlated variables under study.  $\mathbf{B}$  and  $\mathbf{D}$  are nonsingular parameter matrices with  $\mathbf{B}=\{b_{ij}\}$  and  $b_{ii}=1$  for all  $i$ .  $\mathbf{e}$  is a  $n$ -by-1 vector of residuals that are spatially independent of each other. It should have mean zero and covariance matrix  $s^2\mathbf{V}_e$ . Note that since its elements are independent of each other,  $\mathbf{V}_e$  must be a diagonal matrix. In addition, if we assume a constant variance among them, the diagonal elements of  $\mathbf{V}_e$  also should have identical values of one.

Practically the model says, after certain operations on  $(\mathbf{Y}-\boldsymbol{\mu})$ , we should be able to filter out its spatial autocorrelation structure. In other words, we should be able to do something to it so that it becomes a new vector  $\mathbf{D}\mathbf{e}$  that does not contain any spatial autocorrelation. This transformation is done through matrix  $\mathbf{B}$ . The theoretical information we are seeking should be revealed by the format of such a transformation. One particular format starts with the assumption that  $\mathbf{D}=\mathbf{I}$  and  $\mathbf{B}=(\mathbf{I}-\mathbf{S})$ , where  $\mathbf{I}$  is identity matrix. The meaning of  $\mathbf{S}$  will be clear in latter discussion. For now, suffice it to say that  $(\mathbf{I}-\mathbf{S})$  must be invertable, and the diagonal elements of  $\mathbf{S}$  are all zero. In that situation, equation (2) may be written as  $\mathbf{Y}=\boldsymbol{\mu}+\mathbf{S}(\mathbf{Y}-\boldsymbol{\mu})+\mathbf{e}$ . ... (3) In algebraic notation this is equivalent to  $y_i=\mu_i+\sum_j s_{ij}(y_j-\mu_j)+e_i$ . ... (4) The last equation makes the interdependency structure a little more intuitive. It says that the value at areal unit  $i$  depends on the mean of its distribution at that areal unit ( $\mu_i$ ), plus a function of the errors at the areal units connected to it ( $\sum_j s_{ij}(y_j-\mu_j)$ ).

Until this point, the meaning of matrix  $\mathbf{S}$  has not been fully explained. In fact it is similar to but a little more than the connectivity matrix discussed earlier. It determines whether areal units  $i$  and  $j$  are correlated (the connectivity matrix) and how strong that correlation is ( $r$ ). One particular form of  $\mathbf{S}$  that is widely adopted among spatial statisticians is:  $\mathbf{S}=r\mathbf{W}$ , ... (5) where  $r$  is a constant spatial autocorrelation parameter and  $\mathbf{W}$  is a row standardized version of connectivity matrix  $\mathbf{C}$ , with  $\{w_{ij}\}=c_{ij}(\sum_j c_{ij})^{-1}$ , and  $\{c_{ij}\}=1$ , if  $j$  is connected to  $i$ , and  $\{c_{ij}\}=0$  otherwise. Equation (3) then becomes  $\mathbf{Y}=\boldsymbol{\mu}+r\mathbf{W}(\mathbf{Y}-\boldsymbol{\mu})+\mathbf{e}$ , ... (6) and equation (4) becomes  $y_i=\mu_i+r\sum_{j \in N(i)} w_{ij}(y_j-\mu_j)+e_i$ , ... (7) where  $N(i)$  denotes the set of areal units that are connected to areal unit  $i$ , that is, the set for which  $w_{ij} \neq 0$ . Equations (6) and (7) specify an autoregressive model with no explanatory variable. A full autoregressive model, or an AR model, arises when explanatory variables are required. It states that the

realization of response variable  $Y$  at areal unit  $i$  is a function of its expected realization at  $i$  ( $Xb$ ), plus realizations of  $Y$  at locations connected to  $i$  ( $rWY$ ), plus an error term  $e$ . This may be expressed by the following equation:  $Y = Xb + rWY + e$ , ... (8) where  $X$  is the matrix of explanatory variables,  $X_{(n-by-p)} = (X_{1(n-by-1)}, X_{2(n-by-1)}, \dots, X_{p(n-by-1)})$ , and  $b$  is the parameter vector,  $b = (b_0, b_1, b_2, \dots, b_p)^T$ ,  $p$  being number of explanatory variables, and other variables and parameters being defined as before.

### III A Classical Linear Regression Model of Explicit Variables

The linear regression model, the first of the aforementioned two, depicts the influence of various explicit determinants upon the variation of fertility rates, which is the sole response variable. Multiple fertility measures may be computed based on the data available, but the measure this paper relies on is the United Nation's Age-Sex Adjusted Birth Rate. (Shryock and Siegel, et al; 1980, p483).

The second explicit variable under study is general living standard. Within the constraint of data availability, I choose the Indirect Age-sex Standardized Death Rate (DETHR) as a measure of the general living standard. It is seen as a measure of comprehensive quality of life rather than one particular aspect of it. This particular measure considers the age and sex structure of the underlying population, and at the same time generates a single index for each area. It is based on the indirect standardization formula as summarized by Shryock and Siegel (1980, p.421) with the addition of  $w_a$ , the sex structure adjustment weight.

The third explicit variable is indicator of economic structure. Economic structure may be appropriately characterized by employment structure of a region. The Chinese national census classifies regional economic activities into eleven broad categories. In this study, the thirteen economic variables are subjected to a factor analysis, in order to obtain a few synthetic indicators of a region's economic structure. Three factors are retained from the analysis. These factors may be regarded as indices of the respective region's economic structure. Based on the loadings, it may be stated that the first factor (FACTOR1) measures the general economic strength. A region with a high score on this factor is expected to have less agricultural employment, but more industrial employment in almost all urban economic sectors. Regions having high scores on the second factor (FACTOR2) are supposed to have high geological survey employment, which in this province is mainly in oil exploration. High scores on the third factor (FACTOR3) indicate an especially strong mining sector, which is mainly coal mining.

The fourth explicit variable is educational attainment. They are percentages of the total population holding diplomas from: [1] EDCO: college level institutions, including four year colleges, two-to-three year colleges or equivalents, and professional high schools, [2] EDHI: high school level institutions, including regular high schools and junior-high schools, and [3] EDEL: elementary schools.

Before further analysis, maps were drawn using Arc/Info GIS software for variables FERTR, DETHR, EDCO, EDHI, EDEL, and three common factors extracted from the thirteen economic variables (see Appendix, Map 2 through 9). Areal units are classified according to the four quartiles of each variable. Inspection reveals that major urban areas consistently stay in either the first or the fourth quartile with essentially no relationship to their surrounding regions. To account for this effect, a binary indicator variable (DM) is introduced, with its

## SPATIAL STATISTICAL MODELING

value being one if the areal unit is a major urban center and zero otherwise, increasing the total number of explanatory variables to nine.

Following a linearity examination, we proceed to fit a linear regression model. The first step is to select the right explanatory variables. The PROC REG procedure in SAS is used with the MAXR option (see SAS Institute, Inc., 1989). Adjusted R-square and Milliw's C(p) statistic are used to select the best model among the ones given by MAXR regression procedure (see Daniel and Wood, 1983, Chapter 6). Three educational attainment variables should be included in the model by all standards. The inclusion of DM and DETHR, which only made a slight improvement in the R-square, are debatable. The decision concerning their inclusion has to be based on substantive grounds. Since all three firmly selected variables only provide information on educational development, it is believed that inclusion of the two additional variables should add information that is different in nature. Thus the final linear regression model has five explanatory variables. They are EDCO, EDHI, EDEL, DM, and DETHR.

### Linear Regression Results

Source	DF	Sum of Squares	Mean Square	F Value	Prob>F
Model	5	771.02577	154.20515	34.441	0.0001
Error	124	555.19695	4.47739		
C Total	129	1326.22272			

Root MSE	2.11599	R-square	0.5814
Dep Mean	17.98400	Adj R-sq	0.5645
C.V.	11.76593		

MIDDLE STATES GEOGRAPHER - VOL. 28, 1995

Variable	DF	Parameter Estimate	Standard Error	T for H0: Parameter=0	Prob >  T
INTERCEP	1	31.715507	2.36320283	13.421	0.0001
EDCO	1	-29.840538	8.43103179	-3.539	0.0006
EDHI	1	-9.993002	3.53977444	-2.823	0.0055
EDEL	1	-18.055830	4.36639516	-4.135	0.0001
DM	1	-2.098131	1.63157467	-1.286	0.2009
DETHR	1	-0.167109	0.15227221	-1.097	0.2746

Variable	DF	Standardized Estimate	Variance Inflation
INTERCEP	1	0.00000000	0.00000000
EDCO	1	-0.60095460	8.53930537
EDHI	1	-0.25377528	2.39358564
EDEL	1	-0.36062950	2.25281089
DM	1	-0.22147074	8.78562246
DETHR	1	-0.07568409	1.40878091

After inspection for the potential problems of non-normality, heteroscedasticity, and multicollinearity, the linear regression model is considered acceptable, and the results are reported in the above tables. The overall model seems highly significant. Variables DM and DETHR have regression parameter

## SPATIAL STATISTICAL MODELING

estimates that are not significantly different from zero, but those for all three educational attainment variables, EDCO, EDHI, and EDEL are significant at the 0.05 level. Judging from the standardized parameter estimates, of the three variables that have non-zero impacts upon the response variable, EDCO is the most important factor, followed first by EDEL and then by EDHI. The eventual multiple linear regression model seems to be, with natural parameter estimates:  $FERTR=31.715507-29.840538EDCO-18.055830EDEL-9.993002EDHI +e$

### IV In Search of a Spatial Pattern

Residuals from the linear model represent the part of variation in fertility rates that does not go hand in hand with those of the included socio-economic variables. A natural ensuing question asks whether or not the residuals co-vary with other factors that have not been included in the model. This question may be answered by an examination of the residuals. If significant pattern remains among them, we will have to go on to explain the part that has not been explained by the existing explanatory variables. At this point it seems clear that if we were to accept all assumptions of a linear regression, including that of independence among observations, it would indeed seem that the residuals are random. However, given the strong spatial nature of the observations as well as the subject matter in general, it is only reasonable to ask if the residuals are also random spatially, or if there is a latent spatial pattern. The question leads to a scrutiny of spatial autocorrelation among the residuals. Visual inspection seems to suggest certain degree of positive spatial autocorrelation. Computation based on Griffith (1993, pp24-5) shows that the residuals have a Moran Coefficient (I) of 0.39544 and a Geary Ratio (c) of 0.58742, and both are statistically significant, indicating that there is indeed moderate spatial autocorrelation. In other words, there is indeed a spatial pattern not explained by the linear regression model. Our task, then, is to uncover such a pattern.

To understand the process of parameter estimation and statistical inference about the AR model specified earlier, there needs to be some understanding of the Jacobian term. A Jacobian term in this context serves as a normalizing constant. It is important here because it enables us to derive a joint p.d.f. of the vector  $\mathbf{Y}$  in the original model specified by (1) and (2), on the basis of the joint p.d.f. of  $\mathbf{e}$ , and this derived p.d.f. is the basis for a maximum likelihood estimation of the parameter  $r$  as well as  $b$  and  $s^2$  of  $\mathbf{e}$ . Derivation of the joint p.d.f. of  $\mathbf{Y}$  goes like the following:

Assume all random variables in  $\mathbf{e}=(e_1, \dots, e_n)^T$  have identical probability density function  $f(e_i)$ ,  $i=1, 2, \dots, n$ . ... (9) Since  $\mathbf{e}$  is diagonal, the joint probability density function of  $(e_1, \dots, e_n)$  is  $f_e(e_1, \dots, e_n)= \dots$  (10) Meanwhile,  $\mathbf{Y}$  is a function of  $\mathbf{e}$  (recall (2) and the assumptions that  $\mathbf{C}=\mathbf{I}$ , and  $\mathbf{B}=\mathbf{I}-\mathbf{S}$ ),  $\mathbf{Y}=\boldsymbol{\mu}+\mathbf{B}^{-1}\mathbf{e}$  ... (11) and its reverse is,  $\mathbf{e}=\mathbf{B}(\mathbf{Y}-\boldsymbol{\mu})$  ... (12) Then, using algebraic notation, the p.d.f. of  $\mathbf{Y}$  based on the p.d.f. of  $\mathbf{e}$  is  $f(y_1, \dots, y_n) = |\mathbf{J}|f_e()$ , ... (13) where  $|\mathbf{J}|$  is the Jacobian term for transforming the joint p.d.f. from one of  $\mathbf{e}$  to one of  $\mathbf{Y}$ .

Suppose  $\mathbf{e}$  has a standard normal distribution with no latent spatial autocorrelation. Then the joint p.d.f. of  $\mathbf{e}$  should be of the following form: ... (14) Based on (13), this results in a joint p.d.f. of  $\mathbf{Y}$  as:  $f(y_1, \dots, y_n) = |\mathbf{J}| \dots$  (15) The above is also the effective part of the maximum likelihood function of  $\mathbf{Y}$ , and some spatial autocorrelation information is contained in the Jacobian term  $|\mathbf{J}|$ . Note that here we see each areal unit as a variable instead of an observation. In other words the single univariate sample of  $n$  is viewed here as a multivariate sample of size one with  $n$  variables.

In this study, the ordinary least squares method is inappropriate because the estimation equation for  $r$  is non-linear by construction. Maximum likelihood seems to be the only appropriate method. However, as Riply (1990) points out, precisely because of this non-linearity, analytical solutions to the equations are not implementable. The most problematic complication is the excessive numerical intensity involved, which is introduced mainly by the Jacobian term. Various scholars have attempted to simplify the procedure mainly by simplifying the Jacobian term (e.g., Ord, 1975, Gasim, 1988). But the most practical method so far has been proposed by Griffith (1988, 1992, 1993). For irregular lattices, he suggests that the Jacobian term may be approximated quite accurately by the following equation: ... (16) where  $a_1, a_2, d_1$  and  $d_2$  are parameters that are functions of the eigenvalues of  $W$ . After examining a number of empirical cases, Griffith concludes that these four parameters display remarkable consistency over geographical configurations and numbers of areal units, and the values should be generally in the neighborhoods of  $a_1=0.22, a_2=0.12, d_1=1.75,$  and  $d_2=1.05$  ... (17) This finding leads to the generalized Jacobian approximation equation  $J=0.22\ln(1.75)+0.12\ln(1.05)-0.22\ln(1.75+r)-0.12\ln(1.05-r)$  ... (18)

In this study, an AR model is considered more appropriate because it is reasonable to assume that a regions' fertility rate is influenced by its determinants and surrounding region's fertility rates, rather than that the fertility residuals are influenced by the surrounding region's residuals. The necessary SAS programs for model estimation are provided in Griffith (1993, pp76-8). The explanatory variables are EDCO, EDHI, EDEL, DM and DETHR. The response variable is FERTR. Final results of the estimation of the AR model are

Source	DF	Sum of Squares	Mean Square
Regression	7	46736.230166	6676.604309
Residual	123	377.338534	3.067793
Uncorrected Total	130	47113.568700	
(Corrected Total)	129	1440.652591	

SPATIAL STATISTICAL MODELING

Parameter	Estimate	Asymptotic		Confidence Interval	
		Std. Error	95% Lower	Upper	
RHO	0.58900906	0.0851184910	0.420520878	0.757497232	
B0	34.68613633	4.6845827856	25.413218047	43.959054617	
B1	-	31.44874314	6.7308356110	-44.772125706	-18.125360577
B2	-1.96976518	3.0303019580	-7.968109905	4.028579542	
B3	-7.61638378	3.7393701808	-15.018296759	-0.214470801	
B4	-1.43805071	1.2979453169	-4.007274354	1.131172937	
B5	-0.10732190	0.1211423307	-0.347117631	0.132473822	

As the results show, RHO, B0, B1, B3, are significantly different from zero. RHO is about 0.6, which indicates a moderate to strong degree of positive spatial autocorrelation, confirming previous conclusions. When spatial autocorrelation is considered, only two of the five original explanatory variables, plus the sample mean, provided significant influence on the realization of the response variable in an area. They are EDCO and DM. The eventual model is  $(\text{FERTR})=34.69-31.45*(\text{EDCO})- 7.62*(\text{DM}) + 0.589*W*(\text{FERTR}) + e \dots (19)$



V Interpretation and Conclusions

Two sets of calculation results demand our understanding in substantive terms, one from the linear regression model, the other from spatial autocorrelation analysis. The linear regression model is highly significant ( $\text{Prob} > F = 0.0001$ ), with a fit that is very good among social science studies ( $\text{adj. } R\text{-sq.} = 0.5645$ ). Generally speaking, the results are in support of the conventional wisdom that the factors under the definition of modernization drive down fertility rates, at least in this particular province around 1990.

The most outstanding feature of the linear regression model is the overwhelming importance of educational attainment variables, especially the percentage of college graduates among a county's population. Each of the three exhibits a negative relationship with fertility rates, when the other two are included in the model but held constant. That is evidence that education is indeed a crucial factor in fertility reduction. In this case college education is more important than elementary school education, which is more important than high school education.

Population policies are issued by central government and carried out through its administrative hierarchy, of which the counties are the most critical and active nodes. Since there has not been any evidence of discriminatory implementation in this part of the country, we will have to assume that one round of implementation effort takes effect throughout the region at about the same time. In other words, we may expect the neighboring counties to show a simultaneous response in the form of a fertility rate change. For each individual county, this effect is originally mixed with those of explicit socio-economic factors. However, this latter part is filtered out by the linear regression, and what we detect in the spatial model should be purely the effect of the governmental policies. The spatial statistical model defined by equations (19) shows a pattern that takes into account the spatial effects of such policies. It says that when a county's college graduate percentage is held constant, and when it has been classified as one of the major urban areas, its fertility rate is higher when its surrounding counties have high fertility rates. This confirms that, indeed, governmental population policies are effective.

Bibliography

Cliff, A. and K. Ord, 1981, Spatial Processes, London: Pion

Daniel, C., and F.S. Wood, 1983, Fitting Equations to Data: Computer Analysis of Multifactor Data, New York: John Wiley and Sons

Feng, H.M., 1993, The Contextual Determinants of Contemporary Chinese Fertility Transition, unpublished manuscript, Department of Geography, Syracuse University, Syracuse, NY

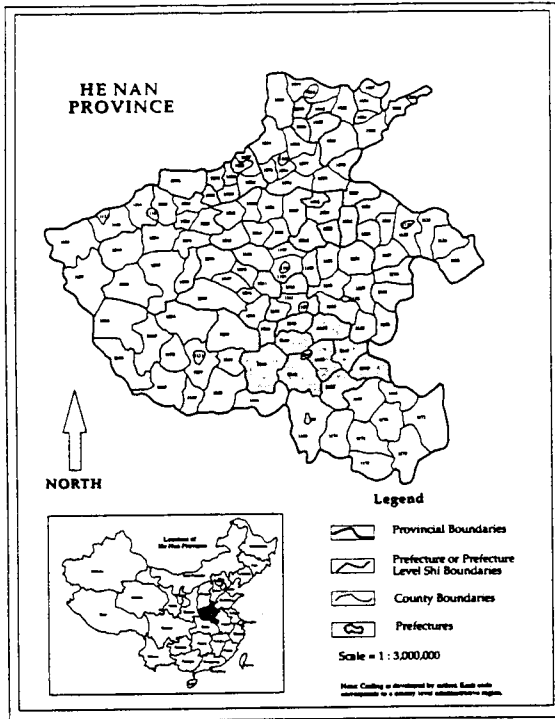
Gasim, A., 1988, "First-order autoregressive models: a method for obtaining eigenvalues for weighting matrices", Journal of Statistical Planning and Inference, Vol.18, pp.391-98.

Griffith, D.A. and C. Amrhein, 1991, Statistical Analysis for Geographers, Englewood, NJ: Prentice Hall

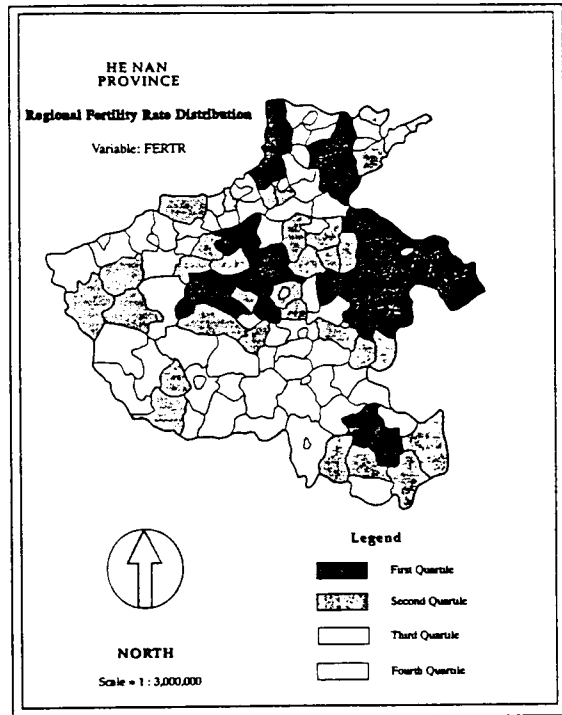
## SPATIAL STATISTICAL MODELING

- Griffith, D.A., 1987, Spatial Autocorrelation: A Primer, Washington, D.C. : Association of American Geographers
- Griffith, D.A., 1988, "Estimating Spatial Autoregressive Model Parameters with Commercial Statistical Packages", Geographical Analysis, Vol.20, No.2, pp.176-86.
- Griffith, D.A., 1992, "Simplifying the normalizing factor in spatial autoregressions for irregular lattices", Papers in Regional Science, Vol.71, pp.71-86.
- Griffith, D.A., 1993, Spatial Regression Analysis on the PC: Spatial Statistics Using SAS, Washington, D.C. : Association of American Geographers
- He-Nan Province Census Bureau, 1992, He-Nan Province Population Census Statistics, 1990, Beijing: China Statistics Publishing Company
- Ord, K., 1975, "Estimation methods for models of spatial interaction", Journal of American Statistical Association, Vol.70, pp.120-26
- SAS Institute Inc., 1989, SAS/STAT User's Guide (two volumes), Cary, N.C.: SAS Institute Inc.
- Shryock, H.S., and J.S. Siegel, et al, 1980, The Methods and Materials of Demography, Washington D.C.: U.S.Bureau of Census. Vol.2
- Wolf, A.P., 1986, "The preeminent role of government intervention in China's family revolution", Population and Development Review, Vol.12, pp.106-16.

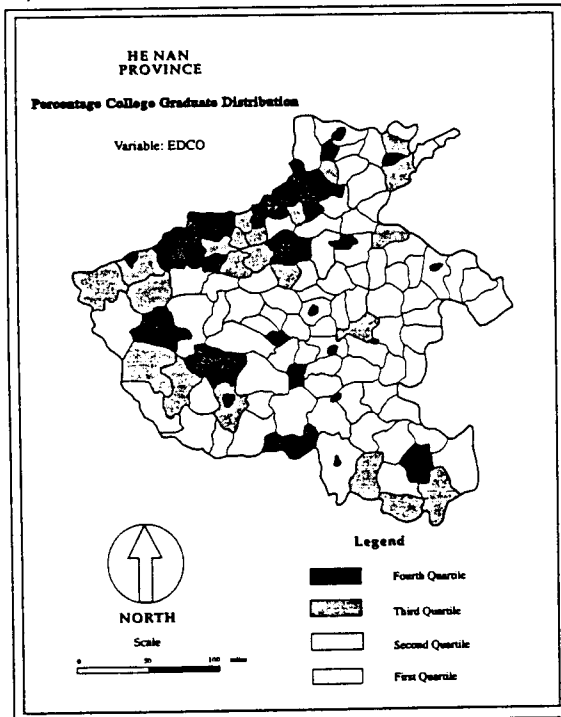
map 1



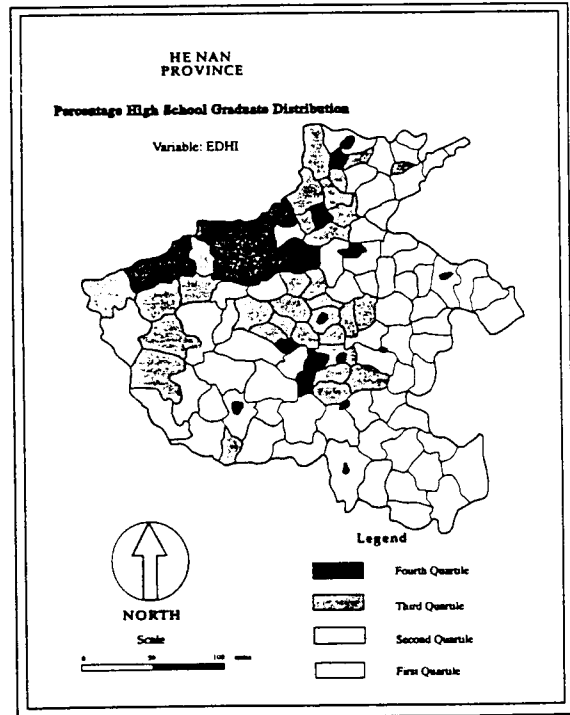
map 2



map 3

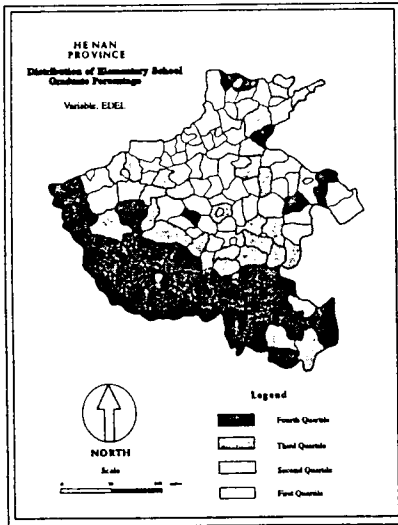


map 4

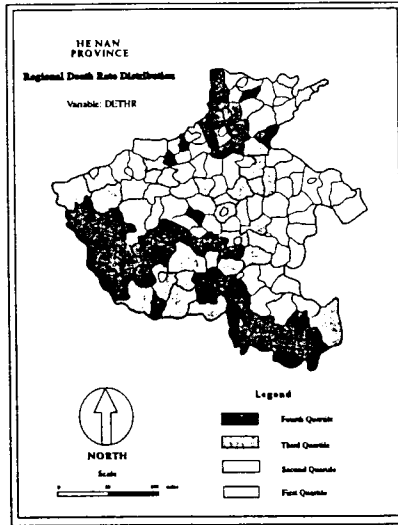


# SPATIAL STATISTICAL MODELING

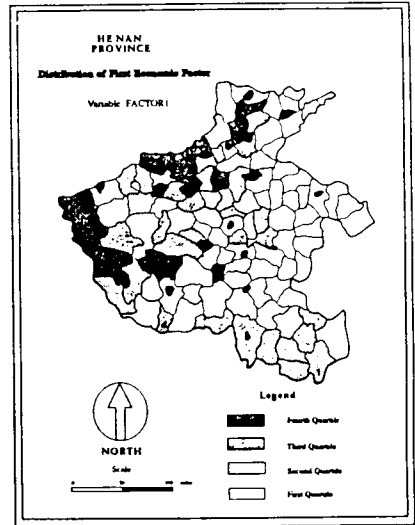
map 5



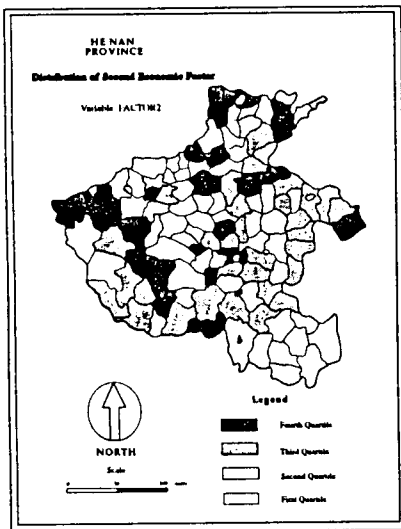
map 6



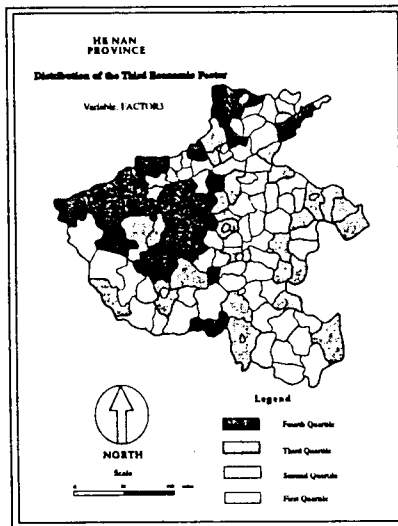
map 7



map 8



map 9



map 10

