

Spatial Contextual Classification and Prediction Models for Mining Geospatial Data

Shashi Shekhar, *Senior Member, IEEE*, Paul R. Schrater, Ranga R. Vatsavai, Weili Wu, *Member, IEEE*, and Sanjay Chawla

Abstract—Modeling spatial context (e.g., autocorrelation) is a key challenge in classification problems that arise in geospatial domains. Markov random fields (MRF) is a popular model for incorporating spatial context into image segmentation and land-use classification problems. The spatial autoregression (SAR) model, which is an extension of the classical regression model for incorporating spatial dependence, is popular for prediction and classification of spatial data in regional economics, natural resources, and ecological studies. There is little literature comparing these alternative approaches to facilitate the exchange of ideas (e.g., solution procedures). We argue that the SAR model makes more restrictive assumptions about the distribution of feature values and class boundaries than MRF. The relationship between SAR and MRF is analogous to the relationship between regression and Bayesian classifiers. This paper provides comparisons between the two models using a probabilistic and an experimental framework.

Index Terms—Markov random fields (MRF), spatial autoregression (SAR), spatial context, spatial data mining.

I. INTRODUCTION

SPATIAL databases (e.g., remote sensing imagery, maps, census data) are an important subclass of multimedia databases for several reasons. First, the industry-wide structured query language multimedia standard (SQL/MM) [20] includes spatial data types along with traditional image, audio, and video data types. Second, spatial concepts and techniques are often crucial in the indexing and retrieval of image and video databases. Finally, according to several estimates, spatial data constitutes almost 80% of all digital data including multimedia data.

Widespread use of spatial databases [28], [29] is leading to an increasing interest in mining interesting and useful but implicit spatial patterns [10], [14], [19], [26]. Traditional data-mining algorithms [1] often make assumptions (e.g., independent, identical distributions) which violate Tobler's first law of geography: *Everything is related to everything else but nearby things are more related than distant things* [31]. In other words, the values

of attributes of nearby spatial objects tend to systematically affect each other. In spatial statistics, an area within statistics devoted to the analysis of spatial data, this is called *spatial autocorrelation* [7]. Knowledge discovery techniques which ignore spatial autocorrelation typically perform poorly in the presence of spatial data. Often the spatial dependencies arise due to the inherent characteristics of the phenomena under study; but, in particular, they arise due to the fact that the spatial resolution of imaging sensors are finer than the size of the object being observed. For example, remote sensing satellites have resolutions ranging from 30 m (e.g., NASA's enhanced thematic mapper of the Landsat 7 satellite) to 1 m (e.g., the IKONOS satellite from SpaceImaging), while the objects under study (e.g., urban, forest, water) are often much larger than 30 m. As a result, per-pixel-based classifiers, which do not take spatial context into account, often produce classified images with *salt and pepper* noise. These classifiers also suffer in terms of classification accuracy.

There are two major approaches for incorporating spatial dependence into classification/prediction models: 1) spatial autoregression (SAR) models [2], [15]–[17], [23], [24] and 2) Markov random field (MRF) models [5], [6], [9], [13], [18], [30], [32]. Here, we want to make a note regarding the terms *spatial dependence* and *spatial context*. These words originated in two different communities. Natural resource analysts and statisticians use *spatial dependence* to refer to *spatial autocorrelation* and the image processing community uses spatial context to mean the same thing. We use *spatial context*, *spatial dependence*, and *spatial autocorrelation* interchangeably to relate to readers of both communities. We also use *classification* and *prediction* interchangeably. Natural resource scientists, ecologists, and economists have incorporated spatial dependence in spatial data analysis by incorporating spatial autocorrelation into the logistic regression models (SAR models). The SAR model states that the class label of a location is partially dependent on the class labels of nearby locations and partially dependent on the feature values. SAR tends to provide better models than logistic regression in terms of achieving higher confidence (R^2). Similarly, the MRF is a popular model for incorporating spatial context into image segmentation and land-use classification problems. Over the last decade, several researchers [13], [30], [32] have exploited spatial context in classification using MRF to obtain higher accuracies over their counterparts (i.e., noncontextual classifiers). MRF provides a uniform framework for integrating spatial context and deriving the probability distribution of interacting objects.

There is little literature comparing alternative models for capturing spatial context, hampering the exchange of ideas across communities. For example, solution procedures [17] for SAR

Manuscript received April 18, 2001; revised February 26, 2002. This work was supported in part by the Army High Performance Computing Research Center under the auspices of the Department of the Army, Army Research Laboratory Cooperative Agreement DAAH04-95-2-0003/Contract DAAH04-95-C-0008. The associate editor coordinating the review of this paper and approving it for publication was Dr. Sankar Basu.

S. Shekhar, R. R. Vatsavai, and W. Wu are with the Department of Computer Science, University of Minnesota, Minneapolis, MN 55455 USA (e-mail: vatsavai@cs.umn.edu; wuw@cs.umn.edu).

P. R. Schrater is with the Department of Psychology, University of Minnesota, Minneapolis, MN 55455 USA (e-mail: schrater@eye.psych.umn.edu).

S. Chawla is with the Vignette Corporation, Boston, MA 02167 USA (e-mail: schawla@vignette.com).

Publisher Item Identifier S 1520-9210(02)04864-2.

tend to be computationally expensive just like the earlier stochastic relaxation approaches [9] for MRF despite optimizations such as sparse-matrix techniques [23], [24]. Recently, new solution procedures, (e.g., graph cuts [5]), have been proposed for MRF. An understanding of the relationship between MRF and SAR may facilitate the development of new solution procedures for SAR. It may also likely lead to cross fertilization of other advances across the two communities.

We compare the SAR and MRF models in this paper using a common probabilistic framework. SAR and MRF use identical models of spatial contexts for spatial locations. However, SAR makes more restrictive assumptions about the probability distributions of feature values as well as the class boundaries. We show that the SAR assumption of the conditional probability of a feature value given a class label means that SAR belongs to the exponential family of models, (e.g., Gaussian, binomial). In contrast, MRF models can work with many other probability distributions. SAR also assumes the linear separability of classes in a transformed feature space resulting from a spatial smoothing of feature values based on autocorrelation parameters. MRF can be used with nonlinear class boundaries. Readers familiar with classification models which ignore spatial context may find the following analogy helpful. The relationship between SAR and MRF is similar to the relationship between logistic regression and Bayesian classifiers.

The rest of the paper is organized as follows. In Section I-A, we introduce a motivating example which will be used throughout the paper. In Section I-B, we formally define the location prediction problem. Section II presents a comparison of classical approaches that do not consider spatial context, namely, logistic regression and Bayesian classifiers. In Section III, we present two modern approaches that model spatial context, namely, SAR [15] and MRF. In Section IV, we compare and contrast the SAR and MRF models in a common probabilistic framework and provide experimental results. Finally, Section V provides conclusions and future research directions.

This paper focuses on a comparison of SAR and MRF. Comparisons of other models of spatial context, and evaluation and translation of new solution procedures for MRF (e.g., graph cuts) to new solution procedures for SAR are beyond the scope of this paper. We plan to address these issues in future work.

A. An Illustrative Application Domain

First, we introduce an example which will be used throughout this paper to illustrate the different concepts in spatial data mining. We are given data about two wetlands, namely, Darr and Stubble, on the shores of Lake Erie in Ohio, in order to predict the spatial distribution of a marsh-breeding bird, the red-winged blackbird (*Agelaius phoeniceus*) [21], [22]. The data was collected from April to June in two successive years: 1995 and 1996.

A uniform grid was imposed on the two wetlands and different types of measurements were recorded at each cell or pixel. In total, the values of seven attributes were recorded at each cell. Domain knowledge is crucial in deciding which attributes are important and which are not. For example, *vegetation durability* was chosen over *vegetation species* because specialized knowledge about the bird-nesting habits of the red-winged blackbird suggested that the choice of nest location is more dependent on

plant structure and plant resistance to wind and wave action than on the plant species.

An important goal is to build a model for predicting the location of bird nests in the wetlands. Typically, the model is built using a portion of the data, called the *learning data* or *training data*, and then tested on the remainder of the data, called the *testing data*. In this study, we built a model using the 1995 Darr wetland data and then test it on 1995 Stubble wetland data. In the learning data, all of the attributes are used to build the model and in the training data, one value is *hidden* (in our case, the location of the nests). Using knowledge gained from the 1995 Darr data and the value of the independent attributes in the test data, we want to predict the location of the nests in the 1995 Stubble data.

In this paper, we focus on three independent attributes: 1) *vegetation durability (Veg)*; 2) *distance to open water (DOW)*; and 3) *water depth (WD)*. The significance of these three variables was established using classical statistical analysis [22]. The spatial distribution of these variables and the actual nest locations for the Darr wetland in 1995 are shown in Fig. 1. These maps illustrate the following two important properties inherent in spatial data. The value of attributes which are referenced by spatial location tend to vary gradually over space. While this may seem obvious, classical data-mining techniques, either explicitly or implicitly, assume that the data is *independently* generated. For example, the maps in Fig. 2 show the spatial distribution of attributes if they were independently generated. Previous studies have evaluated classical data-mining techniques, such as logistic regression [22], neural networks (NNs) [21], decision trees, and classification rules, to build prediction models for bird-nesting locations. Logistic regression was used because the dependent variable is binary (nest/no-nest) and the logistic function “squashes” the real line onto the unit-interval. The values in the unit-interval can then be interpreted as probabilities. These studies concluded that, with the use of logistic regression, the nests could be classified at a rate 24% better than random [21]. In general, logistic regression and NN models have performed better than decision trees and classification rules on this dataset. The fact that classical data-mining techniques ignore spatial autocorrelation and spatial heterogeneity in the model-building process is one reason why these techniques do a poor job. A second, more subtle, but equally important reason is related to the choice of the objective function to measure classification accuracy. For a two-class problem, the standard way to measure classification accuracy is to calculate the percentage of correctly classified objects. This measure may not be the most suitable in a spatial context. *Spatial accuracy* (i.e., how far the predictions are from the actuals) is as important in this application domain due to the effects of discretizations of a continuous wetland into discrete pixels, as shown in Fig. 3. Fig. 3(a) shows the actual locations of nests and Fig. 3(b) shows the pixels with actual nests. Note the loss of information during the discretization of continuous space into pixels. Many nest locations barely fall within the pixels labeled “A” and are quite close to other blank pixels, which represent “no-nest.” Now let us consider the two predictions shown in Fig. 3(c) and (d). Domain scientists prefer the prediction in Fig. 3(d) over the one in Fig. 3(c), since predicted nest locations are closer on average to some actual nest locations. The classification accuracy measure cannot distinguish

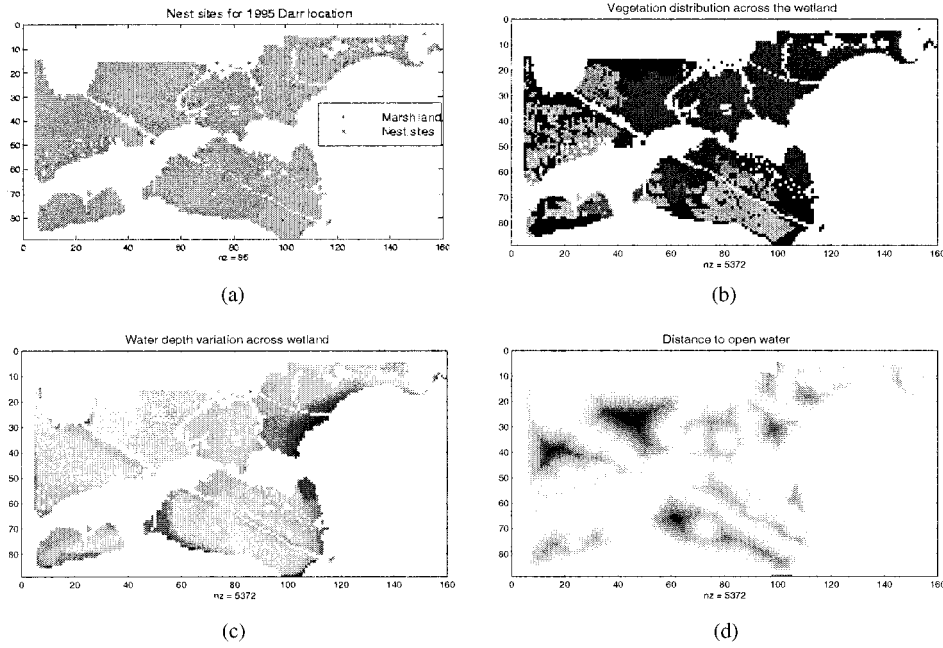


Fig. 1. (a) Learning dataset: the geometry of the Darr wetland and the locations of the nests; (b) the spatial distribution of *vegetation durability* over the marshland; (c) the spatial distribution of *water depth*; and (d) the spatial distribution of *distance to open water*.

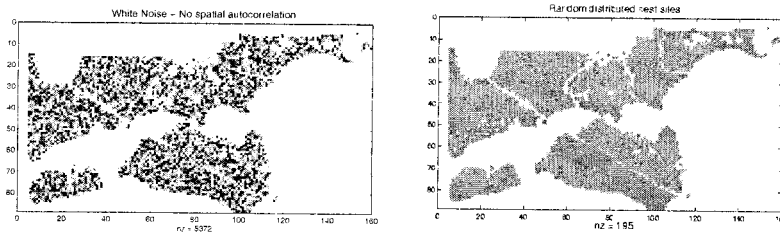


Fig. 2. Spatial distribution satisfying random distribution assumptions of classical regression.

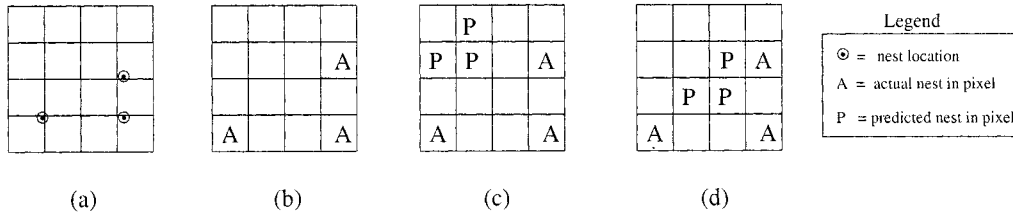


Fig. 3. Example showing different predictions: (a) the actual locations of nests; (b) pixels with actual nests; (c) locations predicted by one model; and (d) locations predicted by another model. Prediction (d) is spatially more accurate than (c).

between Fig. 3(c) and Fig. 3(d), and a measure of spatial accuracy is needed to capture this preference.

B. Location Prediction: Problem Formulation

The location prediction problem is a generalization of the nest location prediction problem. It captures the essential properties of similar problems from other domains including crime prevention and environmental management. The problem is formally defined as follows:

Given:

- 1) a spatial framework S consisting of sites $\{s_1, \dots, s_n\}$ for an underlying geographic space G ;
- 2) a collection X of explanatory functions $f_{X_k}: S \rightarrow R^k, k = 1, \dots, K$. R^k is the range of possible values for the explana-

tory functions. Let $X = [1, X]$, which also includes a constant vector along with explanatory functions;

- 3) a dependent class variable $f_C: S \rightarrow C = \{c_1, \dots, c_M\}$;
- 4) a value for parameter α , relative importance of spatial accuracy.

Find: Classification model: $\hat{f}_C: R^1 \times \dots \times R^k \rightarrow C$.

Objective: Maximize similarity $(\text{map}_{s_i \in S}(\hat{f}_C(f_{X_1}, \dots, f_{X_k})), \text{map}(f_C)) = (1 - \alpha)$; $\text{classification_accuracy}(f_C, \hat{f}_C) + (\alpha)$; $\text{spatial_accuracy}((\hat{f}_C, f_C)$.

Constraints:

- 1) Geographic space S is a multidimensional Euclidean space.¹

¹The entire surface of the earth cannot be modeled as a Euclidean space but locally the approximation holds true.

- 2) The values of the explanatory functions f_{X_1}, \dots, f_{X_k} and the dependent class variable f_C may not be independent with respect to the corresponding values of nearby spatial sites (i.e., spatial autocorrelation exists).
- 3) The domain R^k of the explanatory functions is the one-dimensional (1-D) domain of real numbers.
- 4) The domain of dependent variable $C = \{0, 1\}$.

The above formulation highlights two important aspects of location prediction. It explicitly indicates that: 1) the data samples may exhibit spatial autocorrelation and 2) an objective function (i.e., a map similarity measure) is a combination of classification accuracy and spatial accuracy. The *similarity* between the dependent variable f_C and the predicted variable \hat{f}_C is a combination of the “traditional classification” accuracy and representation-dependent “spatial classification” accuracy. The regularization term α controls the degree of importance of *spatial accuracy* and is typically domain dependent. As $\alpha \rightarrow 0$, the map similarity measure approaches the traditional classification accuracy measure. Intuitively, α captures the spatial autocorrelation present in spatial data.

The study of the nesting locations of red-winged black birds [21], [22] is an instance of the location prediction problem. The underlying spatial framework is the collection of $5 \text{ m} \times 5 \text{ m}$ pixels in the grid imposed on the marshes. Examples of the explanatory variables include water depth, vegetation durability index, and distance to open water, and examples of dependent variables include nest locations. The explanatory and dependent variables exhibit spatial autocorrelation (e.g., gradual variation over space, as shown in Fig. 1). Domain scientists prefer spatially accurate predictions which are closer to actual nests (i.e., $\alpha > 0$).

II. CLASSIFICATION WITHOUT SPATIAL DEPENDENCE

In this section, we briefly review two major statistical techniques that have been commonly used in the classification problem: 1) logistic regression and 2) Bayesian classifiers. These models do not consider spatial dependence. Readers familiar with these two models will find it easier to understand the comparison between SAR and MRF presented later.

A. Logistic Regression Modeling

Logistic regression decomposes \hat{f}_C into two parts, namely, linear regression and logistic transformation. Given an n -vector y of observations and an $n \times m$ matrix X of explanatory data, classical linear regression models the relationship between y and X as

$$y = X\beta + \epsilon$$

where $\beta = (\beta_0, \dots, \beta_m)^T$. The standard assumption on the error vector ϵ is that each component is generated from an independent, identical, zero-mean normal distribution (i.e., $\epsilon_i \sim N(0, \sigma^2)$).

When the dependent variable is binary, as is the case in the “bird-nest” example, the model is transformed via the logistic function and the dependent variable is interpreted as the probability of finding a nest at a given location. Thus, $\Pr(c_i | y) =$

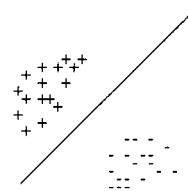


Fig. 4. Two-dimensional feature space, with two classes (+: nest, -: no-nest) that can be separated by a linear surface

$(e^y / (1 + e^y))$. This transformed model is referred to as *logistic regression* [2].

The fundamental limitation of classical regression modeling is that it assumes that the sample observations are independently generated. This may not be true in the case of spatial data. As we have shown in our example application, the explanatory and independent variables show a moderate to high degree of spatial autocorrelation (see Fig. 1). The inappropriateness of the independence assumption shows up in the residual errors, the ϵ_i s. When the samples are spatially related, the residual errors reveal a systematic variation over space (i.e., they exhibit high spatial autocorrelation). This is a clear indication that the model was unable to capture the spatial relationships existing in the data. Thus, the model may be a poor fit to the geospatial data. Incidentally, the notion of spatial autocorrelation is similar to that of time autocorrelation in time series analysis but is more difficult to model because of the multidimensional nature of space. A statistic that quantifies spatial autocorrelation is introduced in the SAR model.

The logistic regression finds a discriminant surface, which is a hyperplane in feature space, as shown in Fig. 4. Formally, a logistic-regression-based classifier is equivalent to a perceptron [11], [12], [27], which can only separate linearly separable classes.

B. Bayesian Classifiers

Bayesian classifiers estimate \hat{f}_C using Bayes’ rule and compute the probability of the class labels c_i given the data X as

$$\Pr(c_i | X) = \frac{\Pr(X | c_i)\Pr(c_i)}{\Pr(X)}. \quad (1)$$

In the case of the location prediction problem, where a single class label is predicted for each location, a decision step can assign the most likely class chosen by Bayes’ rule to be the class for a given location. This solution is often referred to as the *maximum a posteriori estimate (MAP)*.

Given a learning dataset, $\Pr(c_i)$ can be computed as a ratio of the number of locations s_j with $f_C(s_j) = c_i$ to the total number of locations in S . $\Pr(X | c_i)$ can also be estimated directly from the data using histograms or a kernel density estimate over the counts of locations s_j in S for different values X of features and different class labels c_i . This estimation requires a large training set if the domains of features f_{X_k} allow a large number of distinct values. A possible approach is that when the joint-probability distribution is too complicated to be directly estimated, a sufficiently large number of samples from the conditional probability distributions can be used to estimate the *statistics* of the

TABLE I
COMPARISON OF LOGISTIC REGRESSION AND BAYESIAN CLASSIFIERS

Criteria	Classifier	Classifier
Input	Logistic Regression	Bayesian
Input	$f_{x_1}, \dots, f_{x_k}, f_c$	$f_{x_1}, \dots, f_{x_k}, f_c$
Intermediate Result	β	$Pr(c_i), Pr(X c_i)$ using kernel esti.
Output	$Pr(c_i X)$ based on β	$Pr(c_i X)$ based on $Pr(c_i)$ and $Pr(X c_i)$
Decision	Select most likely class for a given feature value	Select most likely class for a given feature value
Assumptions		
- $Pr(X c_i)$	Exponential Family	-
- class boundaries	linearly separable in feature space	-
- autocorrelation in class labels	none	none

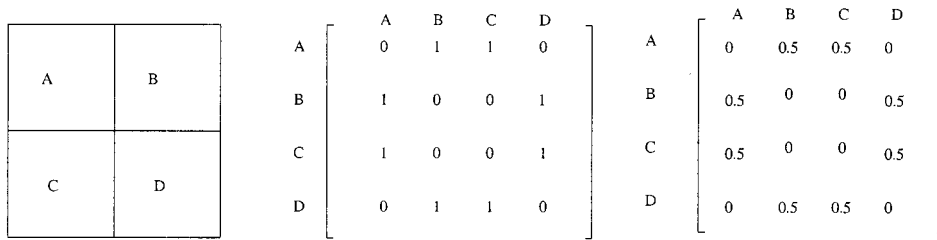


Fig. 5. Spatial framework and its 4-neighborhood contiguity matrix.

full joint-probability distribution.² $Pr(X)$ need not be estimated separately. It can be derived from estimates of $Pr(X|c_i)$ and $Pr(c_i)$. Alternatively, it may be left as unknown, since for any given dataset, $Pr(X)$ is a constant that does not affect the assignment of class labels.

Table I summarizes key properties of logistic-regression-based classifiers and Bayesian classifiers. Both models are applicable to the location prediction problem if spatial autocorrelation is insignificant. However, they differ in many areas. Logistic regression assumes that the $Pr(X|c_i)$ distribution belongs to an exponential family (e.g., binomial, normal), whereas Bayesian classifiers can work with arbitrary distributions. Logistic regression finds a linear classifier specified by β and Bayesian classifier is most effective when classes are not linearly separable in feature space, since it allows nonlinear interaction among features in estimating $Pr(X|c_i)$. Logistic regression can be used with a relatively small training set since it estimates only $(k+1)$ parameters (i.e., β). Bayesian classifiers usually need a larger training set to estimate $Pr(X|c_i)$ due to the potentially large size of the feature space. In many domains, parametric probability distributions (e.g., normal [30], Beta) are used with Bayesian classifiers if large training datasets are not available.

III. MODELING SPATIAL DEPENDENCIES

Several previous studies [13], [30] have shown that modeling of spatial dependency (often called *context*) during the classification process improves overall classification accuracy. Spatial context can be defined by the relationships between spatially adjacent pixels in a small neighborhood. The spatial relationship

²While this approach is very flexible and the workhorse of Bayesian statistics, it is a computationally expensive process. Furthermore, at least for nonstatisticians, it is a nontrivial task to decide what “priors” to choose and what analytic expressions to use for the conditional probability distributions.

among locations in a spatial framework is often modeled via a contiguity matrix. A simple contiguity matrix may represent the neighborhood relationship defined using adjacency, Euclidean distance, etc. Example definitions of neighborhood using adjacency include 4-neighborhood and 8-neighborhood. Given a gridded spatial framework, the 4-neighborhood assumes that a pair of locations influence each other if they share an edge. The 8-neighborhood assumes that a pair of locations influence each other if they share either an edge or a vertex.

Fig. 5(a) shows a gridded spatial framework with four locations, namely, A, B, C, and D. A binary matrix representation of a 4-neighborhood relationship is shown in Fig. 5(b). The row normalized representation of this matrix is called a *contiguity matrix*, as shown in Fig. 5(c). Other contiguity matrices can be designed to model neighborhood relationship based on distance. The essential idea is to specify the pairs of locations that influence each other along with the relative intensity of interaction. More general models of spatial relationships using cliques and hypergraphs are available in [32].

A. Logistic Spatial Autoregression (SAR) Model

Logistic SAR decomposes \hat{f}_C into two parts, namely, SAR and logistic transformation. We first show how spatial dependencies are modeled in the framework of logistic regression analysis. In the SAR model, the spatial dependencies of the error term or the dependent variable, are directly modeled in the regression equation [2]. If the dependent values y_i are related to each other, then the regression equation can be modified as

$$y = \rho W y + X \beta + \epsilon. \quad (2)$$

Here, W is the neighborhood relationship contiguity matrix and ρ is a parameter that reflects the strength of spatial dependencies between the elements of the dependent variable. After

the correction term $\rho W y$ is introduced, the components of the residual error vector ϵ are then assumed to be generated from independent and identical standard normal distributions. As in the case of classical regression, the SAR equation has to be transformed via the logistic function for binary dependent variables.

We refer to this equation as the SAR model. Notice that when $\rho = 0$, this equation collapses to the classical regression model. The benefits of modeling spatial autocorrelation are many. The residual error will have much lower spatial autocorrelation (i.e., systematic variation). With the proper choice of W , the residual error should, at least theoretically, have no systematic variation. If the spatial autocorrelation coefficient is statistically significant, then SAR will quantify the presence of spatial autocorrelation. It will indicate the extent to which variations in the dependent variable (y) are explained by the average of neighboring observation values. Finally, the model will have a better fit, (i.e., a higher R -squared statistic). We compare SAR with linear regression for predicting nest location in Section IV.

A mixed model extends the general linear model by allowing a more flexible specification of the covariance matrix of ϵ . The SAR model can be extended to a mixed model that allows for explanatory variables from neighboring observations [16]. The new model (MSAR) is given by

$$y = \rho W y + X\beta + W X \gamma + \epsilon. \quad (3)$$

The marginal impact of the explanatory variables from the neighboring observations on the dependent variable y can be encoded as a $k * 1$ parameter vector γ .

Solution Procedures: The estimates of ρ and β can be derived using maximum likelihood theory or Bayesian statistics. We have carried out preliminary experiments using the spatial econometrics MATLAB package, which implements a Bayesian approach using sampling-based Markov chain Monte Carlo (MCMC) methods [17]. Without any optimization, likelihood-based estimation would require $O(n^3)$ operations. Recently, [16], [23], and [24] have proposed several efficient techniques to solve SAR. The techniques studied include divide and conquer and sparse matrix algorithms. Improved performance is obtained by using LU decompositions to compute the log-determinant over a grid of values for the parameter ρ by restricting it to $[0, 1]$.

B. Markov Random Field-Based Bayesian Classifiers

MRF-based Bayesian classifiers estimate classification model \hat{f}_C using MRF and Bayes' rule. A set of random variables, the interdependency relationship of which is represented by an undirected graph (i.e., a symmetric neighborhood matrix), is called a *Markov random field (MRF)* [18]. The Markov property specifies that a variable depends only on its neighbors and is independent of all other variables. The location prediction problem can be modeled in this framework by assuming that the class label $l_i = f_C(s_i)$ of different locations, s_i , constitute an MRF. In other words, random variable l_i is independent of l_j if $W(s_i, s_j) = 0$.

The Bayesian rule can be used to predict l_i from feature value vector X and neighborhood class label vector L_i as follows:

$$\Pr(l_i | X, L_i) = \frac{\Pr(X | l_i, L_i) \Pr(l_i | L_i)}{\Pr(X)}. \quad (4)$$

The solution procedure can estimate $\Pr(l_i | L_i)$ from the training data, where L_i denotes a set of labels in the neighborhood of s_i excluding the label at s_i , by examining the ratios of the frequencies of class labels to the total number of locations in the spatial framework. $\Pr(X | l_i, L_i)$ can be estimated using kernel functions from the observed values in the training dataset. For reliable estimates, even larger training datasets are needed relative to those needed for the Bayesian classifiers without spatial context, since we are estimating a more complex distribution. An assumption on $\Pr(X | l_i, L_i)$ may be useful if the training dataset available is not large enough. A common assumption is the uniformity of influence from all neighbors of a location. For computational efficiency, it can be assumed that only local explanatory data $X(s_i)$ and neighborhood label L_i are relevant in predicting class label $l_i = f_C(s_i)$. It is common to assume that all interaction between neighbors is captured via the interaction in the class label variable. Many domains also use specific parametric probability distribution forms, leading to simpler solution procedures. In addition, it is frequently easier to work with a Gibbs distribution specialized by the locally defined MRF through the Hammersley-Clifford theorem [4].

Solution Procedures: Solution procedures for the MRF Bayesian classifier include stochastic relaxation [9], iterated conditional modes [3], dynamic programming [8], highest confidence first [6], and graph cut [5]. We have used the graph-cut method and provided its description in the Appendix .

IV. COMPARISON OF SAR AND MRF BAYESIAN CLASSIFIERS

Both SAR and MRF Bayesian classifiers model spatial context and have been used by different communities for classification problems related to spatial datasets. In this section, we compare these two approaches to modeling spatial context using a probabilistic framework, as well as an experimental framework.

A. Comparison of SAR and MRF Using a Probabilistic Framework

We use a simple probabilistic framework to compare SAR and MRF. We will assume that classes $l_i \in (c_1, c_2, \dots, c_M)$ are discrete and that the class label estimate $\hat{f}_C(s_i)$ for location s_i is a random variable. We also assume that feature values (X) are constant since there is no specified generative model. Model parameters for SAR are assumed to be constant, (i.e., β is a constant vector and ρ is a constant number). Finally, we assume that the spatial framework is a regular grid.

We first note that the basic SAR model can be rewritten as follows:

$$\begin{aligned} y &= X\beta + \rho W y + \epsilon \\ (I - \rho W)y &= X\beta + \epsilon \\ y &= (I - \rho W)^{-1} X\beta + (I - \rho W)^{-1} \epsilon \\ &= (QX)\beta + Q\epsilon \end{aligned} \quad (5)$$

where $Q = (I - \rho W)^{-1}$ and β, ρ are constants (because we are modeling a particular problem). The effect of transforming feature vector X to QX can be viewed as a spatial smoothing operation. The SAR model is similar to the linear logistic model in terms of the transformed feature space. In other words, the SAR

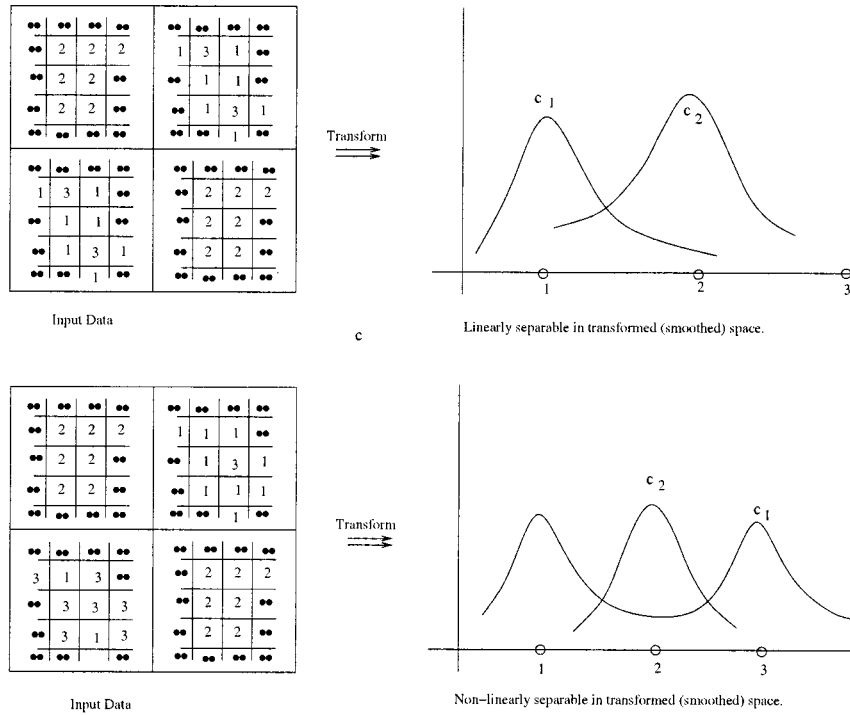


Fig. 6. Spatial datasets with *salt and pepper* spatial patterns.

model assumes the linear separability of classes in transformed feature space.

Fig. 6 shows two datasets with a *salt and pepper* spatial distribution of the feature values. There are two classes, namely, c_1 and c_2 , defined on this feature. Feature values close to 2 map to class c_2 and feature values close to 1 or 3 will map to c_1 . These classes are not linearly separable in the original feature space. Local spatial smoothing can eliminate the *salt and pepper* spatial pattern in the feature values to transform the distribution of the feature values. In the top part of Fig. 6, there are few values of 3 and smoothing revises them close to 1 since most neighbors have values of 1. SAR can perform well with this dataset since classes are linearly separable in the transformed space. However, the bottom part of Fig. 6 shows a different spatial dataset where local smoothing does not make the classes linearly separable. Linear classifiers cannot separate these classes even in the transformed feature space assuming $Q = (I - \rho W)^{-1}$ does not make the classes linearly separable.

Although MRF and SAR classifications have different formulations, they share a common goal, estimating the posterior probability distribution $p(l_i | X)$. However, the posterior for the two models is computed differently with different assumptions. For MRF, the posterior is computed using Bayes' rule. On the other hand, in logistic regression, the posterior distribution is di-

rectly fit to the data. For logistic regression, the probability of the set of labels L is given by

$$\Pr(L | X) = \prod_{i=1}^N p(l_i | X). \quad (6)$$

One important difference between logistic regression and MRF is that logistic regression assumes no dependence on neighboring classes. Given the logistic model, the probability that the binary label takes its first value c_1 at a location s_i is

$$\Pr(l_i | X) = \frac{1}{1 + \exp(-Q_i X \beta)} \quad (7)$$

where the dependence on the neighboring labels exerts itself through the W matrix, and subscript i (in Q_i) denotes the i th row of the matrix Q . Here, we have used the fact that y can be rewritten as in (5).

To find the local relationship between the MRF formulation and the logistic regression formulation (for the two class cases $c_1 = 1$ and $c_2 = 0$), at point s_i , see (8), shown at the bottom of the page, which implies

$$Q_i X \beta = \ln \left(\frac{\Pr(X | l_i = 1, L_i) \Pr(l_i = 1, L_i)}{\Pr(X | l_i = 0, L_i) \Pr(l_i = 0, L_i)} \right) \quad (9)$$

which shows that the spatial dependence is introduced by the W term through Q_i . More importantly, it also shows

$$\begin{aligned} \Pr((l_i = 1) | X, L_i) &= \frac{\Pr(X | l_i = 1, L_i) \Pr(l_i = 1, L_i)}{\Pr(X | l_i = 1, L_i) \Pr(l_i = 1, L_i) + \Pr(X | l_i = 0, L_i) \Pr(l_i = 0, L_i)} \\ &= \frac{1}{1 + \exp(-Q_i X \beta)} \end{aligned} \quad (8)$$

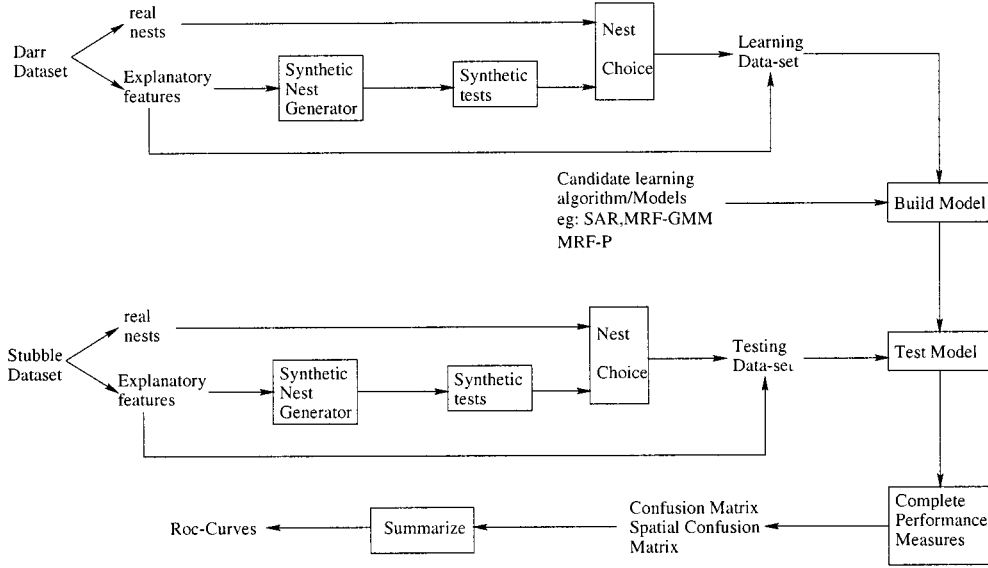


Fig. 7. Experimental method for the evaluation of SAR and MRF.

TABLE II
CONFUSION MATRIX

	Predicted Nest (Present)	Predicted No-nest (Absence)
Actual Nest (Present)	$A_n P_n$	$A_n P_{nn}$
Actual No-nest (Absence)	$A_{nn} P_n$	$A_{nn} P_{nn}$

that, in fitting β , we are trying to simultaneously fit the relative importance of the features and the relative frequency ($\Pr(l_i = 1, L_i)/\Pr(l_i = 0, L_i)$) of the labels. In contrast, in the MRF formulation, we explicitly *model* the relative frequencies in the class prior term. Finally, this relationship shows that we are making distributional assumptions about the class conditional distributions in logistic regression. Logistic regression and logistic SAR models belong to a more general exponential family. The exponential family is given by

$$\Pr(u|v) = e^{A(\theta_v) + B(u, \pi) + \theta_v^T u} \quad (10)$$

where u , and v are location and label, respectively. This exponential family includes many of the common distributions such as Gaussian, binomial, Bernoulli, and Poisson as special cases. The parameters θ_v and π control the form of the distribution. Equation (9) implies that the class conditional distributions are from the exponential family. Moreover, the distributions $\Pr(X|l_i = 1, L_i)$ and $\Pr(X|l_i = 0, L_i)$ are matched in all moments higher than the mean (e.g., covariance, skew, kurtosis, etc.), such that in the difference $\ln(\Pr(X|l_i = 1, L_i)) - \ln(\Pr(X|l_i = 0, L_i))$, the higher order terms cancel out, leaving the linear term ($\theta_v^T u$) in (10) on the left-hand side of (9).

B. Experimental Comparison of SAR and MRF

We carried out experiments to compare the classical regression, SAR, and MRF-based Bayesian classifiers. We compared two families of kernel functions, namely, the Gaussian mixture

model (GMM) and polynomials (P) for MRF-based Bayesian classifiers. We refer to these two families as MRF-GMM and MRF-P, respectively.

The goals of the experiments were:

- 1) to determine whether the real bird habitat datasets follow a Gaussian distribution;
- 2) to evaluate the effect of including a SAR term $\rho W y$ in the logistic regression equation;
- 3) to compare models of spatial context on both real bird habitat datasets and a nonlinear simulated synthetic dataset.

The experimental setup is shown in Fig. 7. The explanatory variables of bird habitat datasets, as described in Section I-A, were used for the learning portion of the experiments. The dependent class variable (i.e., nests) that was used in learning experiments, is of two types, namely, real [see Fig. 1(a)] and synthetic. Synthetic bird datasets were generated using the nonlinear equation (11). All variables in these datasets were defined over a spatial grid of approximately 5000 cells. The 1995 data acquired in the Stubble wetland served as the testing dataset. This data is similar to the learning data except for the spatial locations. We also generated a synthetic dependent class variable Stubble wetlands.

Metrics of Comparison for Classification Accuracy: Consider Boolean vectors $A_n[i] = f_C[s_i]$ representing actual nest locations, and $P_n[i] = \hat{f}_C(s_i)$ representing predicted nest locations and their inverses $A_{nn}[i] = 1 - A_n[i]$ and $P_{nn}[i] = 1 - P_n[i]$. The classification accuracy of various measures for such a binary prediction model is summarized in a matrix, as shown in Table II, using the Boolean vectors.

TABLE III
SPATIAL CONFUSION MATRIX

	Predicted Nest (Present)	Predicted No-nest (Absence)
Actual Nest (Present)	A_nMP_n	A_nMP_{nn}
Actual No-nest (Absence)	$A_{nn}MP_n$	$A_{nn}MP_{nn}$

TABLE IV
DEFINITION OF MEASURES

Measure	Definition	Description
ROC Curve	locus of the pair $(TPR(b), FPR(b))$ for each cut-off probability $TPR = \frac{AnPn}{AnPn + AnPnn}$ $FPR = \frac{AnnPn}{AnnPn + AnnPnn}$	The higher the curve above the straight line $TPR = FPR$, the better the accuracy of the model
Total Error (TE)	$TE = AnPnn + AnnPn$	The lower the value of TE, the better the model
Classification Acc.(CA)	$CA = \frac{AnPnn + AnnPn}{AnPnn + AnnPn + AnnPn + AnPnn}$	
Spatial Acc. Measure	$SAM = A_nMP_n + A_{nn}MP_{nn}$	the higher the value of SAM the better the accuracy of the model
SAM (Normalized)	$SAMN = \frac{A_nMP_n + A_{nn}MP_{nn}}{A_nMP_n + A_{nn}MP_{nn} + A_nMP_{nn} + A_{nn}MP_n}$	
ADNP	$ADNP(A, P) = \frac{1}{K} \sum_{k=1}^K d(A_k, A_k.nearest(P))$	the lower the value of ADNP, the better the model

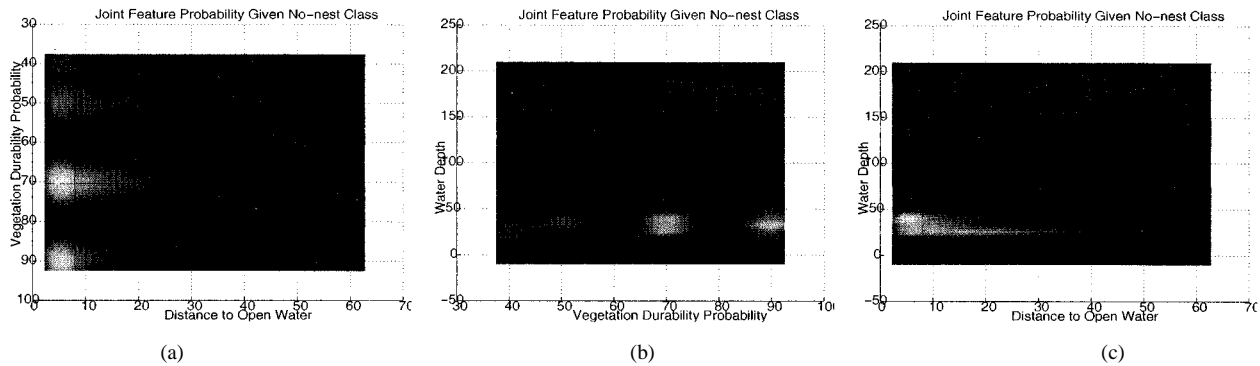


Fig. 8. Joint feature probability distribution for whole datasets: (a) $\Pr(\text{vegetation durability versus distance to open water} \mid \text{no-nest class})$; (b) $\Pr(\text{water depth versus vegetation durability} \mid \text{no-nest class})$; and (c) $\Pr(\text{water depth versus distance to open water} \mid \text{no-nest class})$.

The traditional measure of classification accuracy compares the prediction at location s_i with the actual value at location s_i . This classical measure is not sensitive to the distance between predicted nest and actual nest if the distance is non-zero. We propose new map similarity measures (see Table III). The new map similarity measures compare the prediction at location s_i with the actual value at s_i , as well as the actual values at neighbors of s_i .

In Table III, A_n is an actual nest, A_{nn} is an actual no-nest, P_n is a predicted nest, P_{nn} is a predicted no-nest, and $M = W + I$ is a matrix addition of a contiguity matrix W and an identity matrix I . The spatial accuracy measure (SAM) is defined as $SAM = A_nMP_n + A_{nn}MP_{nn}$.

We summarize various accuracy measures in Table IV.

Average Distance to Nearest Prediction (ADNP) Measure: An orthogonal measure of spatial accuracy is the average distance to nearest prediction (ADNP) from the actual nest sites, which is formulated as $ADNP(A, P)$ in Table IV. A_k represents the actual nest locations, P is the map layer of predicted nest locations, and $A_k \cdots nearest(P)$ denotes the nearest predicted nest location to A_k . K is the number of actual nest sites.

C. Experiments With Real Datasets

We used real datasets from Darr and Stubble wetlands for the results presented in this subsection. The explanatory variables and class labels were described in Section I-A.

1) Characterizing the Probability Distribution ($\Pr(X \mid c_i)$) We analyzed actual wetland datasets to estimate $\Pr(X \mid c_i)$ for the feature values of *Veg*, *DOW*, and *WD*, which were selected as explanatory variables. We explored the statistical probability distribution of each feature given a certain class category (e.g., no-nest class). Fig. 8 illustrates the characteristic probability distribution of each feature value given a nest class for the union of real datasets (learning dataset and testing dataset together). We used the “kernel density estimation toolbox” of MATLAB to fit a smooth function to obtain the observations shown in Fig. 8.

The joint feature probability distribution for a “no-nest” class is displayed in three slices, shown in Fig. 8(a)–(c). Fig. 8(a) shows the slice of the three-dimensional (3-D) joint feature probability of *Veg* versus *DOW* given a “no-nest” class when the other feature (water depth) is fixed at value 38.6. Fig. 8(b) displays the slice of the 3-D joint feature probability of *WD*

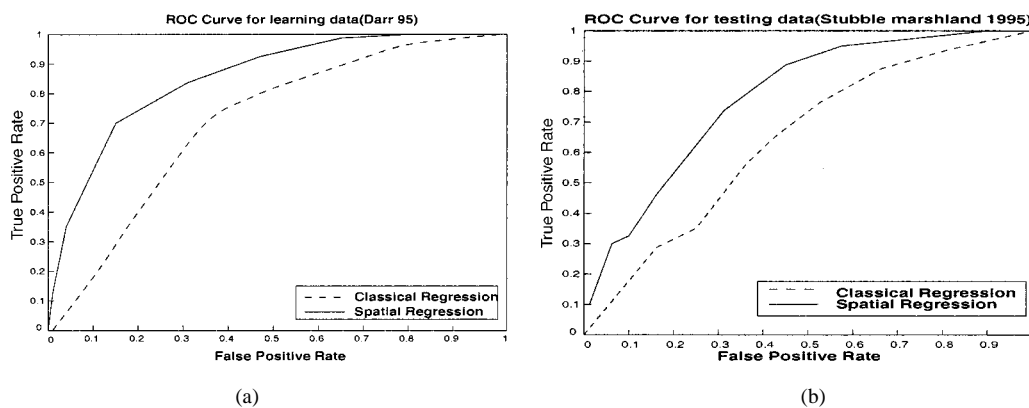


Fig. 9. (a) Comparison of the classical regression model with the SAR model on the Darr learning data. (b) Comparison of the models on the Stubble testing data.

		Learning Data				Test Data			
		Classical Measure		Map Similarity Measure		Classical Measure		Map Similarity Measure	
		Predicted Nest	Predicted No-nest	Predicted Nest	Predicted No-nest	Predicted Nest	Predicted No-nest	Predicted Nest	Predicted No-nest
MRF-P	Actual Nest	42	43	49.06	36.14	9	21	12.42	17.58
	Actual No-nest	96	5191	89.65	5197.3	71	1716	67.58	1719.52
MRF-GMM	Actual Nest	33	52	35.86	49.14	5	25	7.82	22.18
	Actual No-nest	107	5180	96.35	5190.7	73	1714	69.82	1717.18
SAR	Actual Nest	19	66	20.83	64.17	4	26	5.96	24.05
	Actual No-nest	111	5176	109.17	5177.83	76	1711	74.09	1712.8

Fig. 10. Error matrix of real learning and test data.

versus *Veg* given a “no-nest” class when the other feature (*DOW*) is fixed at value 7.97. The slice of the joint feature probability of water depth versus distance to open water given a “no-nest” class when the other feature (vegetation) is fixed at value 70.45 is shown in Fig. 8(c).

It is clear that none of the probability distributions of the real datasets fits a normal distribution, which is a key assumption for regression models (both classical regression and SAR models). However, MRF relaxes this assumption. In the following subsection, we report some experimental results of a comparison of SAR and MRF on both a real bird habitat dataset and a synthetic bird dataset. We used a 11×11 neighborhood matrix in this experimentation.

2) *Comparison of Different Models*: We built a model using the 1995 Darr wetland data and then tested it on the 1995 Stubble wetland data. In the learning data, all of the attributes were used to build the model and in the testing data, one value was hidden (in this case, the location of bird nests). Using the knowledge gained from the 1995 Darr data and the value of the independent attributes in the Stubble test data, we predicted the location of the bird nests in Stubble 1995.

Evaluation of the SAR and Classical Regression Models on Real Datasets: Fig. 9(a) illustrates the ROC curves for SAR and classical regression models built using the real 1995 Darr learning data and Fig. 9(b) displays the ROC curve for the real 1995 Stubble testing data. It is clear that using spatial regression resulted in better predictions at all cutoff probabilities relative to the classical regression model.

Evaluation of the SAR, MRF-GMM, and MRF-P Models: We also compared several spatial contextual models. Fig. 10 illus-

trates learning and testing results for the comparison between SAR, MRF-GMM, and MRF-P kernel density estimation.

The MRF-P model yields better spatial accuracy as well as better classification accuracy than MRF-GMM and SAR in both learning and testing experiments. In this real dataset, the prediction accuracies of MRF-GMM and SAR are very compatible.

We also show maps of the predicted nest locations to visualize the results. Fig. 11(a) shows the actual nest sites for the real learning data (i.e., 1995 Darr bird habitat dataset). Fig. 11(b)–(d) shows the predicted nest locations via the MRF-P kernel density estimation, the MRF-GMM, and the SAR model, respectively. From these maps, we can see that MRF-P yields better prediction. The testing maps are shown in Fig. 11(e)–(h). The ADNP values for each model prediction were also shown in the corresponding figure captions. As can be seen, the SAR predictions are extremely localized, missing actual nests over a large part of the Stubble marsh lands. The SAR predictions in Fig. 11(d) seem to be concentrated on pixels adjacent to water, (i.e., at a small distance to water). This reliance on a single feature is a problem of linear models such as SAR. This is also reflected in the relatively large (two to three times larger than those for MRF models) ADNP values for the predictions from the SAR model.

D. Nonlinear Class Boundary Simulation by Synthetic Bird Datasets

We created a set of synthetic bird datasets based on nonlinear generalization. To generate a set of nonlinear class boundaries, we used the nonlinear equation

$$y = (I - \rho W)^{-1} * (\beta * \cos(X) + c * \text{random}(\epsilon)) \quad (11)$$

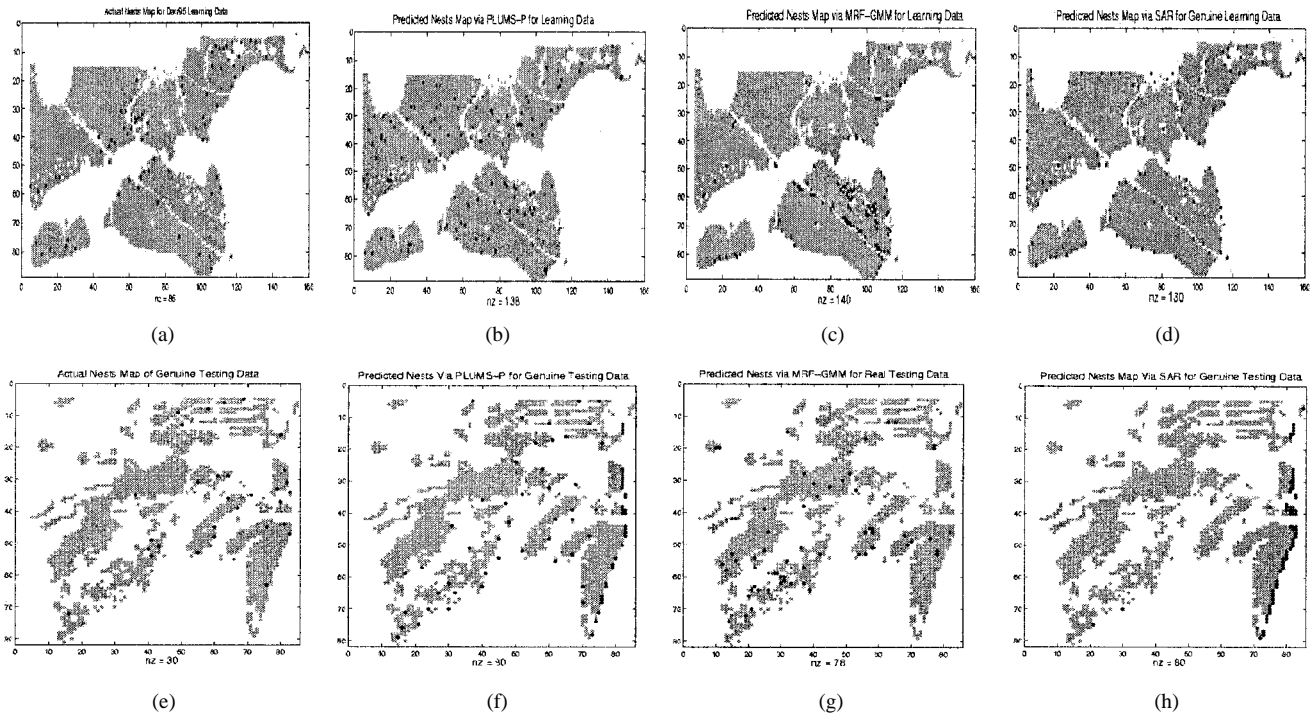


Fig. 11. Predicted nest locations and ADNP for the nonlinear synthetic data. (a) Actual nests on learning. (b) MRF-P learning (ADNP = 1.35). (c) MRF-GMM learning (ADNP = 2.21). (d) SAR learning (ADNP = 6.89). (e) Actual nests on testing. (f) MRF-P testing (ADNP = 2.40). (g) MRF-GMM testing (ADNP = 3.72). (h) SAR testing (ADNP = 5.74).

		Learning Data				Test Data			
		Classical Measure		Map Similarity Measure		Classical Measure		Map Similarity Measure	
		Predicted Nest	Predicted No-nest	Predicted Nest	Predicted No-nest	Predicted Nest	Predicted No-nest	Predicted Nest	Predicted No-nest
MRF-P	Actual Nest	686	866	1007.39	544.61	64	76	72.86	67.14
	Actual No-nest	938	2882	616.61	3203.39	68	1609	59.14	1620.56
MRF-GMM	Actual Nest	522	1030	890.26	661.74	32	108	40.93	99.07
	Actual No-nest	1121	2699	752.74	3067.26	81	1596	72.07	1604.53
SAR	Actual Nest	480	1072	489.36	1062.64	21	119	23.68	116.32
	Actual No-nest	1144	2676	1138.64	2681.36	119	1558	116.32	1560.68

Fig. 12. Error matrix of the nonlinear synthetic learning and testing data generated for Darr95.

where

- X feature values for the independent variables;
- c constant value (we chose 12);
- random(ϵ) random generated error term;
- I identity matrix;
- ρ spatial coefficient (we use $\rho = 0.6$ for both the learning and testing synthetic data);
- W contiguity neighborhood matrix.

To generate synthetic nonlinear learning data, we used the 1995 Darr wetland feature values for X and the contiguity matrix W , and we made the β values the same as SAR's β value. Similarly, using 1995 Stubble wetlands feature values for X , Stubble 95 contiguity matrix W , and the same β values, we generated a synthetic testing dataset on Stubble 1995. For the nonlinear class boundary simulation, we built a model using the nonlinear dataset generated using the Darr wetland and then tested it on the nonlinear synthetic data generated on the 1995 Stubble wetland data. In the learning stage, all of the

feature values of the attributes and spatial dependency are used to build the model and in the testing step, one value is hidden, the location of bird nests. Using the knowledge gained from the learning model and the feature values of the explanatory attributes and spatial dependency in the Stubble test data, we predicted the bird-nest locations in the nonlinear synthetic data on Stubble 1995.

We carried out experiments on these synthetic bird-nesting datasets. Fig. 12 presents accuracy results for MRF-P, MRF-GMM, and SAR models on the nonlinear simulated learning and testing datasets. The confusion matrix shows both classical measure results and map similarity measure results. From Fig. 12, we can easily calculate the total error (TE) of the classical measure and the SAM for the learning model. The TE of MRF-P is $866 + 938 = 1804$, which is significantly less than the TE of MRF-GMM (2151) and SAR (2216). The SAM of MRF-P is $1007.39 + 3203.39 = 4211$, which is greater than those of MRF-GMM (3958) and SAR (3171).

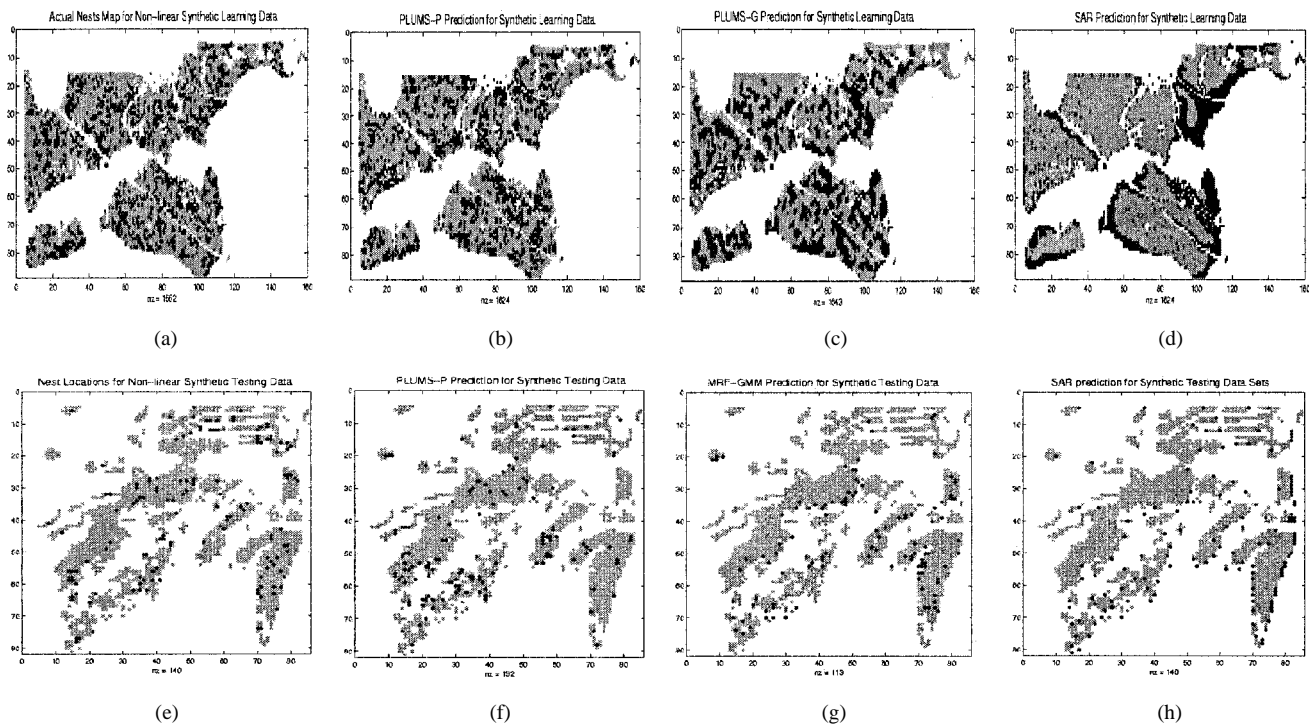


Fig. 13. Predicted nest locations and ADNP for the non-linear synthetic data. (a) Actual nests on learning. (b) MRF-P learning (ADNP = 1.35). (c) MRF-GMM learning (ADNP = 2.21). (d) SAR learning (ADNP = 6.89). (e) Actual nests on testing. (f) MRF-P testing (ADNP = 2.40). (g) MRF-GMM testing (ADNP = 3.72). (h) SAR testing (ADNP = 5.74).

In the nonlinear synthetic dataset, MRF-P achieves better spatial accuracy as well as better classification accuracy than MRF-GMM and SAR in both the learning and testing datasets. The prediction accuracy of MRF-GMM is better than that of SAR in both learning and testing.

We also drew maps of the predicted nest locations to visualize the results (see Fig. 13). Trends were similar to those observed in Fig. 11.

V. CONCLUSION AND FUTURE WORK

In this paper, we have presented two popular classification approaches that model spatial context in the framework of spatial data mining. We have provided theoretical results using a probabilistic framework, as well as experimental results validating the comparison between SAR and MRF. Our paper shows that the SAR model makes more restrictive assumptions about the distribution of features and class shapes (or decision boundaries) than MRF. We also observed an interesting relationship between classical models that do not consider spatial dependence and modern approaches that explicitly model spatial context. The relationship between SAR and MRF is analogous to the relationship between logistic regression and Bayesian classifiers.

In the future, we would like to compare other models that consider spatial context in the classification decision process. We would also like to extend the graph-cut solution procedure for SAR. Finally, we observe that “precision” and “recall” [25] for the learning methods were low (i.e., less than 0.5) for nest predictions, even though classification and spatial accuracies are reasonable. We would like to explore techniques to improve “precision” and/or “recall.”

APPENDIX SOLVING MARKOV RANDOM FIELDS WITH GRAPH PARTITIONING

MRF models generalize Markov chains to multidimensional structures. Since there is no natural order in a multidimensional space, the notion of a transition probability matrix is absent in MRF models.

MRF models have found applications in image processing and spatial statistics, where they have been used to estimate spatially varying quantities like intensity and texture for noisy measurements. Typical images are characterized by piecewise smooth quantities, i.e., they vary smoothly but have sharp jumps (discontinuities) at the boundaries of the homogeneous areas. Because of these discontinuities the least-square approach does not provide an adequate framework for the estimation of these quantities. MRF models provide a mathematical framework to model our *a priori* belief that spatial quantities consist of smooth patches with occasional jumps.

We follow the approach suggested in [5], where it is shown that the MAP estimate of a particular configuration of an MRF can be obtained by solving a suitable min-cut multiway graph partitioning problem. We will formally describe this approach, but first we will illustrate the underlying concept with some examples.

Example 1—A Classification Problem With No Spatial Constraints: Even though MRF models are inherently multidimensional, we will use a simple 1-D example to illustrate the main points. Consider the graph $G = (V, E)$ shown in Fig. 14(a). The node-set V itself consists of two disjoint sets, namely, S and C . The members of S are $\{s_1, s_2, s_3\}$ and the members of C are $\{c_1, c_2\}$. Typically, the $X(s_i)$ s are the feature values

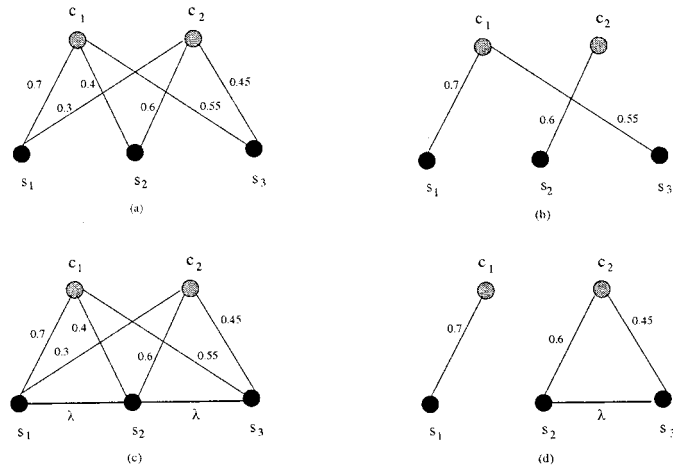


Fig. 14. MRF solution with graph-cut method. (a) Initially each pixel is assigned to both labels with different edge weights. The edge weights correspond to probabilities about assigning each pixel to a different label. (b) A min-cut graph partitioning induces a labeling of the pixel set. Labels which correspond to the maximum probabilities are retained. (c) Spatial autocorrelation is modeled by introducing edges between pixel nodes. (d) A min-cut graph partitioning does not necessarily induce a labeling where the labeling with maximum probabilities are retained. If two neighboring pixels are assigned different labels, then the edge connecting the pixels is added to the cut-set.

at site s_i and the c_i s are the labels, such as *nest* or *no-nest*. There is an edge between each member of the set S and each member of set C . Here, we interpret the edge weights as probabilities. For example, $p_1 = \Pr(X(s_1) = c_1) = 0.7$ and $p_2 = \Pr(X(s_1) = c_2) = 0.3$; $p_1 + p_2 = 1$.

Our goal is to provide a *label* for each location s_i in S using explanatory feature $X(s_i)$. This is done by partitioning the graph into two disjoint sets (not S and C) by removing certain edges, resulting in the following.

- 1) There is a many-to-one mapping from the set S to C . Every element of S must be mapped to one and only one element of C .
- 2) Multiple elements of C cannot belong to a single partition. Thus, there are no edges between elements of C and therefore the number of partitions is equal to the cardinality of C .
- 3) The sum of the weights of the edges removed (the cut-set) is the minimum of all possible cut-sets.

In this example, the cut-set is easily determined. For example, of the two edges connecting each element of S and an element of C , remove the edge with the *smaller* weight. Fig. 14(b) shows the graph with the cut-set removed. Thus, we have just shown that when the weights of the edges are interpreted as probabilities, the min-cut graph partitioning induces a MAP estimate for the pixel labels. We prefer to say that the *min-cut induces a Bayesian classification* on the underlying pixel set. This is because we will use Bayes' theorem to calculate the edge weights of the graphs.

Example 2—Adding Spatial Constraints: In *Example 1*, we did not use any information about the spatial proximity of the pixels relative to each other. We do that now by introducing additional edges in the graph structure.

Consider the graph shown in Fig. 14(c), in which we have added two extra edges (s_1, s_2) and (s_2, s_3) with a weight λ . In this example, we have chosen $\lambda = 0.2$.

Now, if we want to retain the same partitions of the graph as in *Example 1*, then the cut-set has two extra edges, namely, (s_1, s_2) and (s_2, s_3) . Thus, the sum of the weights of the edges in the cut-set W_{C1} is

$$W_{C1} = 0.3 + 0.4 + 0.45 + 2\lambda.$$

However, depending upon λ , the cut-set weight may now not be minimal. For example, if $\lambda = 0.2$, then the weight of the cut-set W_{C2} , consisting of the edges $\{(s_1, c_2), (s_2, c_1), (s_3, c_1), (s_1, s_2)\}$, is

$$W_{C2} = 0.3 + 0.4 + 0.55 + 0.2.$$

Thus, $W_{C2} < W_{C1}$. If two neighboring pixels are assigned to different labels, then the edge between the two neighbors is added to the cut-set. Thus, there is a penalty associated with two neighboring nodes being assigned to different labels every time. Therefore, we can model spatial autocorrelation by adding edges between the pixel nodes of the graph. We can also model spatial heterogeneity by assigning different *weights*, the λ s, to the pixel edges.

Formal Description: Using the terminology introduced in [5], we now formalize the observations made in the previous two examples. Again, consider a graph $G = (V, E)$ with nonnegative edge weights. The set V consists of two types of nodes, namely, *pixels* and *labels*. We will denote the set of pixels as S and the set of labels as C . There are also two types of edges, namely, *n-links* and *l-links*. An *n-link* connects two pixels and an *l-link* connects a pixel with a label. There are no edges between labels. The *l-link* (c_i, s_j) essentially represents the conditional probability $\Pr(l_j = c_i | X(s_j))$.

Definition: A set $K \subset E$ is a *multiway cut* if the label nodes C are completely separated in the graph $G(K) = (V, E - K)$. The sum of the weights of edges in the cut-set K is denoted as $|K|$. A cut-set is a *min cut-set* if its weight is the minimum of all possible cut-sets.

Definition: A cut-set is *feasible* if it induces a many-to-one mapping from S to C and no elements of C can belong to the same set. (From now on, we will only consider feasible cut-sets).

Lemma 1: If a graph G (as defined above) has no *n-links* and the weights on the *l-links* are the *posteriori* probabilities $\Pr(c_i | s_j)$, then the min-cut induces a Bayesian classification on the pixel set S .

Proof: A cut-set K induces a graph in which each pixel is assigned to one and only one label. Thus, every cut-set induces a classification f on the pixel set S . Now

$$|K| = \sum_{s_j \in S} \sum_{c_i \in C, c_i \neq f(s_j)} \Pr(f(s_j) = c_i | X(s_j)).$$

Thus

$$\begin{aligned} \min_f |K| &= \min_f \sum_{s_j \in S} \sum_{c_i \in C, c_i \neq f(s_j)} \Pr(f(s_j) = c_i | X(s_j)) \\ &= \sum_{s_j \in S} \min_f \sum_{c_i \in C, c_i \neq f(s_j)} \Pr(f(s_j) = c_i | X(s_j)). \end{aligned}$$

We can pass the minimum through the first summation because there are no n -links and the cut-sets are feasible. Now, for a given $s_j \in S$

$$\sum_{c_i \in C} \Pr(f(s_j) = c_i) = 1.$$

Therefore

$$\begin{aligned} \sum_{s_j \in S} \min_f \sum_{c_i \in C, c_i \neq f(s_j)} \Pr(f(s_j) = c_i | X(s_j)) \\ = \sum_{s_j \in S} \min_f (1 - \Pr(f(s_j) = c_i)). \end{aligned}$$

The last term is minimized when we choose the maximum probabilities $\Pr(f(s_j) = c_i)$ for each $s_j \in S$. Therefore, $\min |K|$ induces a classifier f which corresponds to the Bayesian classification of the pixel set S , since Bayes' rule was used to determine the edge weights $(s_j, c_i) = \Pr(f(s_j) = c_i)$. The classification f minimizing $|K|$ is chosen as the (\hat{f}_c) solution to the location prediction problem.

Definition: A neighborhood system N of a multiway graph G , as defined above, consists of all unordered pixel pairs $\{s_i, s_j\}$ such that there is an n -link between s_i and s_j . $N(s_i)$ consists of all pixels in G which are n -linked to s_i .

Definition: Let f be the classifier on the pixel set S of a graph G . Then, the energy E associated with f is defined as

$$\begin{aligned} E(f) = \sum_{s_j \in S} \sum_{c_i \in C, c_i \neq f(s_j)} \Pr(f(s_j) = c_i | X(s_j)) \\ + \frac{\lambda}{2} \sum_{s_j \in S} \sum_{s_k \in N(s_j)} (1 - \delta(f(s_j) - f(s_k))) \end{aligned}$$

where δ is the impulse function such that

$$\delta(s_j - s_k) = \begin{cases} 1, & \text{if } s_j = s_k \\ 0, & \text{if } s_j \neq s_k \end{cases}$$

Lemma 2: Let G be a graph, as defined above, where the weights of the l -links are $\Pr(f(s_j) = c_i | X(s_j))$ and the weights of the n -links are λ . Then, a min-cut-set of G induces a classifier f on S , which minimizes the energy function E .

Proof: By construction of the graph G , the weight of the cut-set is E . A min-cut induces an f which minimizes E .

Minimizing E is equivalent to a MAP estimate of the MRF model [5].

How the Edge-Weights of the Graph Are Generated: We use a training set in conjunction with Bayes' theorem to generate the edge weights of the l -links of the graph. In general, the labels of the pixels are not directly observable (that is, what we want to calculate), but we do have an estimate of the "independent" variables, Y . Thus, given a label set C and an observation X at s_j , we can compute the required *posteriori* $\Pr(c_i | X(s_j))$ using Bayes' formulae.

ACKNOWLEDGMENT

The authors would like to thank their collaborator U. Ozesmi for his help and for providing the bird habitat datasets. They would also like to thank J. Lesage for providing the MATLAB

toolbox, and C.-T. Lu and H. Yan for their useful comments. Finally, they would like to thank K. Koffolt, whose comments greatly improved the readability of this paper.

REFERENCES

- [1] R. Agrawal, "Tutorial on database mining," in *Proc. 13th ACM Symp. Principles of Databases Systems*, Minneapolis, MN, 1994, pp. 75–76.
- [2] L. Anselin, *Spatial Econometrics: Methods and Models*. Norwell, MA: Kluwer, 1988.
- [3] J. Besag, "On the statistical analysis of dirty pictures," *J. R. Stat. Soc.*, vol. 48, pp. 259–302, 1986.
- [4] J. E. Besag, "Spatial interaction and statistical analysis of lattice systems," *J. R. Stat. Soc.*, ser. B, vol. 36, pp. 192–236, 1974.
- [5] Y. Boykov, O. Veksler, and R. Zabih, "Fast approximate energy minimization via graph cuts," presented at the Proc. Int. Conf. Computer Vision, Sept. 1999.
- [6] P. B. Chou, P. R. Copper, M. J. Swain, C. M. Brown, and L. E. Wixson, "Probabilistic network inference for cooperative high and low level vision," in *Markov Random Field, Theory and Applications*. New York: Academic, 1993.
- [7] N. A. Cressie, *Statistics for Spatial Data*, Revised ed. New York: Wiley, 1993.
- [8] H. Derin and H. Elliott, "Modeling and segmentation of noisy and textured images using Gibbs random fields," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. PAMI-9, pp. 39–55, 1987.
- [9] S. Geman and D. Geman, "Stochastic relaxation, Gibbs distributions and the Bayesian restoration of images," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. PAMI-6, pp. 721–741, 1984.
- [10] C. Greenman. (2000, Jan.) Turning a map into a cake layer of information. *New York Times* [Online]. Available: <http://www.nytimes.com/library/tech/00/01/circuits/articles/20giss.html>
- [11] S. Haykin, *Neural Networks—A Comprehensive Foundation*. New York: MacMillan, 1994. ISBN-002-352761-7.
- [12] D. W. Hosmer and S. Lemeshow, *Applied Logistic Regression*. New York: Wiley, 1989.
- [13] Y. Jung and P. H. Swain, "Bayesian contextual classification based on modified M -estimates and Markov random fields," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 34, no. 1, pp. 67–75, 1996.
- [14] K. Koperski, J. Adhikary, and J. Han, "Spatial data mining: Progress and challenges," in *Proc. Workshop Research Issues on Data Mining and Knowledge Discovery*, Montreal, Canada, 1996, pp. 1–10.
- [15] J. LeSage, "Regression analysis of spatial data," *J. Regional Anal. Policy*, vol. 27, no. 2, pp. 83–94, 1997.
- [16] J. P. LeSage and R. K. Pace, "Spatial dependence in data mining," in *Data Mining for Scientific and Engineering Applications*. Norwell, MA: Kluwer, 2001.
- [17] J. P. LeSage, "Bayesian estimation of spatial autoregressive models," *Int. Reg. Sci. Rev.*, vol. 20, pp. 113–129, 1997.
- [18] S. Li, "Markov random field modeling," in *Computer Vision*. New York: Springer-Verlag, 1995.
- [19] D. Mark, "Geographical information science: Critical issues in an emerging cross-disciplinary research domain," presented at the NSF Workshop, Feb. 1999.
- [20] J. Melton and A. Eisenberg, "Sql multimedia and application packages (sql/mm)," presented at the SIGMOD Record, vol. 30, Dec. 2001.
- [21] S. Ozesmi and U. Ozesmi, "An artificial neural network approach to spatial habitat modeling with interspecific interaction," in *Ecological Modeling*. Amsterdam, The Netherlands: Elsevier, 1999, vol. 116, pp. 15–31.
- [22] U. Ozesmi and W. Mitsch, "A spatial habitat model for the marsh-breeding red-winged blackbird (*agelaius phoeniceus* L.)," in *Coastal Lake Erie Wetlands Ecological Modeling*. Amsterdam, The Netherlands: Elsevier, 1997, vol. 101, pp. 139–152.
- [23] R. Pace and R. Barry, "Quick computation of regressions with a spatially autoregressive dependent variable," *Geograph. Anal.*, 1997.
- [24] —, "Sparse spatial autoregressions," in *Statistics and Probability Letters*. Amsterdam, The Netherlands: Elsevier, 1997, pp. 291–297.
- [25] B. Ribeiro-Neto and R. Baeza-Yates, *Modern Information Retrieval*. Reading, MA: ACM Press and Addison-Wesley, 1999.
- [26] J. F. Roddick and M. Spiliopoulou, "A bibliography of temporal, spatial and spatio-temporal data mining research," presented at the ACM Special Interest Group on Knowledge Discovery in Data Mining (SIGKDD) Explorations, 1999.
- [27] W. S. Sarle, "Neural networks and statistical models," presented at the 9th Annu. SAS User Group Conf., 1994.

- [28] S. Shekhar, S. Chawla, S. Ravada, A. Fetterer, X. Liu, and C. T. Lu, "Spatial databases: Accomplishments and research needs," *IEEE Trans. Knowledge Data Eng.*, vol. 11, Jan./Feb. 1999.
- [29] S. Shekhar and S. Chawla, *A Tour of Spatial Databases*. Englewood Cliffs, NJ: Prentice-Hall, 2002. ISBN 0-7484-0064-6.
- [30] A. H. Solberg, T. Taxt, and A. K. Jain, "A Markov random field model for classification of multisource satellite imagery," *IEEE Trans. Geosci. Remote Sensing*, vol. 34, no. 1, pp. 100–113, 1996.
- [31] W. R. Tobler, *Cellular Geography, Philosophy in Geography*, W. R. Gale and W. R. Olsson, Eds, Amsterdam, The Netherlands: Reidel, 1979.
- [32] C. E. Warrender and M. F. Augusteijn, "Fusion of image classifications using Bayesian techniques with Markov rand fields," *Int. J. Remote Sens.*, vol. 20, no. 10, pp. 1987–2002, 1999.



Shashi Shekhar (S'86–M'89–SM'96) received the B.Tech degree in computer science from the Indian Institute of Technology, Kanpur, in 1985, and the M.S. degree in business administration and the Ph.D. degree in computer science from the University of California, Berkeley, in 1989.

He is currently a Professor of Computer Science, the University of Minnesota, Minneapolis. His research interests include spatial databases, geographic information systems (GIS), and intelligent transportation systems. He has published over 100

research papers in journals, books, conferences, and workshops. He is a member of the editorial board of *Geo-Informatica: An International Journal on Advances in Computer Science for GIS*.

Dr. Shekhar has served on the editorial board of IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, as well as the IEEE-CS Computer Science and Engineering Practice Board. He is a member of the ACM. He also served as a program Co-chair of the ACM International Workshop on Advances in Geographic Information Systems (1996).



Paul Schrater received the Ph.D. degree in 1999 from the University of Pennsylvania, Philadelphia. His dissertation work involved Bayesian approaches to image motion analysis and their application to understanding human perception.

His current research involves investigating how regularities in natural image and scene statistics can be used to solve problems in human and computer vision, and the application of predictive pattern recognition methods to data-mining problems, including functional MRI imaging data.



Ranga R. Vatsavai received the B.S. degree in applied mathematics from Osmania University, Hyderabad, India, in 1990, and the M.S. degree in statistics from HNBG University, Dehradun, India, in 1996.

He is currently a Research Fellow at the Remote Sensing and Geospatial Analysis Laboratory at the University of Minnesota, Minneapolis. He was a Project Leader of the team that designed and implemented a parallel photogrammetry system for Indian Remote Sensing Satellite (IRS-1C) on a PARAM series of supercomputers at the Center for Development of Advanced Computing (C-DAC), India. He also worked on several remote sensing and GIS research projects funded by FAO, UNDP, and UNEP. His research interests include high-performance signal and image processing algorithms for remote sensing, pattern recognition, spatio-temporal data mining, and semi-structured and spatial databases.



Weili Wu (M'98) received the Ph.D. degree in computer science from the University of Minnesota in 2002.

Her main research interest is in database systems. She has produced a number of research papers in spatial data mining, distributed database systems, and algorithm design. She was also a co-editor of a book on clustering and information retrieval. She will join the Department of Computer Science at the University of Texas at Dallas as Assistant Professor in the Fall of 2002.



Sanjay Chawla received the Ph.D. degree in mathematics from the University of Tennessee, Knoxville, in 1995.

He was a Postdoctoral Associate in the Computer Science Department at the University of Minnesota, Minneapolis. He is currently a Senior Technical Instructor at Vignette Corporation, Boston, MA. His research interests include spatial database management, data mining, geographical information systems (GIS), and optimal control theory.