

Collective Classification in Network Data

Prithviraj Sen, Galileo Namata, Mustafa Bilgic,
Lise Getoor, Brian Gallagher,
and Tina Eliassi-Rad

■ Many real-world applications produce networked data such as the worldwide web (hypertext documents connected through hyperlinks), social networks (such as people connected by friendship links), communication networks (computers connected through communication links), and biological networks (such as protein interaction networks). A recent focus in machine-learning research has been to extend traditional machine-learning classification techniques to classify nodes in such networks. In this article, we provide a brief introduction to this area of research and how it has progressed during the past decade. We introduce four of the most widely used inference algorithms for classifying networked data and empirically compare them on both synthetic and real-world data.

Networks have become ubiquitous. Communication networks, financial transaction networks, networks describing physical systems, and social networks are all becoming increasingly important in our day-to-day life. Often, we are interested in models of how nodes in the network influence each other (for example, who infects whom in an epidemiological network), models for predicting an attribute of interest based on observed attributes of objects in the network (for example, predicting political affiliations based on online purchases and interactions), or we might be interested in identifying important nodes in the network (for example, critical nodes in communication networks). In most of these scenarios, an important step in achieving our final goal is classifying, or labeling, the nodes in the network.

Given a network and a node v in the network, there are three distinct types of correlations that can be utilized to determine the classification or label of v : (1) The correlations between the label of v and the observed attributes of v . (2) The correlations between the label of v and the observed attributes (including observed labels) of nodes in the neighborhood of v . (3) The correlations between the label of v and the unobserved labels of objects in the neighborhood of v . *Collective classification* refers to the combined classification of a set of interlinked objects using all three types of information just described.

Many applications produce data with correlations between labels of intercon-

nected nodes. The simplest types of correlation can be the result of homophily (nodes with similar labels are more likely to be linked) or the result of social influence (nodes that are linked are more likely to have similar labels), but more complex dependencies among labels often exist.

Within the machine-learning community, classification is typically done on each object independently, without taking into account any underlying network that connects the nodes. Collective classification does not fit well into this setting. For instance, in the web page classification problem where web pages are interconnected with hyperlinks and the task is to assign each web page with a label that best indicates its topic, it is common to assume that the labels on interconnected web pages are correlated. Such interconnections occur naturally in data from a variety of applications such as bibliographic data, email networks, and social networks. Traditional classification techniques would ignore the correlations represented by these interconnections and would be hard pressed to produce the classification accuracies possible using a collective classification approach.

Although traditional exact probabilistic inference algorithms such as variable elimination and the junction tree algorithm harbor the potential to perform collective classification, they are practical only when the graph structure of the network satisfies certain conditions. In general, exact inference is known to be

NP-hard and there is no guarantee that real-world network data satisfy the conditions that make exact inference tractable for collective classification. As a consequence, most of the research in collective classification has been devoted to the development of approximate inference algorithms.

In this article we provide an introduction to four popular approximate inference algorithms used for collective classification: iterative classification, Gibbs sampling, loopy belief propagation, and mean-field relaxation labeling. We provide an outline of the basic algorithms by providing pseudocode, explain how one could apply them to real-world data, provide theoretical justifications (if any exist), and discuss issues such as feature construction and various heuristics that may lead to improved classification accuracy. We provide case studies, on both real-world and synthetic data, to demonstrate the strengths and weaknesses of these approaches. All of these algorithms have a rich history of development and application to various problems relating to collective classification, and we provide a brief discussion of this when we examine related work. Collective classification has been an active field of research for the past decade, and as a result, there are numerous other approximate inference algorithms besides the four we describe here. We provide pointers to these works in the related work section. In the next section, we begin by introducing the required notation and define the collective classification problem formally.

Collective Classification: Notation and Problem Definition

Collective classification is a combinatorial problem, in which we are given a set of nodes, $\mathcal{V} = \{V_1, \dots, V_n\}$ and a neighborhood function \mathcal{N} , where $\mathcal{N}_i \subseteq \mathcal{V} \setminus \{V_i\}$, which describes the underlying network structure. Each node in \mathcal{V} is a random variable that can take a value from an appropriate domain. \mathcal{V} is further divided into two sets of nodes: \mathcal{X} , the nodes for which we know the correct values (observed variables), and \mathcal{Y} , the nodes whose values need to be determined. Our task is to label the nodes $Y_i \in \mathcal{Y}$ with one of a small number of labels, $\mathcal{L} = \{L_1, \dots, L_q\}$; we will use the shorthand y_i to denote the label of node Y_i .

We explain the notation further using a web page classification example, motivated by Craven and colleagues (1998), that will serve as a running example throughout the article. Figure 1 shows a network of web pages with hyperlinks. In this example, we will use the words (and phrases) contained in the web pages as node attributes. For brevity, we abbreviate the node attributes, thus, “ST” stands for “student,” “CO” stands for “course,” “CU” stands for “curriculum,” and “AI”

stands for “artificial intelligence.” Each web page is indicated by a box, the corresponding topic of the web page is indicated by an ellipse inside the box, and each word in the web page is represented using a circle inside the box. The observed random variables \mathcal{X} are shaded, whereas the unobserved ones \mathcal{Y} are not. We will assume that the domain of the unobserved label variables, \mathcal{L} , in this case, is a set of two values: “student homepage” (abbreviated to “SH”) and “course homepage” (abbreviated to “CH”). Figure 1 shows a network with two unobserved variables (Y_1 and Y_2), which require prediction, and seven observed variables ($X_3, X_4, X_5, X_6, X_7, X_8$, and X_9). Note that some of the observed variables happen to be labels of web pages (X_6 and X_8) for which we know the correct values. Thus, from the figure, it is easy to see that the web page W_1 , whose unobserved label variable is represented by Y_1 , contains two words “ST” and “CO” and hyperlinks to web pages W_2, W_3 , and W_4 .

As mentioned in the introduction, due to the large body of work done in this area of research, we have a number of approaches for collective classification. At a broad level of abstraction, these approaches can be divided into two distinct types, one in which we use a collection of unnormalized local conditional classifiers and one in which we define the collective classification problem as one global objective function to be optimized. We next describe these two approaches, and for each approach, we describe two approximate inference algorithms. For each topic of discussion, we will try to mention the relevant references so that the interested reader can follow up for a more in-depth view.

Approximate Inference Algorithms for Approaches Based on Local Conditional Classifiers

Two of the most commonly used approximate inference algorithms following this approach are the iterative classification algorithm (ICA) (Neville and Jensen 2000, Lu and Getoor 2003) and Gibbs sampling (GS), and we next describe these in turn.

Iterative Classification

The basic premise behind ICA is extremely simple. Consider a node $Y_i \in \mathcal{Y}$ whose value we need to determine and suppose we know the values of all the other nodes in its neighborhood \mathcal{N}_i (note that \mathcal{N}_i can contain both observed and unobserved variables). Then, ICA assumes that we are given a local classifier f that takes the values of \mathcal{N}_i as arguments and returns the best label value for Y_i from the class label set \mathcal{L} . For local classifiers f that do not return a class label but a goodness or likelihood value given a set of attribute values and a label, we simply choose the label that corresponds to the maximum

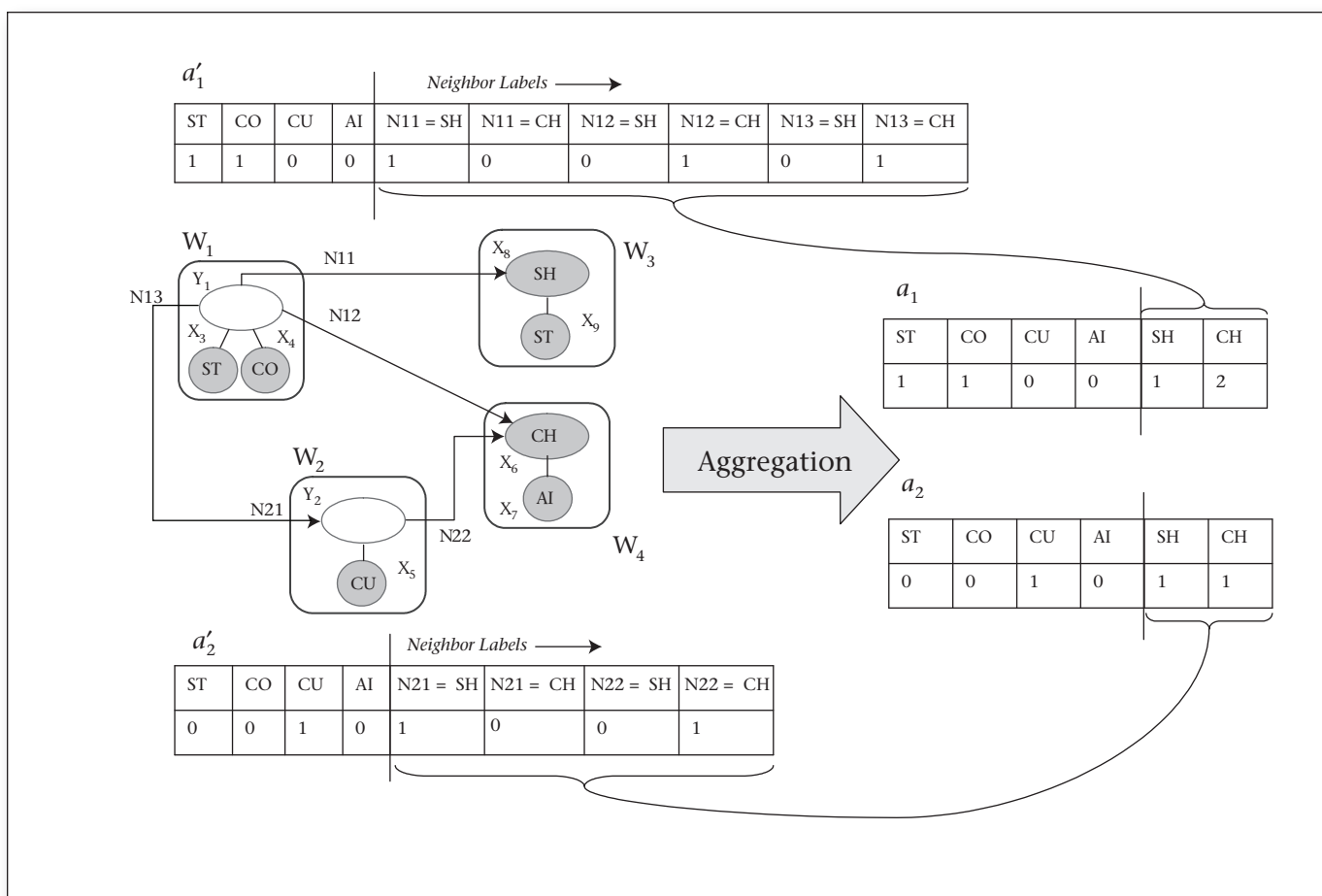


Figure 1. A Small Web Page Classification Problem.

Each box denotes a web page, each directed edge between a pair of boxes denotes a hyperlink, each oval node denotes a random variable, each shaded oval denotes an observed variable, whereas an unshaded oval node denotes an unobserved variable whose value needs to be predicted. Assume that the set of label values is $\mathcal{L} = \{SH, CH\}$. The figure shows a snapshot during a run of ICA. Assume that during the last ICA labeling iteration we chose the following labels: $\gamma_1 = SH$ and $\gamma_2 = CH$. a'_1 and a'_2 show what may happen if we try to encode the respective values into vectors naively, that is, we get variable-length vectors. The vectors a_1 and a_2 , obtained after applying count aggregation, show one way of getting around this issue to obtain fixed-length vectors. See the article for more explanation.

goodness or likelihood value; in other words, we replace f with $\text{argmax}_{l \in \mathcal{L}} f$. This makes the local classifier f an extremely flexible function, and we can use anything ranging from a decision tree to an SVM in its place. Unfortunately, it is rare in practice that we know all values in \mathcal{N}_i , which is why we need to repeat the process iteratively, in each iteration, labeling each Y_i using the current best estimates of \mathcal{N}_i and the local classifier f_i , and continuing to do so until the assignments to the labels stabilize.

Most local classifiers are defined as functions whose argument consists of a fixed-length vector of attribute values. Going back to the example we introduced in the last section in figure 1, assume that we are looking at a snapshot of the state of the labels after a few ICA iterations have elapsed and the label values assigned to Y_1 and Y_2 in the last it-

eration are "SH" and "CH," respectively. a'_1 in figure 1 denotes one attempt to pool all values of \mathcal{N}_1 into one vector. Here, the first entry in a'_1 that corresponds to the first neighbor of Y_1 is a "1," denoting that Y_1 has a neighbor that is the word "ST" (X_3), and so on. Unfortunately, this not only requires putting an ordering on the neighbors, but since Y_1 and Y_2 have a different number of neighbors, this type of encoding results in a'_1 consisting of a different number of entries than a'_2 . Since the local classifier can take only vectors of a fixed length, this means we cannot use the same local classifier to classify both a'_1 and a'_2 .

A common approach to circumvent such a situation is to use an aggregation operator such as *count*, *mode*, or *prop*. Figure 1 shows the two vectors a_1 and a_2 that are obtained by applying the count operator to a'_1 and a'_2 , respectively. The *count* op-

```

for each node  $Y_i \in \mathcal{Y}$  do // bootstrapping
  // compute label using only observed nodes in  $\mathcal{N}_i$ 
  compute  $\mathbf{a}_i$  using only  $\mathcal{X} \cap \mathcal{N}_i$ 
   $y_i \leftarrow f(\mathbf{a}_i)$ 
end for
repeat // iterative classification
  generate ordering  $O$  over nodes in  $\mathcal{Y}$ 
  for each node  $Y_i \in O$  do
    // compute new estimate of  $y_i$ 
    compute  $\mathbf{a}_i$  using current assignments to  $\mathcal{N}_i$ 
     $y_i \leftarrow f(\mathbf{a}_i)$ 
  end for
until all class labels have stabilized or a threshold number of
iterations have elapsed

```

Algorithm 1. Iterative Classification Algorithm (ICA).

erator simply counts the number of neighbors assigned “SH” and “CH” and adds these entries to the vectors, while the *prop* operator returns the proportion of neighbors with each label. Thus we get one new value for each entry in the set of label values. Assuming that the set of label values does not change from one unobserved node to another, this results in fixed-length vectors that can now be classified using the same local classifier. Thus, \mathbf{a}_1 , for instance, contains two entries beside the entries corresponding to the local word attributes, encoding that Y_1 has one neighbor labeled “SH” (X_8), and two neighbors currently labeled “CH” (Y_2 and X_6). Algorithm 1 depicts the ICA algorithm as pseudocode where we use \mathbf{a}_i to denote the vector encoding the values in \mathcal{N}_i obtained after aggregation. Note that in the first ICA iteration, all labels y_i are undefined, and to initialize them we simply apply the local classifier to the observed attributes in the neighborhood of Y_i ; this is referred to as “bootstrapping” in algorithm 1.

Gibbs Sampling

Gibbs sampling (GS) (Gilks, Richardson, and Spiegelhalter 1996) is widely regarded as one of the most accurate approximate inference procedures. It was originally proposed by Geman and Geman (1984) in the context of image restoration. Unfortunately, it is very slow, and a common issue while implementing GS is to determine when the procedure has converged. Even though there are tests that can help one determine convergence, they are usually expensive or complicated to implement.

Due to the issues with traditional GS, researchers in collective classification (Macskassy and Provost 2007; McDowell, Gupta, and Aha 2007) use a simplified version where they assume, just like in the

case of ICA, that we have access to a local classifier f that can be used to estimate the conditional probability distribution of the labels of Y_i given all the values for the nodes in \mathcal{N}_i . Note that, unlike traditional GS, there is no guarantee that this conditional probability distribution is the correct conditional distribution to sample from. At best, we can only assume that the conditional probability distribution given by the local classifier f is an approximation of the correct conditional probability distribution. Neville and Jensen (2007) provide more discussion and justification for this line of thought in the context of relational dependency networks, where they use a similar form of GS for inference.

The pseudocode for GS is shown in algorithm 2. The basic idea is to sample for the best label estimate for Y_i given all the values for the nodes in \mathcal{N}_i using local classifier f for a fixed number of iterations (a period known as “burn-in”). After that, not only do we sample for labels for each $Y_i \in \mathcal{Y}$ but we also maintain count statistics as to how many times we sampled label l for node Y_i . After collecting a predefined number of such samples, we output the best label assignment for node Y_i by choosing the label that was assigned the maximum number of times to Y_i while collecting samples. For all our experiments (that we report later) we set burn-in to 200 iterations and collected 800 samples.

Feature Construction and Further Optimizations

One of the benefits of both ICA and GS is the fact that it is fairly simple to make use of any local classifier. There is some evidence to indicate that some local classifiers tend to produce higher accuracies than others, at least in the application domains where such experiments have been conducted. For instance, Lu and Getoor (2003) report that on bibliography data sets and web page classification problems logistic regression outperforms naïve Bayes.

Recall that, to represent the values of \mathcal{N}_i , we described the use of an aggregation operator. In the example, we used the *count* operator to aggregate values of the labels in the neighborhood, but *count* is by no means the only aggregation operator available. Past research has used a variety of aggregation operators including *minimum*, *maximum*, *mode*, *exists*, and *proportion*. The choice of which aggregation operator to use depends on the application domain and relates to the larger question of relational feature construction where we are interested in determining which features to use so that classification accuracy is maximized. In addition, values derived from the graph structure of the network in the data, such as the clustering coefficient and betweenness centrality, may be beneficial to

the accuracy of the classification task. Within the inductive logic programming community, aggregation has been studied as a means for propositionalizing a relational classification problem (Knobbe, deHaas, and Siebes 2001; Kramer, Lavrac, and Flach 2001). Within the statistical relational learning community, Perlich and Provost (2003, 2006) have studied aggregation extensively, and Popescul and Ungar (2003) have worked on feature construction using techniques from inductive logic programming. In addition, Macskassy and Provost (2007) have investigated approaches that only make use of the labels of neighbors.

Other aspects of ICA that have been the subject of investigation include the ordering strategy to determine in which order to visit the nodes to re-label in each ICA iteration. There is some evidence to suggest that ICA is fairly robust to a number of simple ordering strategies such as random ordering, visiting nodes in ascending order of diversity of its neighborhood class labels, and labeling nodes in descending order of label confidences. However, there is also some evidence that certain modifications to the basic ICA procedure tend to produce improved classification accuracies. For instance, both Neville and Jensen (2005) and McDowell, Gupta, and Aha (2007) propose a strategy where only a subset of the unobserved variables are utilized as inputs for feature construction. More specifically, in each iteration, they choose the top- k most confident predicted labels and use only those unobserved variables in the following iteration's predictions, thus ignoring the less confident predicted labels. In each subsequent iteration they increase the value of k so that in the last iteration all nodes are used for prediction. McDowell and colleagues report that such a "cautious" approach leads to improved accuracies.

Approximate Inference Algorithms for Approaches based on Global Formulations

An alternate approach to performing collective classification is to define a global objective function to optimize. In what follows, we will describe one common way of defining such an objective function and this will require some more notation.

We begin by defining a pairwise Markov random field (pairwise MRF) (Taskar, Abbeel, and Koller 2002). Intuitively, a pairwise MRF is a graph where each node denotes a random variable and every edge denotes a correlation between a pair of random variables. Let $G = (\mathcal{V}, E)$ denote a graph of random variables as before where \mathcal{V} consists of two types of random variables, the unobserved variables, \mathcal{Y} , which need to be assigned values from label set \mathcal{L} , and observed variables \mathcal{X} whose values we know.

```

for each node  $Y_i \in \mathcal{Y}$  do // bootstrapping
  // compute label using only observed nodes in  $\mathcal{N}_i$ 
  compute  $\mathbf{a}_i$  using only  $\mathcal{X} \cap \mathcal{N}_i$ 
   $y_i \leftarrow f(\mathbf{a}_i)$ 
end for
for  $n = 1$  to  $B$  do // burn-in
  generate ordering  $O$  over nodes in  $\mathcal{Y}$ 
  for each node  $Y_i \in O$  do
    compute  $\mathbf{a}_i$  using current assignments to  $\mathcal{N}_i$ 
     $y_i \leftarrow f(\mathbf{a}_i)$ 
  end for
end for
for each node  $Y_i \in \mathcal{Y}$  do // initialize sample counts
  for each label  $l \in \mathcal{L}$  do
     $c[i, l] = 0$ 
  end for
end for
for  $n = 1$  to  $S$  do // collect samples
  generate ordering  $O$  over nodes in  $\mathcal{Y}$ 
  for each node  $Y_i \in O$  do
    compute  $\mathbf{a}_i$  using current assignments to  $\mathcal{N}_i$ 
     $y_i \leftarrow f(\mathbf{a}_i)$ 
     $c[i, y_i] \leftarrow c[i, y_i] + 1$ 
  end for
end for
for each node  $Y_i \in \mathcal{Y}$  do // compute final labels
   $y_i \leftarrow \operatorname{argmax}_{l \in \mathcal{L}} c[i, l]$ 
end for
```

Algorithm 2. Gibbs Sampling Algorithm (GS).

To quantify the exact nature of each correlation in an MRF we use small functions usually referred to as *clique potentials*. A clique potential is simply a function defined over a small set of random variables that maps joint assignments of the set of random variables to goodness values, numbers that reflect how favorable the assignment is. Let Ψ denote a set of clique potentials. Ψ contains three distinct types of functions:

- (1) For each $Y_i \in \mathcal{Y}$, $\psi_i \in \Psi$ is a mapping $\psi_i: \mathcal{L} \rightarrow \mathfrak{R}_{\geq 0}$, where $\mathfrak{R}_{\geq 0}$ is the set of nonnegative real numbers.
- (2) For each $(Y_i, X_j) \in E$, $\psi_{ij} \in \Psi$ is a mapping $\psi_{ij}: \mathcal{L} \rightarrow \mathfrak{R}_{\geq 0}$.
- (3) For each $(Y_i, Y_j) \in E$, $\psi_{ij} \in \Psi$ is a mapping $\psi_{ij}: \mathcal{L} \times \mathcal{L} \rightarrow \mathfrak{R}_{\geq 0}$.

Let \mathbf{x} denote the values assigned to all the observed variables in G and let x_i denote the value assigned to X_i . Similarly, let \mathbf{y} denote any assignment to all the unobserved variables in G and let y_i denote a value assigned to Y_i . For brevity of notation we will denote by ϕ_i the clique potential obtained by computing

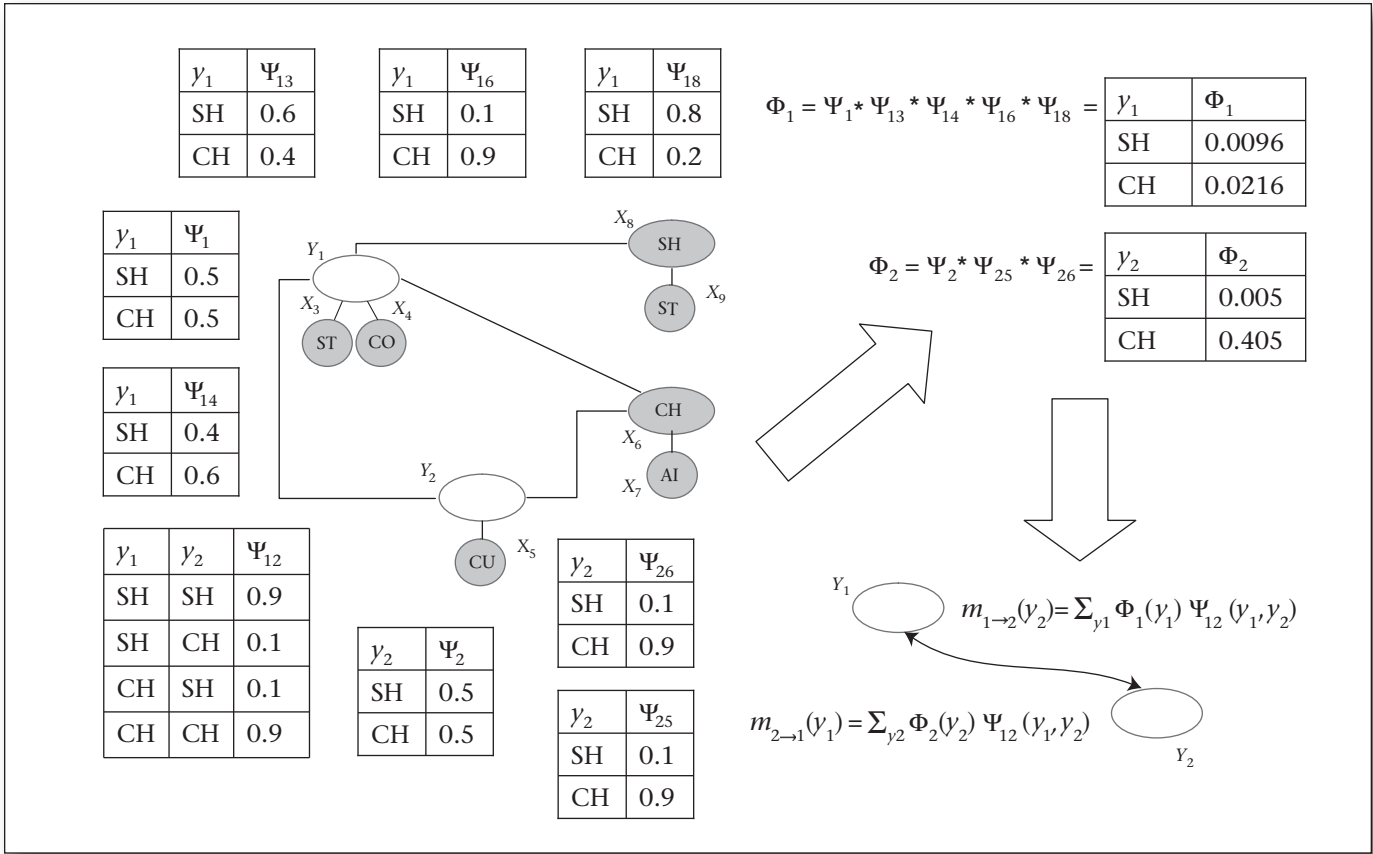


Figure 2. A Small Web Page Classification Problem Expressed as a Pairwise Markov Random Field with Clique Potentials. The figure also shows the message-passing steps followed by LBP. See text for more explanation.

$$\phi_i(y_i) = \psi_i(y_i) \prod_{(y_i, x_j) \in E} \psi_{ij}(y_i, x_j).$$

We are now in a position to define a pairwise MRF.

Definition 1. A pairwise Markov random field (MRF) is given by a pair (G, Ψ) where G is a graph and Ψ is a set of clique potentials with ϕ_i and ψ_{ij} as defined previously. Given an assignment \mathbf{y} to all the unobserved variables \mathcal{Y} , the pairwise MRF is associated with the probability distribution

$$P(\mathbf{y} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{y_i \in \mathcal{Y}} \phi_i(y_i) \prod_{(y_i, y_j) \in E} \psi_{ij}(y_i, y_j)$$

where \mathbf{x} denotes the observed values of \mathcal{X} and

$$Z(\mathbf{x}) = \sum_{\mathbf{y}} \prod_{y_i \in \mathcal{Y}} \phi_i(y_i) \prod_{(y_i, y_j) \in E} \psi_{ij}(y_i, y_j).$$

In figure 2, we show our running example augmented with clique potentials. MRFs are defined on undirected graphs and thus we have dropped the directions on all the hyperlinks in the example. Thus, ψ_1 and ψ_2 denote two clique potentials defined on the unobserved variables (Y_1 and Y_2) in the network. Similarly, we have one ψ defined for each edge that involves at least one unobserved variable as an end point. For instance, ψ_{13} defines a map-

ping from \mathcal{L} (which is set to $\{SH, CH\}$ in the example) to nonnegative real numbers. There is only one edge between the two unobserved variables in the network, and this edge is associated with the clique potential ψ_{12} that is a function over two arguments. Figure 2 also shows how to compute the ϕ clique potentials. Essentially, given an unobserved variable Y_i , one collects all the edges that connect it to observed variables in the network and multiplies the corresponding clique potentials along with the clique potential defined on Y_i itself. Thus, as the figure shows, $\phi_2 = \psi_2 \times \psi_{25} \times \psi_{26}$.

Given a pairwise MRF, it is conceptually simple to extract the best assignments to each unobserved variable in the network. For instance, we may adopt the criterion that the best label value for Y_i is simply the one corresponding to the highest marginal probability obtained by summing over all other variables from the probability distribution associated with the pairwise MRF. Computationally, however, this is difficult to achieve since computing one marginal probability requires summing over an exponentially large number of terms, which is why we need approximate inference algorithms.

We describe two approximate inference algorithms in this article. Both of them adopt a similar approach to avoiding the computational complexity of computing marginal probability distributions. Instead of working with the probability distribution associated with the pairwise MRF directly (definition 1) they both use a simpler “trial” distribution. The idea is to design the trial distribution so that once we fit it to the MRF distribution then it is easy to extract marginal probabilities from the trial distribution (as easy as reading off the trial distribution). This is a general principle that forms the basis of a class of approximate inference algorithms known as variational methods (Jordan et al. 1999).

We are now in a position to discuss loopy belief propagation (LBP) and mean-field relaxation labeling (MF).

Loopy Belief Propagation

Intuitively, loopy belief propagation (LBP) is an iterative message-passing algorithm in which each random variable Y_i in the MRF sends a message $m_{i \rightarrow j}$ to a neighbouring random variable Y_j depicting its belief of what Y_j 's label should be. Thus, a message $m_{i \rightarrow j}$ is simply a mapping from \mathcal{L} to $\mathfrak{R}_{\geq 0}$. The point that differentiates LBP from other gossip-based message protocols is that when computing a message $m_{i \rightarrow j}$ to send from Y_i to Y_j , LBP discounts the message $m_{i \rightarrow j}$ (the message Y_j sent to Y_i in the previous iteration) in an attempt to stop a message computed by Y_j from going back to Y_i . LBP applied to pairwise MRF $\langle G, \Psi \rangle$ can be concisely expressed as the following set of equations:

$$m_{i \rightarrow j}(y_j) = \alpha \sum_{y_i \in \mathcal{L}} \Psi_{ij}(y_i, y_j) \phi_i(y_i) \prod_{Y_k \in \mathcal{N}_i \cap \mathcal{Y} \setminus Y_j} m_{k \rightarrow i}(y_i), \quad \forall y_j \in \mathcal{L} \quad (1)$$

$$b_i(y_i) = \alpha \phi_i(y_i) \prod_{Y_j \in \mathcal{N}_i \cap \mathcal{Y}} m_{j \rightarrow i}(y_i), \quad \forall y_i \in \mathcal{L} \quad (2)$$

where $m_{i \rightarrow j}$ is a message sent by Y_i to Y_j and α denotes a normalization constant that ensures that each message and each set of marginal probabilities sum to 1; more precisely,

$$\sum_{y_j} m_{i \rightarrow j}(y_j) = 1 \text{ and } \sum_{y_i} b_i(y_i) = 1.$$

The algorithm proceeds by making each $Y_i \in Y_j$ communicate messages with its neighbors in $\mathcal{N}_i \cap \mathcal{Y}$ until the messages stabilize (equation 1). After the messages stabilize, we can calculate the marginal probability of assigning Y_i with label y_i by computing $b_i(y_i)$ using equation 2. The algorithm is described more precisely in algorithm 3. Figure 2 shows a sample round of message-passing steps followed by LBP on the running example.

LBP has been shown to be an instance of a variational method. Let $b_i(y_i)$ denote the marginal

```

for each  $(Y_i, Y_j) \in E(G)$  s.t.  $Y_i, Y_j \in \mathcal{Y}$  do
  for  $y_j \in \mathcal{L}$  do
     $m_{i \rightarrow j}(y_j) \leftarrow 1$ 
  end for
end for
repeat // perform message passing
  for each  $(Y_i, Y_j) \in E(G)$  s.t.  $Y_i, Y_j \in \mathcal{Y}$  do
    for each  $y_j \in \mathcal{L}$  do
       $m_{i \rightarrow j}(y_j) \leftarrow \alpha \sum_{y_i} \Psi_{ij}(y_i, y_j) \phi_i(y_i) \prod_{Y_k \in \mathcal{N}_i \cap \mathcal{Y} \setminus Y_j} m_{k \rightarrow i}(y_i)$ 
    end for
  end for
until  $m_{i \rightarrow j}(y_j)$  stop showing any change
for each  $Y_i \in \mathcal{Y}$  do // compute beliefs
  for each  $y_i \in \mathcal{L}$  do
     $b_i(y_i) \leftarrow \alpha \phi_i(y_i) \prod_{Y_j \in \mathcal{N}_i \cap \mathcal{Y}} m_{j \rightarrow i}(y_i)$ 
  end for
end for

```

Algorithm 3. Loopy Belief Propagation (LBP).

probability associated with assigning unobserved variable Y_i the value y_i and let $b_{ij}(y_i, y_j)$ denote the marginal probability associated with labeling the edge (Y_i, Y_j) with values (y_i, y_j) . Then Yedidia, Freeman, and Weiss (2005) showed that the following choice of trial distribution

$$b(\mathbf{y}) = \frac{\prod_{(Y_i, Y_j) \in E} b_{ij}(y_i, y_j)}{\prod_{Y_i \in \mathcal{Y}} b_i(y_i)^{|\mathcal{Y} \cap \mathcal{N}_i| - 1}}$$

and subsequently minimizing the Kullback-Leibler divergence between the trial distribution from the distribution associated with a pairwise MRF gives us the LBP message-passing algorithm with some qualifications. Note that the trial distribution explicitly contains marginal probabilities as variables. Thus, once we fit the distribution, extracting the marginal probabilities is as easy as reading them off.

Relaxation Labeling through Mean-Field Approach

Another approximate inference algorithm that can be applied to pairwise MRFs is mean-field relaxation labeling (MF). The basic algorithm can be described by the following fixed-point equation:

```

for each  $Y_j \in \mathcal{Y}$  do // initialize messages
  for each  $y_j \in \mathcal{L}$  do
     $b_j(y_j) \leftarrow 1$ 
  end for
end for
repeat // perform message passing
  for each  $Y_j \in \mathcal{Y}$  do
    for each  $y_j \in \mathcal{L}$  do
       $b_j(y_j) \leftarrow \alpha \phi_j(y_j) \prod_{Y_i \in \mathcal{N}_j \cap \mathcal{Y}, y_i \in \mathcal{L}} \psi_{ij}^{b_j(y_j)}(y_i, y_j)$ 
    end for
  end for
until all  $b_j(y_j)$  stop changing

```

Algorithm 4. Mean-Field Relaxation Labeling (MF).

$$b_j(y_j) = \alpha \phi_j(y_j) \prod_{Y_i \in \mathcal{N}_j \cap \mathcal{Y}} \prod_{y_i \in \mathcal{L}} \psi_{ij}^{b_j(y_j)}(y_i, y_j),$$

where $b_j(y_j)$ denotes the marginal probability of assigning $Y_j \in \mathcal{Y}$ with label y_j and α is a normalization constant that ensures

$$\sum_{y_j} b_j(y_j) = 1.$$

The algorithm simply computes the fixed-point equation for every node Y_j and keeps doing so until the marginal probabilities $b_j(y_j)$ stabilize. When they do, we simply return $b_j(y_j)$ as the computed marginals. The pseudocode for MF is shown in algorithm 4.

MF can also be justified as a variational method in almost exactly the same way as LBP. In this case, however, we choose a simpler trial distribution:

$$b(\mathbf{y}) = \prod_{Y_i \in \mathcal{Y}} b_i(y_i).$$

We refer the interested reader to the work of Weiss and colleagues (Weiss 2001; Yedidia, Freeman, and Weiss 2005) for more details.

Experiments

In our experiments, we compared the four collective classification algorithms (CC) discussed in the previous sections and a content-only classifier (CO), which does not take the network into account, along with two choices of local classifiers on document classification tasks. The two local classifiers we tried were naïve Bayes (NB) and logistic regression (LR). This gave us eight different classifiers: CO with NB, CO with LR, ICA with NB, ICA with LR, GS with NB, GS with LR, MF, and LBP. The data sets we used for the experiments included both real-world and synthetic data sets.

Features Used

For CO classifiers, we used the words in the documents for observed attributes. In particular, we used a binary value to indicate whether or not a word appears in the document. In ICA and GS, we used the same local attributes (that is, words) followed by the *count* aggregation to count the number of each label value in a node's neighborhood. Finally, for LBP and MF, we used pairwise MRF with clique potentials defined on the edges and unobserved nodes in the network.

Experimental Setup

Because we are dealing with network data, traditional approaches to experimental data preparation such as using random sampling for constructing training and test sets may not be directly applicable. To test the effectiveness of collective classification, we want to have splits such that a reasonable number of neighbors are unlabeled. In order to achieve this, we use a snowball sampling evaluation strategy (SS). In this strategy, we construct splits for test data by randomly selecting the initial node and expanding around it. We do not expand randomly; instead, we select nodes based on the class distribution of the whole corpus; that is, the test data is stratified. The nodes selected by the SS are used as the test set while the rest are used for training. We repeat this process k times to obtain k test-train pairs of splits. Besides experimenting on test splits created using SS, we also experimented with splits created using the standard k -fold cross-validation methodology where we choose nodes randomly to create splits and refer to this as RS.

When using SS, some of the objects may appear in more than one test split. In that case, we need to adjust the accuracy computation so that we do not overcount those objects. A simple strategy is to average the accuracy for each instance first and then take the average of the averages. Further, to help the reader compare the SS results with the RS results, we also provide accuracies averaged per instance and across all instances that appear in at least one SS split. We denote these numbers using the term matched cross-validation (M).

Learning the Classifiers

One aspect of the collective classification problem that we have not discussed so far is how to learn the various classifiers described in the previous sections. Learning refers to the problem of determining the parameter values for the local classifier, in the case of ICA and GS, and the entries in the clique potentials, in the case of LBP and MF, which can then be subsequently used to classify unseen test data. For all our experiments, we learned the parameter values from fully labeled data sets created through the splits generation methodology de-

scribed above using gradient-based optimization approaches. For a more detailed discussion, see, for example, Taskar, Abbeel, and Koller (2002), and Sen and Getoor (2007).

Real World Data Sets

We experimented with two real-world bibliographic data sets: Cora (McCallum et al. 2000) and CiteSeer (Giles, Bollacker, and Lawrence 1998). The Cora data set contains a number of machine-learning papers divided into one of seven classes while the CiteSeer data set has six class labels. For both data sets, we performed stemming and stop word removal beside removing the words with document frequency less than 10. The final corpus has 2708 documents, 1433 distinct words in the vocabulary, and 5429 links in the case of Cora; and 3312 documents, 3703 distinct words in the vocabulary, and 4732 links in the case of CiteSeer. For each data set, we performed both RS evaluation (with 10 splits) and SS evaluation (averaged over 10 runs).

Results. The accuracy results for the real-world data sets are shown in table 1. The accuracies are separated by sampling method and base classifier. The highest accuracy at each partition is in bold. We performed a *t*-test (paired where applicable, and Welch *t*-test otherwise) to test statistical significance between results. Here are the main results:

First, do CC algorithms improve over CO counterparts? In both data sets, CC algorithms outperformed their CO counterparts in all evaluation strategies (SS, RS, and M). The performance differences were significant for all comparisons except for the NB (M) results for CiteSeer.

Second, does the choice of the base classifier affect the results of the CC algorithms? We observed a similar trend for the comparison between NB and LR. LR (and the CC algorithms that used LR as a base classifier) outperformed NB versions in all data sets, and the difference was statistically significant for Cora.

Third, is there any CC algorithm that dominates the other? The results for comparing CC algorithms are less clear. In the NB partition, ICA-NB outperformed GS-NB significantly for Cora using SS and M, and GS-NB outperformed ICA-NB for CiteSeer SS. Thus, there was no clear winner between ICA-NB and GS-NB in terms of performance. In the LR portion, again the differences between ICA-LR and GS-LR were not significant for all data sets. As for LBP and MF, they often slightly outperformed ICA-LR and GS-LR, but the differences were not significant.

Fourth, how do SS results and RS results compare? Finally, we take a look at the numbers under the columns labeled M. First, we would like to remind the reader that even though we are comparing the results on the test set that is the intersec-

tion of the two evaluation strategies (SS and RS), different training data could have been potentially used for each test instance, thus the comparison can be questioned. Nonetheless, we expected the matched cross-validation results (M) to outperform SS results simply because each instance had more labeled data around it from RS splitting. The differences were not big (around 1 or 2 percent); however, they were significant. These results tell us that the evaluation strategies can have a big impact on the final results, and care must be taken while designing an experimental setup for evaluating CC algorithms on network data (Gallagher and Eliassi-Rad 2007).

Synthetic Data

We implemented a synthetic data generator following Sen and Getoor (2007). The process proceeds as follows: At each step, we either add a link between two existing nodes or create a node based on the *ld* parameter (such that higher *ld* value means higher link density, that is, more links in the graph) and link this new node to an existing node; when we are adding a link, we choose the source node randomly, but we choose the destination node using the *dh* parameter (which varies homophily by specifying the percentage, on average, of a node's neighbor that is of the same type) as well as the degree of the candidates (preferential attachment). When we are generating a node, we sample a class for it and generate attributes based on this class using a binomial distribution. Then, we add a link between the new node and one of the existing nodes, again using the homophily parameter and the degree of the existing nodes. In all of these experiments, we generated 1000 nodes, where each node is labeled with one of five possible class labels and has 10 attributes. We experimented with varying degree of homophily and link density in the graphs generated.

Results. The results for varying values of *dh*, homophily, are shown in figure 3a. When homophily in the graph was low, both CO and CC algorithms performed equally well, which was the expected result based on similar work. As we increased the amount of homophily, all CC algorithms drastically improved their performance over CO classifiers. With homophily at *dh* = .9, for example, the difference between our best performing CO algorithm and our best performing CC algorithm is about 40 percent. Thus, for data sets that demonstrate some level of homophily, using CC can significantly improve performance.

We present the results for varying the *ld*, link density, parameter in figure 3(b). As we increased the link density of the graphs, we saw that accuracies for all algorithms went up, possibly because the relational information became more significant and useful. However, the LBP accuracy had a

Algorithm	Cora			CiteSeer		
	SS	RS	M	SS	RS	M
CO-NB	0.7285	0.7776	0.7476	0.7427	0.7487	0.7646
ICA-NB	0.8054	0.8478	0.8271	0.7540	0.7683	0.7752
GS-NB	0.7613	0.8404	0.8154	0.7596	0.7680	0.7737
CO-LR	0.7356	0.7695	0.7393	0.7334	0.7321	0.7532
ICA-LR	0.8457	0.8796	0.8589	0.7629	0.7732	0.7812
GS-LR	0.8495	0.8810	0.8617	0.7574	0.7699	0.7843
LBP	0.8554	0.8766	0.8575	0.7663	0.7759	0.7843
MF	0.8555	0.8836	0.8631	0.7657	0.7732	0.7888

Table 1. Accuracy Results for the Cora and CiteSeer Data Sets.

For Cora, the CC algorithms outperformed their CO counterparts significantly. LR versions significantly outperformed NB versions. ICA-NB outperformed GS-NB for SS and M; the other differences between ICA and GS were not significant (both NB and LR versions). Even though MF outperformed ICA-LR, GS-LR, and LBP, the differences were not statistically significant. For CiteSeer, the CC algorithms significantly outperformed their CO counterparts except for ICA-NB and GS-NB for matched cross-validation. CO and CC algorithms based on LR outperformed the NB versions, but the differences were not significant. ICA-NB outperformed GS-NB significantly for SS, but the rest of the differences between LR versions of ICA and GS, LBP, and MF were not significant.

sudden drop when the graph became very dense. The reason behind this result is the well known fact that LBP has convergence issues when there are many short loops in the graph.

Practical Issues

In this section, we discuss some of the practical issues to consider when applying the various CC algorithms. First, although MF and LBP performance is in some cases a bit better than that of ICA and GS, MF and LBP were also the most difficult to work with in both learning and inference. Choosing the initial weights so that the weights will converge during training is nontrivial. Most of the time, we had to initialize the weights with the weights we got from ICA in order to get the algorithms to converge. Thus, the results reported from MF and LBP needed some extra help from ICA to get started. We also note that of the two, we had the most trouble with MF being unable to converge or, when it did, not converging to the global optimum. Our difficulty with MF and LBP is consistent with previous work (Weiss 2001, Mooij and Kappen 2004, Yanover and Weiss 2002) and should be taken into consideration when choosing to apply these algorithms.

In contrast, ICA and GS parameter initializations worked for all data sets we used, and we did not have to tune the initializations for these two algorithms. They were the easiest to train and test among all the collective classification algorithms evaluated.

Finally, while ICA and GS produced very similar

results for almost all experiments, ICA is a much faster algorithm than GS. In our largest data set, CiteSeer, for example, ICA-NB took 14 minutes to run while GS-NB took over 3 hours. The large difference is due to the fact that ICA converges in just a few iterations, whereas GS has to go through significantly more iterations per run due to the initial burn-in stage (200 iterations), as well as the need to run a large number of iterations to get a sufficiently large sampling (800 iterations).

Related Work

Even though collective classification has gained attention only in the past five to seven years, initiated by the work of Jennifer Neville and David Jensen (Neville and Jensen 2000; Jensen, Neville, and Gallagher 2004; Neville and Jensen 2007) and the work of Ben Taskar and colleagues (Taskar, Segal, and Koller 2001; Getoor et al. 2001; Taskar, Abbeel, and Koller 2002; Taskar, Guestrin, and Koller 2003), the general problem of inference for structured output spaces has received attention for a considerably longer period of time from various research communities including computer vision, spatial statistics, and natural language processing. In this section, we attempt to describe some of the work that is most closely related to the work described in this article; however, due to the widespread interest in collective classification, our list is sure to be incomplete.

One of the earliest principled approximate inference algorithms, relaxation labeling (Hummel

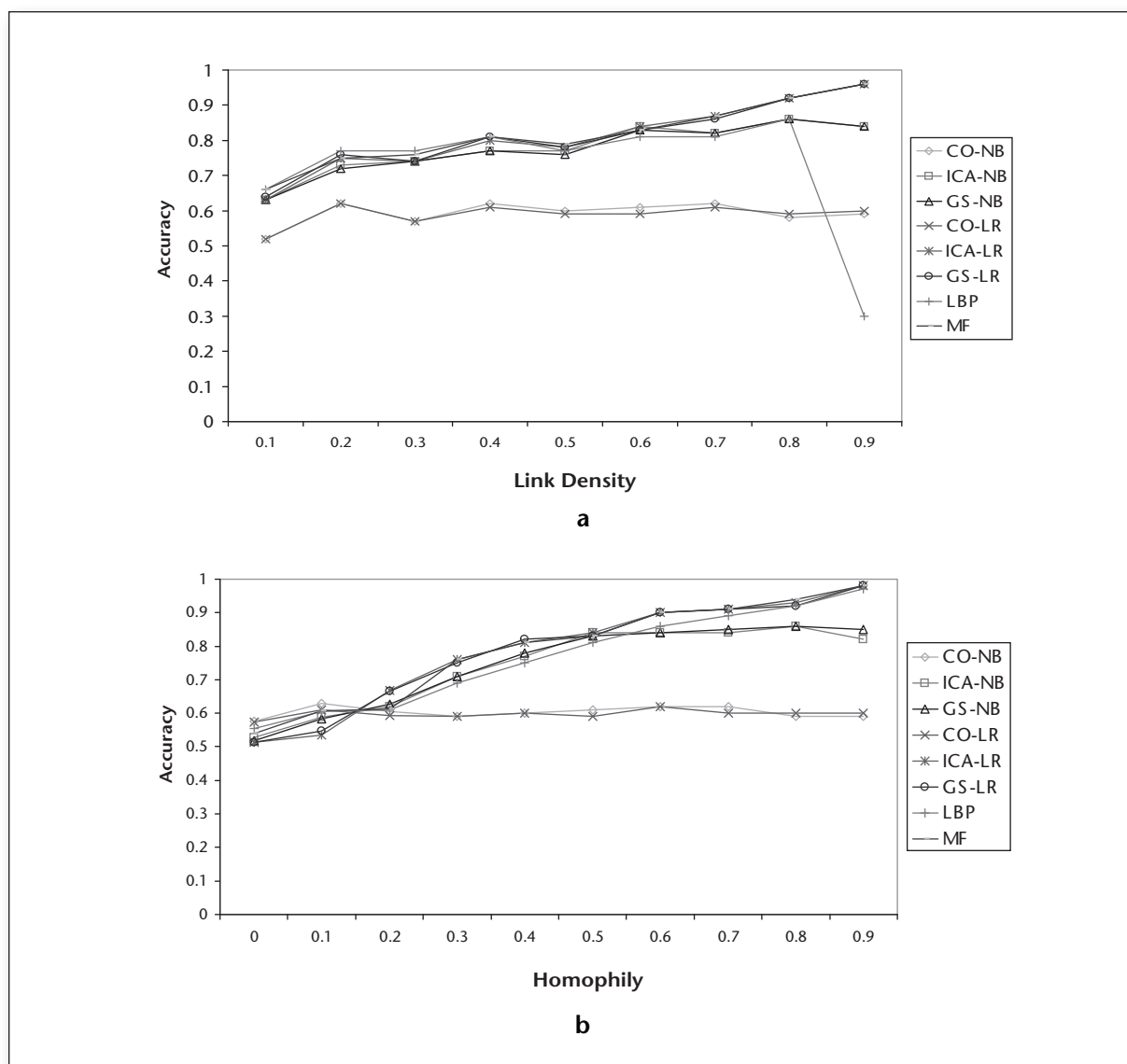


Figure 3. Algorithm Accuracy.

(a) Accuracy of algorithms through different values for dh varying the levels of homophily. When the homophily is very low, both CO and CC algorithms perform equally well, but as we increase homophily, CC algorithms improve over the CO classifier. (b) Accuracy of algorithms through different values for ld varying the levels of link density. As we increase link density, ICA and GS improve their performance the most. Next comes MF. However, LBP has convergence issues due to increased cycles and in fact performs worse than CO for high link density.

and Zucker 1983), was developed by researchers in computer vision in the context of object labeling in images. Due to its simplicity and appeal, relaxation labeling was a topic of active research for some time, and many researchers developed different versions of the basic algorithm (Li, Wang, and Petrou 1994). Mean-field relaxation labeling (Weiss 2001; Yedidia, Freeman, and Weiss 2005), discussed in this article, is a simple instance of this general class of algorithms. Besag (1986) also considered statistical analysis of images and proposed

a particularly simple approximate inference algorithm called *iterated conditional modes*, which is one of the earliest descriptions and a specific version of the iterative classification algorithm presented in this article. Besides computer vision, researchers working with an iterative decoding scheme known as “Turbo Codes” came up with the idea of applying Pearl’s belief propagation algorithm on networks with loops. This led to the development of the approximate inference algorithm that we, in this article, refer to as *loopy belief*

propagation (LBP) (also known as sum product algorithm).

Another area that often uses collective classification techniques is document classification. Chakrabarti, Dom, and Indyk (1998) were among the first to apply collective classification to a corpora of patents linked by hyperlinks and reported that considering attributes of neighboring documents actually hurts classification performance. Slattery and Craven (1998) also considered the problem of document classification by constructing features from neighboring documents using an inductive logic programming rule learner (Slattery and Craven 1998). Yang, Slattery, and Ghani (2002) conducted an in-depth investigation over multiple data sets commonly used for document classification experiments and identified different patterns. Since then, collective classification has also been applied to various other applications such as part-of-speech tagging (Lafferty, McCallum, and Pereira 2001), classification of hypertext documents using hyperlinks (Taskar, Abbeel, and Koller 2002), link prediction in friend-of-a-friend networks (Taskar et al. 2003), classification of email “speech acts” (Carvalho and Cohen 2005), counterterrorism analysis (Macskassy and Provost 2005), and targeted marketing (Hill, Provost, and Volinsky 2007).

Besides the four approximate inference algorithms discussed in this article, there are other algorithms that we did not discuss, such as graphcuts based formulations, and formulations based on linear programming relaxations. More recently, there have been some attempts to extend collective classification techniques to the semisupervised learning scenario (Xu et al. 2006; Macskassy 2007).

Conclusion

In this article, we gave a brief description of four popular collective classification algorithms. We explained the algorithms, showed how to apply them to various applications using examples and highlighted various issues that have been the subject of investigation in the past. Most of the inference algorithms available for practical tasks relating to collective classification are approximate. We believe that a better understanding of when these algorithms perform well will lead to more widespread application of these algorithms to more real-world tasks and that this should be a subject of future research. Most of the current applications of these algorithms have been on homogeneous networks with a single type of unobserved variable that share a common domain. Even though extending these ideas to heterogeneous networks is conceptually simple, we believe that a further investigation into techniques that do so will lead to novel approaches to feature construction and a

deeper understanding of how to improve the classification accuracies of approximate inference algorithms. Collective classification has been a topic of active research for the past decade, and we hope that articles such as this one will help more researchers gain introduction to this area, thus promoting further research into the understanding of existing approximate inference algorithms, and perhaps help develop new, improved inference algorithms.

Acknowledgements

We would like to thank Luke McDowell for his useful and detailed comments. We would also like to thank Jen Neville, David Jensen, Foster Provost, and the reviewers for their comments. This material is based upon work supported in part by the National Science Foundation under Grant No. 0308030. In addition, this work was partially performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344.

References

- Besag, J. 1986. On the Statistical Analysis of Dirty Pictures. *Journal of the Royal Statistical Society Series B*. 48(3): 259–302.
- Carvalho, V., and Cohen, W. W. 2005. On the Collective Classification of Email Speech Acts. In *Proceedings of the 28th Annual International ACM Special Interest Group on Information Retrieval Conference on Research and Development in Information Retrieval*, 345–352. New York: Association for Computing Machinery.
- Chakrabarti, S.; Dom, B.; and Indyk, P. 1998. Enhanced Hypertext Categorization Using Hyperlinks. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, 308–318. New York: Association for Computing Machinery.
- Craven, M.; DiPasquo, D.; Freitag, D.; McCallum, A; Mitchell, T; Nigam, K.; Slattery, S. 1998. Learning to Extract Symbolic Knowledge from the World Wide Web. In *Proceedings of the Fifteenth National Conference on Artificial Intelligence*. Menlo Park, CA: AAAI Press.
- Gallagher, B., and Eliassi-Rad, T. 2007. An Evaluation of Experimental Methodology for Classifiers of Relational Data. Paper presented at the Workshop on Mining Graphs and Complex Structures, IEEE International Conference on Data Mining, Omaha, NE, 28–31 October.
- Geman, S. and Geman, D. 1984. Stochastic Relaxation, Gibbs Distributions and the Bayesian Restoration of Images. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 6(6): 721–742.
- Getoor, L.; Segal, E.; Taskar, B.; and Koller, D. 2001. Probabilistic Models of Text and Link Structure for Hypertext Classification. Paper presented at the IJCAI Workshop on Text Learning: Beyond Supervision. Seattle, WA, 6 August.
- Giles, C. L.; Bollacker, K.; and Lawrence, S. 1998. CiteSeer: An Automatic Citation Indexing System. In *Digital Libraries 98: The Third ACM Conference on Digital Libraries*. New York: Association for Computing Machinery.

- Gilks, W. R.; Richardson, S.; and Spiegelhalter, D. J. 1996. *Markov Chain Monte Carlo in Practice: Interdisciplinary Statistics*. Boca Raton, FL: Chapman and Hall/CRC Press.
- Hill, S.; Provost, F.; and Volinsky, C. 2007. Learning and Inference in Massive Social Networks. Paper presented at the Fifth International Workshop on Mining and Learning with Graphs, Firenze, Italy, August 1–3.
- Hummel, R., and Zucker, S. 1983. On the Foundations of Relaxation Labeling Processes. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 5(3): 267–287.
- Jensen, D.; Neville, J.; and Gallagher, B. 2004. Why Collective Inference Improves Relational Classification. In *Proceedings of the 10th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. New York: Association for Computing Machinery.
- Jordan, M. I.; Ghahramani, Z.; Jaakkola, T. S.; and Saul, L. K. 1999. An Introduction to Variational Methods for Graphical Models. *Machine Learning* 37(2): 183–23.
- Knobbe, A.; deHaas, M.; and Siebes, A. 2001. Propositional-ization and Aggregates. Paper presented at the Fifth European Conference on Principles of Data Mining and Knowledge Discovery, Freiburg, Germany, September 3–5.
- Kramer, S.; Lavrac, N.; and Flach, P. 2001. Propositional-ization Approaches to Relational Data Mining. In *Relational Data Mining*, ed. S. Dzeroski and N. Lavrac. New York: Springer-Verlag.
- Lafferty, J. D.; McCallum, A.; and Pereira, F. C. N. 2001. Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data. In *Proceedings of the Eighteenth International Conference on Machine Learning (ICML 2001)*. San Francisco: Morgan Kaufmann Publishers.
- Li, S. Z.; Wang, H.; Petrou, M. 1994. Relaxation Labeling of Markov Random Fields. In *Proceedings of the 12th IAPR International Conference Pattern Recognition*. Los Alamitos, CA: IEEE Computer Society.
- Lu, Q., and Getoor, L. 2003. Link Based Classification. In *Machine Learning, Proceedings of the Twentieth International Conference (ICML 2003)*. Menlo Park, CA: AAAI Press.
- Macskassy, S. 2007. Improving Learning in Networked Data by Combining Explicit and Mined Links. In *Proceedings of the Twenty-Second Conference on Artificial Intelligence*. Menlo Park, CA: AAAI Press.
- Macskassy, S., and Provost, F. 2005. Suspicion Scoring Based on Guilt-by-Association, Collective Inference, and Focused Data Access. Paper presented at the International Conference on Intelligence Analysis, McLean, VA, 2 May.
- Macskassy, S., and Provost, F. 2007. Classification in Networked Data: A Toolkit and a Univariate Case Study. *Journal of Machine Learning Research*. 8(May): 935–983
- McCallum, A. K.; Nigam, K.; Rennie, J.; and Seymore, K. 2000. Automating the Construction of Internet Portals with Machine Learning. *Information Retrieval Journal* 3(2): 127–163.
- McDowell, L. K.; Gupta, K. M.; and Aha, D. W. 2007. Cautious Inference in Collective Classification. In *Proceedings of the Twenty-Second Conference on Artificial Intelligence*. Menlo Park, CA: AAAI Press.
- Mooij, J., and Kappen, H. J. 2004. Validity Estimates for Loopy Belief Propagation on Binary Real-World Networks. Paper presented at the Conference on Advances in Neural Information Processing Systems 17. December 13–18, 2004, Vancouver, British Columbia, Canada.
- Neville, J., and Jensen, D. 2000. Iterative Classification in Relational Data. In *Learning Statistical Models from Relational Data: Papers from the AAAI Workshop*. Technical Report WS-00-06. Menlo Park, CA: AAAI Press.
- Neville, J., and Jensen, D. 2007. Relational Dependency Networks. *Journal of Machine Learning Research* 8: 653–692.
- Perlich, C., and Provost, F. 2003. Aggregation-Based Feature Invention and Relational Concept Classes. Paper presented at the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Washington, D.C., August.
- Perlich, C., and Provost, F. 2006. Distribution-Based Aggregation for Relational Learning with Identifier Attributes. *Machine Learning Journal* 62(1–2): 65–105.
- Popescul, A., and Ungar, L. 2003. Structural Logistic Regression for Link Analysis. Paper presented at the Second KDD Workshop on Multi-Relational Data Mining, August 27, Washington, D.C.
- Sen, P., and Getoor, L. 2007. Link-based Classification, Technical Report, CS-TR-4858, University of Maryland, College Park, MD.
- Slattery, S., and Craven, M. 1998. Combining Statistical and Relational Methods for Learning in Hypertext Domains. In *Inductive Logic Programming, 8th International Workshop, ILP-98*. Lecture Notes in Computer Science 1446. Berlin: Springer.
- Taskar, B.; Abbeel, P.; and Koller, D. 2002. Discriminative Probabilistic Models for Relational Data. In *Proceedings of the 18th Conference in Uncertainty in Artificial Intelligence*. San Francisco: Morgan Kaufmann, Publishers.
- Taskar, B.; Guestrin, C.; and Koller, D. 2003. Max-Margin Markov Networks. *Advances in Neural Information Processing Systems 16*. Cambridge, MA: The MIT Press.
- Taskar, B.; Segal, E.; and Koller, D. 2001. Probabilistic Classification and Clustering in Relational Data. In *Proceedings of the Seventeenth International Joint Conference on Artificial Intelligence*. San Francisco: Morgan Kaufmann Publishers.
- Taskar, B.; Wong, M. F.; Abbeel, P.; and Koller, D. 2003. Link Prediction in Relational Data. In *Proceedings of the Seventeenth Annual Neural Information Processing Systems Conference*. Cambridge, MA: The MIT Press.
- Weiss, Y. 2001. Comparing the Mean Field Method and Belief Propagation for Approximate Inference in MRFs. In *Advanced Mean Field Methods*, ed. M. Opper and D. Saad. Cambridge, MA: The MIT Press.
- Xu, L.; Wilkinson, D.; Southey, F.; and Schuurmans, D. 2006. Discriminative Unsupervised Learning of Structured Predictors. In *Machine Learning, Proceedings of the Twenty-Third International Conference (ICML 2006)*. ACM International Conference Proceeding Series 148. New York: Association for Computing Machinery.
- Yang, Y.; Slattery, S.; and Ghani, R. 2002. A Study of Approaches to Hypertext Categorization. *Journal of Intelligent Information Systems* 18(2): 219–241.
- Yanover, C., and Weiss, Y. 2002. Approximate Inference and Protein-Folding. In *Proceedings of the Sixteenth Annual Neural Information Processing Systems Conference*. Cambridge, MA: The MIT Press.
- Yedidia, J. S.; Freeman, W. T.; Weiss, Y. 2005. Constructing Free-Energy Approximations and Generalized Belief

Spring Symposium Series Call for Participation

AAAI presents the 2009 Spring Symposium Series, to be held Monday - Wednesday, March 23-25, 2008, at Stanford University. The topics of the nine symposia will be:

- Agents that Learn from Human Teachers. Andrea L. Thomaz (aaaiss09lfh at easychair.org)
- Benchmarking of Qualitative Spatial and Temporal Reasoning Systems. Bernhard Nebel (aaai09bench-qsr at informatik.uni-freiburg.de)
- Experimental Design for Real-World Systems. Katherine Tsui (aaai-sss-2009 at cs.uml.edu)
- Human Behavior Modeling. Tanzeem Choudhury (aaai_ss09_hbm at cs.dartmouth.edu)
- Intelligent Event Processing. Nenad Stojanovic (nstojano at fzi.de)
- Intelligent Narrative Technologies II. Sandy Louchart, Manish Mehta, and David Roberts (int2 at cc.gatech.edu)
- Learning by Reading and Learning to Read. Sergei Nirenburg and Tim Oates (sergei, oates at umbc.edu)
- Social Semantic Web: Where Web 2.0 Meets Web 3.0. Jie Bao (dingl at cs.rpi.edu)
- Technosocial Predictive Analytics. Antonio Sanfilippo (antonio.sanfilippo at pnl.gov)

Submissions for the symposia are due on October 3, 2008. Notification of acceptance will be given by November 7, 2008. Material to be included in the working notes of the symposium must be received by January 16, 2009. The complete Call for Participation is available at www.aaai.org/Symposia/Spring/sss09.php. Registration information will be available by December 15, 2008.

Please contact AAAI at
sss09@aaai.org
with any questions.

Propagation Algorithms. *IEEE Transactions on Information Theory*. 51(7): 2282-2312.



Prithviraj Sen (www.cs.umd.edu/~sen) is a Ph.D. student in the Department of Computer Science at the University of Maryland, College Park. His primary interests lie in the fields of machine learning and probabilistic databases. He received a BTech (Comp. Sc.) from the Indian Institute of Technology, Bombay, in 2003 and an MS (Comp. Sc.) from University of Maryland in 2006.



Galileo Mark Namata (www.cs.umd.edu/~namatag) is a Ph.D. student at the University of Maryland, College Park. He received his B.S. in computer science from Georgetown University in 2003 and his M.S. in computer science from the University of Maryland in 2008. His research interests include artificial intelligence, machine learning, knowledge discovery, and data mining with a particular focus on network data.



Mustafa Bilgic received his MS degree in computer science from the University of Maryland at College Park in 2006, where he is currently a Ph.D. candidate. His research interests are in the areas of machine learning, information acquisition, and decision theory.



Lise Getoor (www.cs.umd.edu/~getoor) received her Ph.D. in computer science from Stanford University in 2001. She is an associate professor in the Department of Computer Science and a member of the Institute for Advanced Computer Studies, University of Maryland, College Park. Her research interests include machine learning, reasoning under uncertainty, databases, and artificial intelligence.



Brian Gallagher is a computer scientist in the Science and Technology Computing Division at Lawrence Livermore National Laboratory. He received a B.A. in computer science from Carleton College in 1998 and an M.S. in computer science from the University of Massachusetts Amherst in 2004. His research interests include artificial intelligence, machine learning, and knowledge discovery and data mining. His current research focuses on classification and analysis in network-structured data.



Tina Eliassi-Rad is a staff computer scientist at the Center for Applied Scientific Computing within Lawrence Livermore National Laboratory. She received her Ph.D. in computer sciences at the University of Wisconsin-Madison in 2001. Broadly speaking, her research interests include artificial intelligence, machine learning, knowledge discovery, and data mining. Her work has been applied to the world wide web, scientific simulation data, and complex networks.