

---

# Collective Object Identification

---

**Parag**      **Pedro Domingos**

Department of Computer Science and Engineering  
University of Washington  
Seattle, WA 98185-2350, U.S.A.  
{*parag, pedrod*}@cs.washington.edu

## Abstract

Object identification is the problem of determining whether different observations correspond to the same object. It occurs in a wide variety of fields, including vision, natural language, citation matching, and information integration. Traditionally, the problem is solved separately for each pair of observations, followed by transitive closure. We propose solving it collectively, performing simultaneous inference for all candidate match pairs, and allowing information to propagate from one candidate match to another via the attributes they have in common. Our formulation is based on conditional random fields, and allows an optimal solution to be found in polynomial time using a graph cut algorithm. Parameters are learned using a voted perceptron algorithm. Experiments on real and synthetic datasets show that collective objective identification outperforms the standard approach.

## 1 Introduction

In many domains, the objects of interest are not uniquely identified, and the problem arises of determining which observations correspond to the same object. For example, in vision we may need to determine whether two similar shapes appearing at different times in a video stream are in fact the same object. In natural language processing and information extraction, a key task is determining which noun phrases are co-referent (i.e., refer to the same entity). When creating a bibliographic database from reference lists in papers, we need to determine which citations refer to the same papers in order to avoid duplication. When merging multiple databases, a problem of keen interest to many large scientific projects, businesses, and government agencies, we need to determine which records represent the same entity and should therefore be merged. This problem, originally defined by Newcombe et al. [13] and placed on a firm statistical footing by Fellegi and Sunter [6], is known by the name of object identification, record linkage, de-duplication, merge/purge, identity uncertainty, hardening soft information sources, co-reference resolution, and others. There is a large literature on it, including [20], [8], [3], [12], [4], [16], [19], [2], etc. Most approaches are variants of the original Fellegi-Sunter model, in which object identification is viewed as a classification problem: given a vector of similarity scores between the attributes of two observations, classify it as “Match” or “Non-match.” A separate match decision is made for each candidate pair, followed by transitive closure to eliminate inconsistencies. Typically, a logistic regression model is used [1].

Making match decisions separately ignores that information gleaned from one match decision may be useful in others. For example, if we find that a paper appearing in *Proceedings*

of *NIPS-2001* is the same as a paper appearing in *Advances in NIPS 14*, this implies that these two strings refer to the same venue, which in turn can help match other pairs of NIPS papers. In this paper, we propose an approach that accomplishes this propagation of information. It is based on conditional random fields, which are discriminatively trained, undirected graphical models [9]. Our formulation allows us to find the globally optimal match in polynomial time using a graph cut algorithm. The parameters of the model are learned using a voted perceptron [5].

Recently, Pasula et al. proposed an approach to the citation matching problem that has collective inference features [14]. This approach is based on directed graphical models, uses a different representation of the matching problem, also includes parsing of the references into fields, and is quite complex. It is a generative rather than discriminative approach, requiring modeling of all dependencies among all variables, and the learning and inference tasks are correspondingly more difficult. A collective discriminative approach has been proposed by McCallum and Wellner [11], but the only inference it performs across candidate pairs is the transitive closure that is traditionally done as a post-processing step. Our model can be viewed as a form of relational Markov network [17], except that it involves the creation of new nodes for match pairs, and consequently cannot be directly created by queries to the databases of interest. Max-margin Markov networks [18] can also be viewed as collective discriminative models, and applying their type of margin-maximizing training to our model is an interesting direction for future research.

We first describe in detail our approach, which we call the collective model. We then report experimental results on real and semi-artificial datasets, which illustrate the advantages of our model relative to the standard Fellegi-Sunter one.

## 2 Collective Model

Using the original database-oriented nomenclature, the input to the problem is a database of records (set of observations), with each record being a tuple of fields (attributes). We now describe the graphical structure of our model, its parameterization, and inference and learning algorithms for it.

### 2.1 Model Structure

Consider a database relation  $R = \{r_1, r_2, \dots, r_n\}$ , where  $r_i$  is the  $i^{th}$  record in the relation. Let  $F = \{F^1, F^2, \dots, F^m\}$  denote the set of fields in the relation. For each field  $F^k$ , we have a set  $FV^k$  of corresponding field values appearing in the relation,  $FV^k = \{f_1^k, f_2^k, \dots, f_{l_k}^k\}$ . We will use the notation  $r_i.F^k$  to refer to the value of  $k^{th}$  field of record  $r_i$ . The goal is to determine, for each pair of records  $(r_i, r_j)$ , whether they refer to the same underlying entity. Our graphical model contains three types of nodes:

**Record-match nodes.** The model contains a Boolean node  $R_{ij}$  for each pairwise question of the form: “Is record  $r_i$  the same as record  $r_j$ ?”

**Field-match nodes.** The model contains a Boolean node  $F_{xy}^k$  for each pairwise question of the form: “Do field values  $f_x^k$  and  $f_y^k$  represent the same underlying property?” For example, for the venue field in a bibliography database, the model might contain a node for the question: “Do the strings ‘Adv. NIPS 14’ and ‘Proc. NIPS 2001’ represent the same venue?”

**Field-similarity nodes.** For pair of field values  $f_x^k, f_y^k \in FV^k$ , the model contains a node  $S_{xy}^k$  whose domain is the  $[0, 1]$  interval. This node encodes how similar the two field values are, according to a pre-defined similarity measure. For example, for textual fields

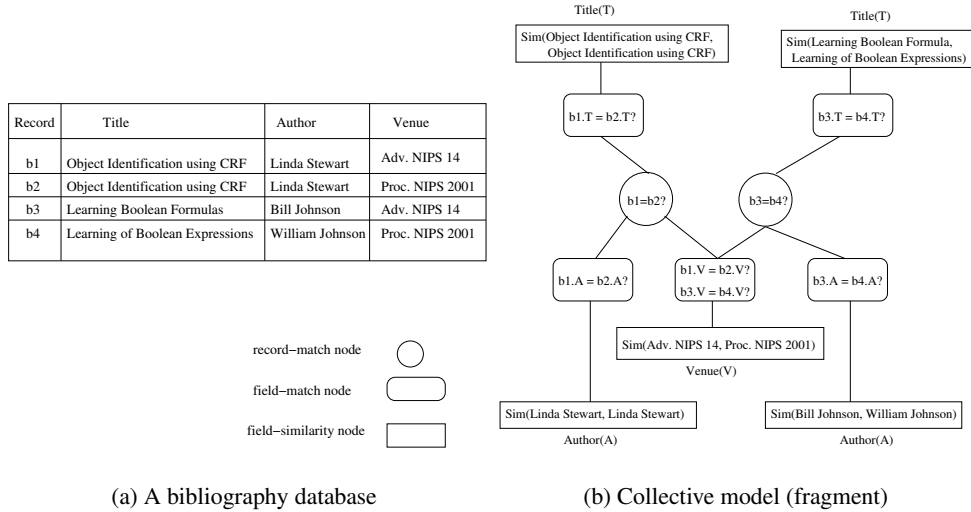


Figure 1: Example of collective object identification

this could be the TF/IDF score [15]. Since their values are computed directly from the data, we will also call these nodes *evidence nodes*.

Because of the symmetric nature of their semantics,  $R_{ij}$ ,  $F_{xy}^k$  and  $S_{xy}^k$  represent the same nodes as  $R_{ji}$ ,  $F_{yx}^k$  and  $S_{yx}^k$ , respectively.

The structure of the model is as follows. Each record-match node  $R_{ij}$  is connected by an edge to each corresponding field-match node  $F_{xy}^k$ ,  $1 \leq k \leq m$ . Formally,  $R_{ij}$  is connected to  $F_{xy}^k$  iff  $r_i.F^k = f_x^k$  and  $r_j.F^k = f_y^k$ . Each field-match node  $F_{xy}^k$  is connected to the corresponding field-similarity node  $S_{xy}^k$ . In general, a field-match node will be linked to many record-match nodes, as the same pair of field values can be shared by many record pairs. This sharing lies at the heart of our model. The field-match nodes allow information to propagate from one candidate record pair to another. Notice that directly connecting the record-match nodes to the evidence nodes and merging the evidence nodes corresponding to the same field value pairs, without introducing field-match nodes, would not work. This is because evidence nodes have known values at inference time, rendering the record-match nodes independent and reducing our approach to the standard one. Figure 1(a) shows a four-record bibliography database, and 1(b) shows the corresponding graphical representation for the candidate pairs  $(b_1, b_2)$  and  $(b_3, b_4)$ . Note how dependencies flow through the shared field-match node corresponding to the venue field. Inferring that  $b_1$  and  $b_2$  refer to the same underlying paper will lead to the inference that the corresponding venue strings “Adv. NIPS 14” and “Proc. NIPS 2001” refer to the same underlying venue, which in turn might provide sufficient evidence to merge  $b_3$  and  $b_4$ . In general, our model can capture complex interactions between candidate pair decisions, potentially leading to better object identification.

## 2.2 Conditional Random Fields

Conditional random fields, introduced by Lafferty et al. [9], define the conditional probability of a set of output variables  $\mathbf{Y}$  given a set of input or evidence variables  $\mathbf{X}$ . Formally,

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z_{\mathbf{x}}} \sum_{c \in C} \exp \sum_l \lambda_{lc} f_{lc}(y_c, x_c) \quad (1)$$

where  $C$  is the set of cliques in the graph,  $x_c$  and  $y_c$  denote the subset of variables participating in clique  $c$ , and  $Z_{\mathbf{x}}$  is a normalization factor.  $f_{lc}$ , known as a feature function, is a function of variables involved in clique  $c$ , and  $\lambda_{lc}$  is the corresponding weight. In many domains, rather than having different parameters (feature weights) for each clique in the graph, the parameters of a conditional random field are tied across repeating clique patterns in the graph, called clique templates [17]. The probability distribution can then be specified as

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z_{\mathbf{x}}} \sum_{t \in T} \sum_{c \in C_t} \exp \sum_l \lambda_{lt} f_{lt}(y_c, x_c) \quad (2)$$

where  $T$  is the set of all the templates,  $C_t$  is the set of cliques which satisfy template  $t$ , and  $f_{lt}$  and  $\lambda_{lt}$  are respectively a feature function and a feature weight, pertaining to template  $t$ .

### 2.3 Model Parameters

Our model has a singleton clique for each record-match node and one for each field-match node, a two-way clique for each edge linking a record-match node and a field-match node, and a two-way clique between each field-match node and the corresponding field-similarity node. The parameters for all cliques of the same type are tied; there is a template for the singleton record-match cliques, one for each type of singleton field-match clique (e.g., in a bibliography database, one for author fields, one for title fields, one for venue fields, etc.), and so on. The probability of a particular assignment  $\mathbf{r}$  to the record-match and field-match nodes, given that the field-similarity (evidence) node values are  $\mathbf{s}$ , is

$$P(\mathbf{r}|\mathbf{s}) = \frac{1}{Z_{\mathbf{s}}} \exp \sum_{i,j} \left[ \sum_l \lambda_l f_l(r_{ij}) + \sum_k \left( \sum_l \phi_{kl} f_l(r_{ij}, F^k) \right. \right. \\ \left. \left. + \sum_l \gamma_{kl} g_l(r_{ij}, r_{ij}, F^k) + \sum_l \delta_{kl} h_l(r_{ij}, F^k, r_{ij}, S^k) \right) \right] \quad (3)$$

where  $(i, j)$  ranges over all candidate pairs and  $k$  ranges over all fields.  $r_{ij}, F^k$  and  $r_{ij}, S^k$  refer to the  $k^{th}$  field-match node and field-similarity node, respectively, for the record pair  $(r_i, r_j)$ .  $\lambda_l$  and  $\phi_{kl}$  denote the feature weights for singleton cliques.  $\gamma_{kl}$  denotes the feature weights for a two-way clique between a record-match node and a field-match node.  $\delta_{kl}$  denotes the feature weights for a two-way clique between a field-match node and an evidence node. Cliques have one feature per possible state. Singleton cliques thus have two (redundant) features:  $f_0(x) = 1$  if  $x = 0$ , and  $f_0(x) = 0$  otherwise;  $f_1(x) = 1$  if  $x = 1$ , and  $f_1(x) = 0$  otherwise. Two-way cliques involving Boolean variables have four features:  $g_0(x, y) = 1$  if  $(x, y) = (0, 0)$ ;  $g_1(x, y) = 1$  if  $(x, y) = (0, 1)$ ;  $g_2(x, y) = 1$  if  $(x, y) = (1, 0)$ ;  $g_3(x, y) = 1$  if  $(x, y) = (1, 1)$ ; each of these features is zero in all other states. Two-way cliques between a field-match node  $q$  and a field-similarity node  $s$  have two features, defined as follows:  $h_0(q, s) = 1 - s$  if  $q = 0$ , and  $h_0(q, s) = 0$  otherwise;  $h_1(q, s) = s$  if  $q = 1$ , and  $h_1(q, s) = 0$  otherwise. This captures the fact that, the more similar two field values are, the more likely they are to match.

Notice that a particular field-match node appears in Equation 3 once for each pair of records containing the corresponding field values. This reflects the fact that that node is effectively

the result of merging the field-match nodes from each of the individual record-match decisions.

## 2.4 Inference and Learning

Inference in our model corresponds to finding the configuration  $\mathbf{r}^*$  of non-evidence nodes that maximizes  $P(\mathbf{r}^*|\mathbf{s})$ . For random fields where maximum clique size is two and all non-evidence nodes are binary-valued, this problem can be reduced to a graph min-cut problem, provided certain constraints on the parameters are satisfied [7]. Our model is of this form, and it can be shown that satisfying the following constraints suffices for the min-cut reduction to hold:  $\gamma_{k0} + \gamma_{k3} - \gamma_{k1} - \gamma_{k2} \geq 0, \forall k, 1 \leq k \leq m$ , where the  $\gamma_{kl}, 0 \leq l \leq 3$ , are the parameters of the clique template for edges linking record-match nodes to field-match nodes of type  $F^k$  (see Equation 3).<sup>1</sup> Our learning algorithm ensures that the learned parameters satisfy these constraints. Since min-cut can be solved exactly in polynomial time, we have a polynomial-time exact inference algorithm for our model.

Learning involves finding maximum-likelihood parameters from data. The partial derivative of the log-likelihood  $L$  (see Equation 3) with respect to the parameter  $\gamma_{kl}$  is

$$\frac{\partial L}{\partial \gamma_{kl}} = \sum_{i,j} g_l(r_{ij}, r_{ij}.F^k) - \sum_{\mathbf{r}'} P_{\Lambda}(\mathbf{r}'|\mathbf{s}) \sum_{i,j} g_l(r'_{ij}, r'_{ij}.F^k) \quad (4)$$

where  $\mathbf{r}'$  varies over all possible configurations of the non-evidence nodes in the graph, and  $P_{\Lambda}(\mathbf{r}'|\mathbf{s})$  denotes the probability distribution according to the current set of parameters. In words, the derivative of the log-likelihood with respect to a parameter is the difference between the empirical and expected counts of the corresponding feature, with the expectation taken according to the current model. The other components of the gradient are found analogously. To satisfy the constraint  $\gamma_{k0} + \gamma_{k3} - \gamma_{k1} - \gamma_{k2} \geq 0$ , we perform the following reparameterization:  $\gamma_{k0} = f(\beta_1) + \beta_2$ ,  $\gamma_{k1} = f(\beta_1) - \beta_2$ ,  $\gamma_{k2} = -f(\beta_3) + \beta_4$ ,  $\gamma_{k3} = -f(\beta_3) - \beta_4$ , where  $f(x) = \log(1 + e^x)$ . We then learn the  $\beta$  parameters using the appropriate transformation of Equation 4. The second term in this equation involves the expectation over an exponential number of configurations, and its computation is intractable. We use a voted perceptron algorithm [5], which approximates this expectation by the feature counts of the most likely configuration, which we find using our polynomial-time inference algorithm with the current parameters. The final parameters are the average of the ones learned during each iteration of the algorithm. Notice that, because parameters are learned at the template level, we are able to propagate information through field values that did not appear in the training data.

## 3 Experiments

We performed experiments on real and semi-artificial datasets, comparing the performance of the collective model and the standard Fellegi-Sunter model using logistic regression. If we consider every possible pair of records for a match, the potential number of matches is  $O(n^2)$ , which is a very large number even for datasets of moderate size. Therefore, we used the technique of first clustering the dataset into possibly-overlapping *canopies* using an inexpensive distance metric, as described by McCallum et al. [10], and then applying our inference and learning algorithms only to record pairs which fall in the same canopy. This reduced the number of potential matches to on the order 1% of all possible matches. In our experiments we used this technique with both our model and the standard model. The results that we report are inclusive of the canopy process, i.e., they are over all the

---

<sup>1</sup>The constraint mentioned in Greig et al. [7] translates to  $\gamma_{k0}, \gamma_{k3} \geq 0, \gamma_{k1}, \gamma_{k2} \leq 0$ , which is a more restrictive version of the constraint above.

Table 1: Experimental results on the Cora dataset.

Model	Before transitive closure			After transitive closure		
	F-measure	Recall	Precision	F-measure	Recall	Precision
Standard	84.4%	81.4%	88.5%	80.7%	92.0%	73.6%
Collective	86.9%	89.0%	85.8%	86.9%	90.8%	84.1%

possible  $O(n^2)$  candidate match pairs. The field-similarity nodes were computed using cosine similarity with TF/IDF [15].

### 3.1 Real-World Data

Our primary source of data was the hand-labeled subset of the Cora dataset provided by McCallum<sup>2</sup> and previously used by Bilenko and Mooney [2] and others. This dataset is a collection of 1295 different citations to 112 computer science research papers from the Cora Computer Science Research Paper Engine. The original data set contains only unsegmented citation strings. Bilenko and Mooney [2] used a segmented version of the data for their experiments, with each bibliographic reference split into its constituent fields (author, venue, title, publisher, year, etc.) using an information extraction system. We used this processed version of the Cora dataset for our experiments. We used only the three most informative fields: author, title and venue (with venue encompassing different types of publication venue, such as conferences, journals, workshops, etc.). Two field values were labeled as matching if they appeared in at least one pair of matching records, and as non-matching otherwise. We divided the data into equal-sized training and test sets, ensuring that no true set of matching records was split among the two, to avoid contamination of the test data by the training set. We performed two-fold cross-validation, and report the average F-measure, recall and precision over twenty random splits of the data. Next, we took the transitive closure over the matches produced by each model as a post-processing step to remove any inconsistent decisions. Table 1 shows the results obtained before and after taking the transitive closure step.

Before taking the transitive closure, the collective model gives an F-measure gain of about 2.5% over the standard model, which is the result of a substantial gain in recall that outweighs a smaller loss in precision. After taking the transitive closure, the recall of the standard model is greatly improved, but the precision is reduced even more drastically, resulting in a substantial deterioration in F-measure. On the other hand, the collective model is relatively stable to the transitive closure step, with its F-measure remaining the same as a result of small increase in recall and a small loss in precision. We attribute this to the fact that the flow of information facilitated by the collective model not only improves performance but also tends to make match decisions more consistent. The net F-measure gain of the collective model over the standard model after transitive closure step is about 6.2%. These results indicate the promise of the collective approach.

### 3.2 Semi-Artificial Data

To further observe the behavior of the algorithms, we generated variants of the Cora dataset by taking distinct field values from the original dataset and randomly combining them to generate distinct papers. This allowed us to control various factors like the number of clusters, level of distortion, etc., and observe how these factors affect the performance of our algorithm. To generate the semi-artificial dataset, we created eight distorted duplicates of each field value taken from the Cora dataset. The number of distortions within each

<sup>2</sup><http://www.cs.umass.edu/~mccallum/data/cora-refs.tar.gz>

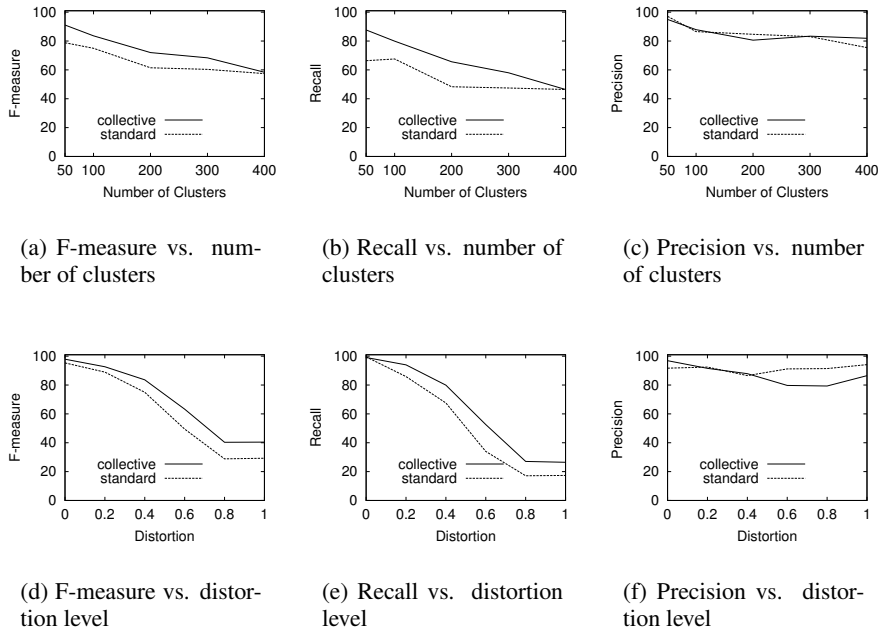


Figure 2: Experimental results on semi-artificial datasets

duplicate was chosen according to a binomial distribution whose “probability of success” parameter we varied in our experiments; a single Bernoulli trial corresponds to the distortion of a single word in the original string. The total number of records was kept constant at 1000 in all the experiments with semi-artificial data. To generate the records in the dataset, we first decided the number of clusters, and then created duplicate records for each cluster by randomly choosing the duplicates for each field value in the cluster. The results reported for semi-artificial data are obtained by performing two-fold cross validation over five random splits of the data. All the results reported are before taking the transitive closure step.

The first set of experiments compared the relative performance of the standard model and the collective model as we varied the number of clusters from 50 to 400, with the first two cluster sizes being 50 and 100 and then varying the size at an interval of 100. The binomial distortion parameter was kept at 0.4. Figures 2(a) to 2(c) show the results. The F-measure (Figure 2(a)) drops as the number of clusters is increased, but the collective model always outperforms the standard model. The recall curve (Figure 2(b)) shows similar behavior. Precision (Figure 2(c)) seems to drop with increasing number of clusters, with neither of the models emerging as the clear winner.

The second set of experiments compared the relative performance of the two models as we varied the level of distortion from 0 to 1, at intervals of 0.2. (0 means no distortion, and 1 means that every word in the string is distorted.) The number of clusters in the dataset was kept constant at 100. Figures 2(d) to 2(f) show the results. As expected, the F-measure (Figure 2(d)) drops as the level of distortion in the data increase, with the collective model dominating throughout. The recall curve (Figure 2(b)) shows similar behavior. Precision (Figure 2(c)) seems to fluctuate with increasing distortion, with neither of the models emerging as the clear winner. Overall, the collective model clearly dominates the standard model over a broad range of the number of clusters and level of distortion in the data.

## 4 Conclusion and Future Work

Determining which observations correspond to the same object is a key problem in information integration, citation matching, natural language, vision, and other areas. It is traditionally solved by making a separate decision for each pair of observations. In this paper, we proposed a collective approach, where information is propagated among related decisions via the attribute values they have in common. In our experiments, this produced better results than the standard method. Directions for future work include enriching the model with more complex dependencies (which will entail moving to approximate inference), scaling it to larger data sources, and applying it in a variety of domains.

## References

- [1] A. Agresti. *Categorical Data Analysis*. Wiley, New York, NY, 1990.
- [2] M. Bilenko and R. Mooney. Adaptive duplicate detection using learnable string similarity measures. In *Proc. 9th SIGKDD*, pages 7–12, 2003.
- [3] W. Cohen, H. Kautz, and D. McAllester. Hardening soft information sources. In *Proc. 6th SIGKDD*, pages 255–259, 2000.
- [4] W. Cohen and J. Richman. Learning to match and cluster large high-dimensional data sets for data integration. In *Proc. 8th SIGKDD*, pages 475–480, 2002.
- [5] M. Collins. Discriminative training methods for hidden Markov models: Theory and experiments with perceptron algorithms. In *EMNLP*, 2002.
- [6] I. Fellegi and A. Sunter. A theory for record linkage. *Journal of the American Statistical Association*, 64:1183–1210, 1969.
- [7] D. M. Greig, B. T. Porteous, and A. H. Seheult. Exact maximum a posteriori estimation for binary images. *Journal of the Royal Statistical Society, Series B*, 51:271–279, 1989.
- [8] M. Hernandez and S. Stolfo. The merge/purge problem for large databases. In *Proc. 1995 SIGMOD*, pages 127–138, 1995.
- [9] J. Lafferty, A. McCallum, and F. Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proc. 18th ICML*, pages 282–289, 2001.
- [10] A. McCallum, K. Nigam, and L. Ungar. Efficient clustering of high-dimensional data sets with application to reference matching. In *Proc. 6th SIGKDD*, pages 169–178, 2000.
- [11] A. McCallum and B. Wellner. Object consolidation by graph partitioning with a conditionally trained distance metric. In *Proc. SIGKDD-2003 Workshop on Data Cleaning, Record Linkage, and Object Consolidation*, pages 19–24, 2003.
- [12] A. Monge and C. Elkan. An efficient domain-independent algorithm for detecting approximately duplicate database records. In *Proc. SIGMOD-1997 Workshop on Research Issues in Data Mining and Knowledge Discovery*, 1997.
- [13] H. Newcombe, J. Kennedy, S. Axford, and A. James. Automatic linkage of vital records. *Science*, 130:954–959, 1959.
- [14] H. Pasula, B. Marthi, B. Milch, S. Russell, and I. Shpitser. Identity uncertainty and citation matching. In *Adv. NIPS 15*, 2003.
- [15] G. Salton and M. McGill. *Introduction to Modern Information Retrieval*. McGraw-Hill, New York, NY, 1983.
- [16] S. Sarawagi and A. Bhamidipaty. Interactive deduplication using active learning. In *Proc. 8th SIGKDD*, pages 269–278, 2002.
- [17] B. Taskar, P. Abbeel, and D. Koller. Discriminative probabilistic models for relational data. In *Proc. 18th UAI*, pages 485–492, 2002.
- [18] B. Taskar, C. Guestrin, B. Milch, and D. Koller. Max-margin Markov networks. In *Adv. NIPS 16*, 2004.
- [19] S. Tejada, C. Knoblock, and S. Minton. Learning domain-independent string transformation weights for high accuracy object identification. In *Proc. 8th SIGKDD*, pages 350–359, 2002.
- [20] W. Winkler. The state of record linkage and current research problems. Technical report, Statistical Research Division, U.S. Census Bureau, 1999.