

# An Approach to Collective Entity Linking

Ashish Kulkarni   Kanika Agarwal   Pararth Shah   Sunny Raj Rathod   Ganesh Ramakrishnan

Department of Computer Science  
Indian Institute of Technology Bombay  
Mumbai, India

{kulashish, kanika1712, pararthshah717, sunnyrajrathod} @gmail.com  
ganesh@cse.iitb.ac.in

## Abstract

Entity linking is the task of disambiguating entities in unstructured text by linking them to an entity in a catalog. Several collective entity linking approaches exist that attempt to collectively disambiguate all mentions in the text by leveraging both local mention-entity context and global entity-entity relatedness. However, the complexity of these models makes it unfeasible to employ exact inference techniques and jointly train the local and global feature weights. In this work we present a collective disambiguation model, that, under suitable assumptions makes efficient implementation of exact MAP inference possible. We also present an efficient approach to train the local and global features of this model and implement it in an interactive entity linking system. The system receives human feedback on a document collection and progressively trains the underlying disambiguation model.

## 1 Introduction

Search systems proposed today (Chakrabarti et al., 2006; Cheng et al., 2007; Kasneci et al., 2008; Li et al., 2010) are greatly enriched by recognizing and exploiting entities embedded in unstructured pages. In a typical system architecture (Cucerzan, 2007; Dill et al., 2003; Kulkarni et al., 2009; Milne and Witten, 2008) a spotter first identifies short token segments or “spots” as potential mentions of entities from its catalog. For our purposes, a catalog consists of a directed graph of categories, to which entity nodes are attached. Many entities may qualify for a given text segment, e.g., both *Kernel trick* and *Linux Kernel* might qualify for the text segment “...Kernel...”. In the second stage, a disambiguator assigns zero or more entities to

selected mentions, based on mention-entity coherence, as well as entity-entity similarity.

Some of the recent work (Zhou et al., 2010; Lin et al., 2012) shows that several mentions may have no associated sense in the catalog. This is referred to as the no-attachment (NA) problem (or NIL in the TAC-KBP challenge (McNamee, 2009)). The other, relatively lesser addressed challenge is that of multiple attachments (Kulkarni et al., 2014), where a mention might link to more than one entities from the catalog. This might often be a result of insufficient context and has been acknowledged by some of the recent entity disambiguation challenges<sup>1</sup>.

We present an approach to collective disambiguation of several mentions by combining various mention-entity compatibility and entity-entity relatedness features. Also, unlike most of the prior work, we jointly learn the local and global feature weights. Our Markov network-based model, along with suitable assumptions, makes efficient learning possible. The model links mentions to zero or more entities, thus offering a natural solution to the problem of NAs and multiple attachments.

## 2 Prior Work

Earlier works (Dill et al., 2003; Bunescu and Pasca, 2006; Mihalcea and Csomai, 2007) on entity annotation focused on per-mention disambiguation. This involves selecting the best entity to assign to a mention, independent of the assignments to other mentions in the document. Wikify! (Mihalcea and Csomai, 2007) for instance, uses context overlap for disambiguation and combines it with a classifier model that exploits local and topical features. Cucerzan (Cucerzan, 2007) introduced the notion of agreement on categories of entities in addition to the local context overlap, in which the entity context comprised

<sup>1</sup><http://web-ngram.research.microsoft.com/ERD2014/>  
D S Sharma, R Sangal and E Sherly. Proc. of the 12th Intl. Conference on Natural Language Processing, pages 219–228, Trivandrum, India. December 2015. ©2015 NLP Association of India (NLP AI)

out-links from and in-links to their corresponding Wikipedia documents. Milne *et al.* (Milne and Witten, 2008) formulated a “relatedness” measure of similarity between two entities from Wikipedia, based on their in-link overlap. Relatedness, in conjunction with the prior probability of occurrence of an entity, was then used to train a classifier model. Han *et al.* (Han and Zhao, 2009) leveraged the semantic information in Wikipedia to build a large-scale semantic network and developed a similarity measure to be used for disambiguation. Kulkarni *et al.* (Kulkarni et al., 2009) were the first to propose a general collective disambiguation approach, giving formulations for trade-off between mention-entity compatibility and coherence between entities. Several graph-based approaches (Hoffart et al., 2011; Fahrni and Strube, 2012) followed that cast the disambiguation problem as a problem of dense subgraph selection from a graph of mentions and candidate entities, making use of collective signals.

Most of these systems seem to prefer tagging conservatively. Some of them (Cucerzan, 2007; Hoffart et al., 2011) restrict their tagging to named entities, while others use a subset of entities from a background taxonomy such as TAP (Dill et al., 2003) or Wikipedia (Milne and Witten, 2008; Mihalcea and Csomai, 2007). Others (Kataria et al., 2011; Bhattacharya and Getoor, 2006) have proposed LDA-based generative models but focus only on person names. Some of the more recent systems (Kulkarni et al., 2009; Han et al., 2011) do perform aggressive spotting, aided by the anchor dictionary of Wikipedia entities and study the recall-precision tradeoff.

(Kulkarni et al., 2014) propose a joint disambiguation model based on a Markov network of entities as nodes and edges for their relatedness. Disambiguation is achieved by performing a MAP inference on this graph and it naturally handles the NA and multiple attachment cases. Unlike other approaches (Han and Sun, 2012; Ratinov et al., 2011), their graph models candidate entities with binary labels, instead of mentions with multiple labels. A suitable assumption on cliques and their potentials makes efficient computation of exact inference possible. However, it is not clear as to how the node and edge feature weights are set.

To the best of our knowledge, none of the graph-based models above have attempted to jointly learn the node and edge feature weights. While

there is prior work (Wellner et al., 2004; Wick et al., 2009) that applied graphical models to the problem of information extraction and coreference resolution, exact inference and estimation is intractable in these models. Similar approaches have also been applied to the problem of entity disambiguation (Kulkarni et al., 2009; Han and Sun, 2012; Ratinov et al., 2011), but hardly anyone has attempted to jointly learn the feature weights.

## 2.1 Our Contributions

We leverage the disambiguation model of (Kulkarni et al., 2014) and propose an efficient approach to jointly learn the node and edge feature weights. We also develop an interactive active learning framework that progressively improves the model as more training data becomes available. We implemented our approach in an online annotation system<sup>2</sup> and used it to semi-automatically curate labeled data<sup>3</sup>. Our trained model performs better than several other systems including that of Kulkarni *et al.*

## 3 Preliminaries

We borrowed the features and the disambiguation model from the work described in Kulkarni *et al.* and present it in brief here. We first start with the problem definition.

### 3.1 Problem Definition

The primary goal of document annotation is to link entity mentions in an input document to entities in a catalog. Mentions (or “spots”) are contiguous token sequences in a text, *e.g.* *Bush*, that can potentially link to an entity, *e.g.* *George Bush* in the catalog. Let  $\mathcal{M}_d$  be the set of all mentions in a document  $d$  and  $\mathcal{E}$  be the set of all entities in the catalog. Then, the entity linking problem is to find for each mention  $m \in \mathcal{M}_d$ , the set of entities  $\hat{E}_m \subset \mathcal{E} \cup \{NA\}$  that it can link to.

As a first step, the input text  $d$  is processed by a “spotter” to identify the set  $\mathcal{M}_d$  of mentions and the set of candidates  $E_m \subset \mathcal{E}$ ,  $\forall m \in \mathcal{M}_d$ .  $e_m \in E_m$  is called a candidate entity for spot  $m$ . The set  $E_d = \cup_{m \in \mathcal{M}_d} E_m$  forms the candidate entities set for document  $d$ . This is then followed by a “disambiguation” phase that obtains from  $E_m$ , the set  $\hat{E}_m$  of entities that the mention  $m$  can actually link to. When none of the entities in  $E_m$

<sup>2</sup><http://tinyurl.com/entitydisamb-demo>

<sup>3</sup><http://tinyurl.com/entitydisamb-data>

are valid, then  $\hat{E}_m = \{NA\}$ . Alternatively, more than one entities from  $E_m$  might get promoted to  $\hat{E}_m$ . Assuming *one sense per discourse* (Gale et al., 1992), an entity in the candidate set of more than one mentions, links (or does not link) to all those mentions.

### 3.2 Entity Catalog

A catalog is a structured knowledge base comprising categories with entities under them, along with their attributes and relations. Wikipedia has seen an extensive organic growth and covers entities spanning a vast set of domains. We report experimental results using the Wikipedia dump from May 2011, with approximately 4.4 million entities. For evaluation on ERD, we used as catalog, the snapshot of Freebase as provided in the challenge.

### 3.3 Features

We used three types of features - (1) Popularity-based features of an entity: Prior Sense Probability (Mihalcea and Csomai, 2007), In-Link Count, Out-Link Count; (2) Mention-Entity compatibility features (Kulkarni et al., 2009); (3) Entity-Entity relatedness features: Category-based Similarity (Cucerzan, 2007), In-link based Similarity (Milne and Witten, 2008), Out-link based Similarity, Contextual Similarity.

### 3.4 The Disambiguation Model

Having identified the set of candidate entities for each mention, a disambiguator attempts to link each mention to zero or more entities. The label of a candidate is a collective result of the interplay of local mention-entity and global entity-entity relatedness signals.

A *Markov Random Field* (MRF) is an undirected graphical model that captures local correlations between random variables (Taskar and Koller, 2001). A node is instantiated in the MRF

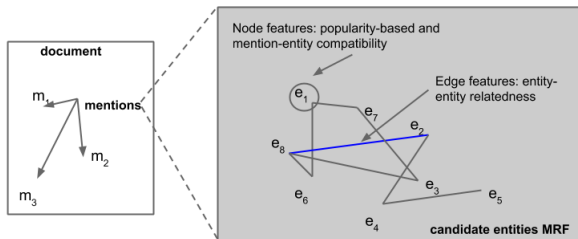


Figure 1: Candidate entities MRF model

graph for each possible entity mapping of each

mention instance in a document. Edges capture entity-entity relatedness. Let  $\mathbf{x}_i$  be the node feature vector of candidate  $i$  and  $\mathbf{x}_{ij}$  be the edge feature vector of the edge joining candidates  $i$  and  $j$ . Each candidate corresponds to a random variable that takes a binary label,  $y_i \in \{0, 1\}$ , based on whether or not it correctly disambiguates the underlying mention. Let  $C$  be the set of all cliques in the MRF and each clique  $c \in C$  be associated with a clique potential  $\phi_c(\cdot)$ . Cliques are restricted to nodes and edges and their potentials are parameterized by log-linear functions of feature vectors; *i.e.*,  $\log \phi_c(\cdot) = \mathbf{w}_c \cdot \mathbf{x}_c$ , where,  $\mathbf{x}_c$  is the feature vector of a clique and  $\mathbf{w}_c$ , the corresponding weight vector. The potentials are assumed to be submodular (Taskar et al., 2004), that is, they are associated with assignments where variables in a clique have the same label. Let  $\mathbf{w}_0$  and  $\mathbf{w}_1$  be the node feature weights influencing node labels 0 and 1 respectively and  $\mathbf{w}_{00}$  and  $\mathbf{w}_{11}$  be the associative edge weights influencing the connected nodes to take the same label. The probability of a complete graph labeling  $y$  is given by  $P(y) = \frac{1}{Z} \prod_{c \in C} \phi_c(y_c)$ , where  $Z$  is the partition function. Disambiguation is achieved by doing MAP inference on this graph.

$$\begin{aligned}
 L(y) &= \arg \max_{y \in \mathcal{Y}} \sum_{i \in \mathcal{N}} \log \phi_i(y_i) + \sum_{ij \in \mathcal{E}} \log \phi_{ij}(y_{ij}) \\
 &= \arg \min_{y \in \mathcal{Y}} - \left( \sum_{i \in \mathcal{N}} \mathbf{w}_0 \cdot \mathbf{x}_i (1 - y_i) + \mathbf{w}_1 \cdot \mathbf{x}_i y_i \right. \\
 &\quad \left. + \sum_{ij \in \mathcal{E}} \mathbf{w}_{00} \cdot \mathbf{x}_{ij} (1 - y_{ij}) + \mathbf{w}_{11} \cdot \mathbf{x}_{ij} y_{ij} \right)
 \end{aligned} \tag{1}$$

$\mathcal{N}$  is the set of nodes and  $\mathcal{E}$  is the set of edges in the MRF,  $y_i \in \{0, 1\}$  and  $y_{ij} = y_i y_j$ . For an MRF with binary labeled nodes and associative edge potentials, MAP inference can be computed exactly in polynomial time, by finding the min cut of an augmented flow graph (Boykov and Kolmogorov, 2004). The MRF is augmented by adding two special terminal nodes *source*,  $s$  and *sink*,  $t$  that correspond to the two labels 0/1. For each node  $i$ , we add terminal edges  $s \rightarrow i$  with weight  $\langle \mathbf{w}_0, \mathbf{x}_i \rangle$  and  $i \rightarrow t$  with weight  $\langle \mathbf{w}_1, \mathbf{x}_i \rangle$ . For each neighborhood edge  $i \rightarrow j$ , we assign weight  $\langle (\mathbf{w}_{00} + \mathbf{w}_{11}), \mathbf{x}_{ij} \rangle$ . We also add weight  $\langle \mathbf{w}_{00}, \mathbf{x}_{ij} \rangle$  to the edge  $s \rightarrow i$  and  $\langle \mathbf{w}_{11}, \mathbf{x}_{ij} \rangle$  to the edge  $j \rightarrow t$ . MAP inference in original MRF corresponds to the  $s/t$  min cut on this augmented

graph, with nodes on the  $s$  side of the cut getting a label of 0, and the nodes on the  $t$  side being assigned a label of 1.

## 4 Learning Feature Weights

---

### Algorithm 1: MRF Learning algorithm

---

**Data:** Training set  $\{X, \hat{q}\}$ , MRF graph  $g$ , Slack penalty  $C$ , Iterations  $T$ , Step size  $\alpha_t$   
**Result:** Weight vector  $w$

```

 $w \leftarrow 0$ 
 $t \leftarrow 1$ 
 $N_n \leftarrow$  number of nodes in  $g$ 
 $f_{opt} \leftarrow \infty$ 
 $w_{opt} \leftarrow 0$ 
while  $t \leq T$  do
   $g \leftarrow$  construct flow network from  $g$ 
   $\hat{q} \leftarrow$   $s/t$  mincut of  $g$ 
   $\nabla_w \xi(w) \leftarrow 2w + C(\hat{q} - \bar{q})^T X$ 
   $w \leftarrow w - \alpha_t \nabla_w \xi(w)$ 
  Project  $w$  onto the positive orthant
  Compute function value  $f$ 
  if  $f < f_{opt}$  then
     $f_{opt} \leftarrow f$ 
     $w_{opt} \leftarrow w$ 
  end
   $t \leftarrow t + 1$ 
end
return  $w_{opt}$ 

```

---

The submodularity restriction and binary labels, make efficient implementation of learning possible. We jointly learn both node and edge feature weights following the general max-margin framework described in (Taskar et al., 2004; Vernaza et al., 2008). Consider a graph with  $\mathcal{N}$  nodes and  $\mathcal{E}$  edges constructed as described above. Following Taskar *et al.*, the learning problem can be formulated in terms of the cut vector, such that, we minimize the norm of the weight vector subject to the constraint that the desired labeling scores better than an arbitrary labeling by an amount that scales with the Hamming distance between the desired and incorrect labelings.

$$\min_{\mathbf{w} \geq 0} \|\mathbf{w}\|^2 \quad (2)$$

subject to

$$\min_{\mathbf{q} \in Q} \sum_{i,j \in \mathcal{E}} \mathbf{w} \cdot \mathbf{x}_{ij} (q_{ij} - \hat{q}_{ij}) - (N_n - \hat{\mathbf{q}}_n^T \cdot \mathbf{q}_n) \geq 0$$

Here,  $Q$  is the set of all valid cuts and  $q_{ij} \in \{0, 1\}$  indicates if edge  $i \rightarrow j$  is cut ( $q_{ij} = 1$ ).  $\mathbf{q}_n$  is the cut vector for terminal edges with components  $q_{si}$  and  $q_{it}$ , where,  $q_{si} = 1$  implies that  $i$  is labeled 1.  $\hat{\mathbf{q}}$  is the cut vector corresponding to the desired labeling. The first component of the constraint captures the difference in cost of the min cut induced by the weights  $\mathbf{w}$  and the desired labeling. The other component corresponds to the

number of labeling disagreements,  $N_n$  being the number of nodes in the graph (excluding  $s$  and  $t$ ).

By rearranging terms, we obtain

$$\min_{\mathbf{w} \geq 0} \|\mathbf{w}\|^2 \quad (3)$$

$$\begin{aligned} \text{subject to } \min_{\mathbf{q} \in Q} \sum_{i,j \in \mathcal{E}} (\mathbf{w}^T \cdot \mathbf{x}_{ij} + \hat{q}_{ij}(\delta_{is} + \delta_{jt}))q_{ij} \\ \geq N_n + \sum_{i,j \in \mathcal{E}} (\mathbf{w}^T \cdot \mathbf{x}_{ij})\hat{q}_{ij} \end{aligned}$$

Here,  $\delta_{ij}$  is the Kronecker delta. It can be shown that the left-hand-side and right-hand-side of the inequality in the constraint are equivalent (Vernaza et al., 2008). Moving the constraint to the objective, we get,

$$\begin{aligned} \min_{\mathbf{w} \geq 0} \|\mathbf{w}\|^2 + C(N_n + \sum_{i,j \in \mathcal{E}} (\mathbf{w}^T \cdot \mathbf{x}_{ij})\hat{q}_{ij} - \\ \min_{\mathbf{q} \in Q} \sum_{i,j \in \mathcal{E}} (\mathbf{w}^T \cdot \mathbf{x}_{ij} + \hat{q}_{ij}(\delta_{is} + \delta_{jt}))q_{ij}) \end{aligned}$$

Summing over all the documents in the training set, we get the final objective,

$$\begin{aligned} \min_{\mathbf{w} \geq 0} \|\mathbf{w}\|^2 + \sum_{d \in D} (C(\mathcal{N}_d + \sum_{i,j \in \mathcal{E}_d} (\mathbf{w}^T \cdot \mathbf{x}_{ij})\hat{q}_{ij} - \\ \min_{\mathbf{q} \in Q} \sum_{i,j \in \mathcal{E}_d} (\mathbf{w}^T \cdot \mathbf{x}_{ij} + \hat{q}_{ij}(\delta_{is} + \delta_{jt}))q_{ij})) \end{aligned} \quad (4)$$

Here,  $\mathbf{w} = [\mathbf{w}_0^T \ \mathbf{w}_1^T \ \mathbf{w}_{00}^T \ \mathbf{w}_{11}^T]^T$ ,  $\mathcal{N}_d$  is the number of nodes (excluding  $s$  and  $t$ ) and  $\mathcal{E}_d$  is the set of edges in the candidate entity MRF graph for a document  $d \in D$ , the set of all training documents,  $s$  and  $t$  are special source and sink nodes, respectively. The term  $\mathcal{N}_d - \hat{q}_{ij}(\delta_{is} + \delta_{jt})q_{ij}$  gives the number of misclassified nodes and  $\sum_{i,j \in \mathcal{E}_d} \mathbf{w}^T \cdot \mathbf{x}_{ij}\hat{q}_{ij} - \mathbf{w}^T \cdot \mathbf{x}_{ij}q_{ij}$  is the total capacity of incorrectly cut edges in the flow graph.  $C$  is the penalty associated with the incorrect labeling. We solved the formulation (4) using the subgradient descent method as described in Algorithm 1.

### 4.1 Handling Unbalanced Training Data

The training data has many more entities labeled 0 as compared to those labeled 1. In our datasets, we observed a skew of about 3 : 1. This results in a bias towards the overrepresented class in the

learning algorithm and the accuracy of the non-dominant class suffers. We addressed this problem by assigning separate misclassification penalties  $C_0$  and  $C_1$  for label 0 and 1 disagreements respectively in equation 4, where, disagreements are defined as below.

**Definition 1.** Let  $l_i \in \{0, 1\}$  and  $\hat{l}_i \in \{0, 1\}$  be the predicted and actual labels of node  $i$ . We say that a node  $i$  has label 0 disagreement if  $l_i \neq \hat{l}_i = 0$ . Similarly it has label 1 disagreement if  $l_i \neq \hat{l}_i = 1$ .

**Proposition 1.** For an edge  $i \rightarrow j$  with  $q_{ij} \neq \hat{q}_{ij}$ , exactly one of the nodes agrees on the label i.e.  $l_i = \hat{l}_i$  (or  $l_j = \hat{l}_j$ ) and the other node disagrees on the label i.e.  $l_j \neq \hat{l}_j$  ( $l_i \neq \hat{l}_i$ ).

*Proof.* Case 1: Let  $q_{ij} \neq \hat{q}_{ij} = 0$ . This implies that the edge is not cut in the actual labeling and therefore  $\hat{l}_i = \hat{l}_j$ . However,  $q_{ij} = 1$  implies that  $l_i \neq l_j$ . It follows that either  $l_i = \hat{l}_i$  or  $l_j = \hat{l}_j$ .

Case 2: Let  $q_{ij} \neq \hat{q}_{ij} = 1$ . Following a similar argument as that for case 1 above, we have that  $\hat{l}_i \neq \hat{l}_j$  and  $l_i = l_j$ . Again, it follows that either  $l_i = \hat{l}_i$  or  $l_j = \hat{l}_j$ .  $\square$

**Definition 2.** An edge  $i \rightarrow j$  with  $q_{ij} \neq \hat{q}_{ij}$ , is said to have a label 0 disagreement if either  $l_i \neq \hat{l}_i = 0$  or  $l_j \neq \hat{l}_j = 0$ . It is said to have a label 1 disagreement if either  $l_i \neq \hat{l}_i = 1$  or  $l_j \neq \hat{l}_j = 1$ .

## 4.2 Active Learning

Our online annotation system presents an opportunity to continuously update the model as more labeled data becomes available. The commonly used *passive learning* approach involves manual annotation of randomly and independently sampled data. Due to the time and cost associated with this process, often there is not enough training data to meet certain level of performance. *Active learning* (Lewis and Catlett, 1994) aims to minimize the labeling effort, by requesting labels for the most informative samples, so as to achieve a desired level of accuracy. While there are several approaches to querying examples for labeling (Li and Sethi, 2006), we follow a pragmatic approach, that can be characterized as *least certain querying* method. The method samples examples with the smallest difference between two highest probability classes. Our binary labeled MRF model labels a node, based on the collective effect of the node potential and the edge potentials on the edges connecting the node to its neighbors. We define cer-

tainty  $C(i)$  at a node  $n_i$  as

$$C(i) = \left| \left( w_0 \cdot x_i + \sum_{(ij) \in \mathcal{E}: j \in N(i)} w_{00} \cdot x_{ij} \right) - \left( w_1 \cdot x_i + \sum_{(ij) \in \mathcal{E}: j \in N(i)} w_{11} \cdot x_{ij} \right) \right| \quad (5)$$

where  $N(i)$  is the set of all neighboring nodes of the node  $i$ . The certainty score  $C(d)$  for a document  $d$  is then computed as the average certainty score across all nodes in that document.

$$C(d) = \frac{1}{|\mathcal{N}|} \sum_{i \in \mathcal{N}} C(i) \quad (6)$$

The active learning algorithm then queries for a document with the lowest  $C(d)$  and presents the document for labeling. It might be possible to further reduce the labeling effort by requesting labels for only top  $k$  entities in the selected document, where the entities are ordered in increasing values of their  $C(i)$ .

## 5 Experiments and Results

### 5.1 Data Sets

We use *Wiki<sub>cur</sub>* (created by (Kulkarni et al., 2014)) for training our model and present cross-validation results. We also evaluate on several other datasets from the entity linking literature. (Kulkarni et al., 2009) had created a dataset (*IITB<sub>part</sub>*) based on aggressive spotting but assuming single attachment. We, therefore, used our annotation system to manually complete annotations (to create *IITB<sub>cur</sub>*) for the documents in this dataset.

### 5.2 Evaluation Measures

We follow the *fuzzy evaluation measure* (Cornolti et al., 2013) that accounts for slight syntactic and semantic variations in the match of a predicted and true annotation, where an annotation  $a$  is defined as the mention-entity pair  $\langle m, e \rangle$ . Using their notion of *weak annotation match*  $M_w(a_1, a_2)$ <sup>4</sup>, we use as performance metrics, *Recall*, *Precision* and *F1* micro-averaged over all documents in a dataset. After factoring out spotter errors, we also separately report the accuracy of our disambiguation model alone (Referred to as ‘‘disambiguator only’’).

<sup>4</sup>which is true iff mentions  $m_1$  and  $m_2$  overlap in the input text and entities  $e_1$  and  $e_2$  are synonyms

Dataset	Disambiguator only			Weak annotation match		
	P	R	F	P	R	F
<i>Wiki<sub>cur</sub></i>	.82	.67	.74	.82	.56	.67
<i>IITB<sub>part</sub></i>	.82	.66	.73	.82	.50	.62

Table 1: Non-collective results (only node features) on *Wiki<sub>cur</sub>* set and *IITB<sub>part</sub>* datasets

### 5.3 Experiments with only Node Features

#### 5.3.1 Is there merit in data curation?

The data curation process presents an opportunity for continuous training where our inference model periodically evolves, as more and more data gets curated. Optionally, in the absence of any curated data to start with (at time  $t = t_0$  when our model is yet untrained), one could use a Logistic Regression model, trained on a large uncurated dataset, to warm-start the data curation process. As data gets curated and our model is trained, we switch to our trained model at time  $t = t_k$ .

We trained binary label LR models using 10000 randomly sampled Wikipedia documents, replacing an original Wikipedia document with its curated version from *Wiki<sub>cur</sub>*, one at a time. Figure 2 plots the training accuracies of these models for an increasing number of curated documents. The improvement in accuracy could be explained by the reduction in false negatives achieved by virtue of aggressive tagging and multiple attachments in the curated dataset. Based on this observation, we claim (and verify it in section 5.4.3) that our MRF model too would benefit from data curation. At the same time, the use of an LR model for warm-starting an online annotation system as ours is strongly recommended.

#### 5.3.2 How does our model perform?

Thereafter, we trained our candidate entity MRF model on *Wiki<sub>cur</sub>* dataset using the node features alone. We report two-fold cross-validation results on *Wiki<sub>cur</sub>* and test results on *IITB<sub>part</sub>* (Refer table 1). These serve as a baseline for our collective approach.

### 5.4 Collective Disambiguation

Next, we trained our model using node features and one or more edge features. Iterations  $T$  were fixed at 600,  $\mathcal{C}$  was tuned as described below, and step size (at iteration  $t$ ),  $\alpha_t = K/\sqrt{t}$ , where  $K$  was empirically set to 0.01.

#### 5.4.1 Effect of $\mathcal{C}$ on accuracy

The  $\mathcal{C}$  parameter in equation 4 acts as a regularizer and is indicative of the tolerance of disagreement between predicted and true labels. It was tuned on the training fold during two-fold cross-validation on *Wiki<sub>cur</sub>*. Also, to account for the skew in label 0 and 1 instances in the dataset, we penalized label 0 and label 1 disagreements separately, using  $\mathcal{C}_0$  and  $\mathcal{C}_1$ , respectively. A higher  $\mathcal{C}_1$  for instance, improves the label 1 recall while adversely impacting the precision. It is this recall-precision tradeoff for varying values of  $\mathcal{C}_0$ ,  $\mathcal{C}_1$ , that we capture in Figure 3. We chose the best  $\mathcal{C}_0$  and  $\mathcal{C}_1$  from these for all our experiments.

#### 5.4.2 Effect of Edge Features

Table 2 shows the effect of different edge features in a collective setting. The model seems to benefit the most from the inlink and outlink relatedness features, while context overlap-based features seem to be noisy. This is understandable as context overlap-based signals are useful only for topically coherent entities, which might not hold true for an aggressively tagged corpus like ours (Kulkarni et al., 2009).

Edge feature	Disambiguator only			Weak annotation match		
	P	R	F	P	R	F
Category	.72	.74	.73	.72	.63	<b>.67</b>
<b>Outlink (O)</b>	.84	.67	<b>.74</b>	.84	.57	<b>.68</b>
<b>Inlink (I)</b>	.80	.73	<b>.76</b>	.80	.62	<b>.70</b>
Frequent (F)	.84	.64	.73	.84	.54	.66
Synopsis	.69	.61	.65	.69	.52	.59
Syn. V/Adj.	.69	.67	.68	.69	.57	.62
Full text	.85	.63	.73	.85	.54	.66
All features	.44	.50	.47	.44	.42	.43
<b>I+O</b>	.85	.67	<b>.74</b>	.85	.56	<b>.68</b>
<b>I+O+F</b>	.79	.74	<b>.76</b>	.79	.63	<b>.70</b>

Table 2: Effect of edge features: two-fold cross validation on *Wiki<sub>cur</sub>*. Edge features that showed improvement over node features are shown in bold.

#### 5.4.3 Does training help?

We sampled 50 documents from the *Wiki<sub>cur</sub>* dataset, 5 at a time and used them for training, applying both passive (PL) and active learning (AL). The  $F_1$  measure evaluated on an independent test set of 30 documents is shown in the plot (Refer to figure 4). The  $F_1$  on training set seems to fluctuate, more so for *Train-AL* (Chen et al., 2006), but the  $F_1$  on test set does show a steady improvement.

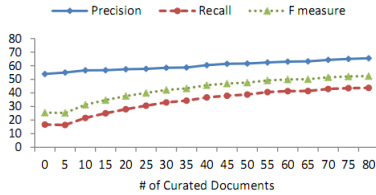


Figure 2: Effect of data curation

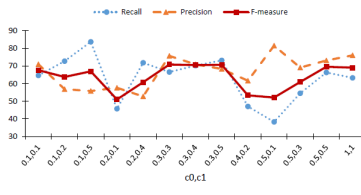


Figure 3: Effect of varying  $C_0$  and  $C_1$

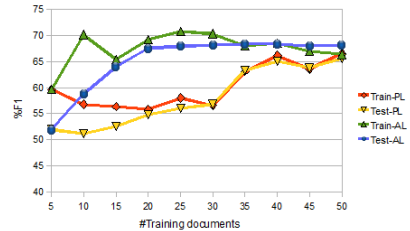


Figure 4: Effect of training

Annotator	<i>IITB<sub>part</sub></i>			AQUAINT			MSNBC		
	F	P	R	F	P	R	F	P	R
AIDA	.07	.66	.04	.21	.35	.15	.47	.75	.35
Wikify!	.37	.55	.28	.34	.29	.42	.41	.34	.51
TagMe	.44	.45	.42	.51	.46	.57	.52	.48	.55
Wikipedia Miner	.52	.57	.48	.47	.38	.63	.46	.55	.36
Illionis Wikifier	.44	.58	.36	.34	.29	.42	.41	.34	.51
<b>Our Model (Node+I)</b>	<b>.67</b>	<b>.76</b>	<b>.60</b>	<b>.78</b>	<b>.81</b>	<b>.74</b>	<b>.67</b>	<b>.68</b>	<b>.66</b>
<b>Our Model (Node+I+O+F)</b>	<b>.65</b>	<b>.69</b>	<b>.61</b>	<b>.79</b>	<b>.82</b>	<b>.75</b>	<b>.66</b>	<b>.63</b>	<b>.69</b>

Table 3: Comparison with publicly available systems (as reported by (Cornolti et al., 2013)) on three datasets

#### 5.4.4 Comparison with collective approaches

We compared our system against several other collective annotation approaches: AIDA (Hofart et al., 2011), Wikify! (Mihalcea and Csomai, 2007), TagMe (Ferragina and Scaiella, 2010), Wikipedia Miner (Milne and Witten, 2008) and Illionis Wikifier (Ratinov et al., 2011) on three datasets *viz.* *IITB<sub>part</sub>*, AQUAINT (Wikipedia Miner) and MSNBC (Cucerzan, 2007). Our system consistently beats all these systems on all the three datasets (Table 3). Some of the other collective annotation systems like Cucerzan ( $F_1 : .45$ ), CSAW (Kulkarni et al., 2009) ( $F_1 : .69$ ), (Han et al., 2011) ( $F_1 : .73$ ), and (Han and Sun, 2012) ( $F_1 : .8$ ) have used CSAW’s evaluation measure to evaluate on *IITB<sub>part</sub>*. We achieved an  $F_1$  of 0.6 using the same measure. The relatively lower  $F_1$  on this dataset could be attributed to inconsistencies between the ground truth and our knowledge base. During our manual annotation of *IITB<sub>part</sub>*, we came across over 8000 annotations that were either added or removed<sup>5</sup> to create the *IITB<sub>cur</sub>* dataset.

We evaluated our system on the ERD dataset and achieved R:.62, P:.66,  $F_1 : .64$ . We believe that our system benefits from model training, thereby performing better than that of (Kulkarni

et al., 2014) ( $F_1 : .61$ ). While some of the other systems (Cornolti et al., 2014) at ERD performed better, this could be attributed to their choice of features. Our system offers an end-end annotation framework that is interactive and jointly trains feature weights.

#### 5.5 Results on *IITB<sub>cur</sub>*

Section 6 shows some examples of incomplete annotations in the *IITB<sub>part</sub>* dataset. It is precisely such cases that we tried to correct during data preparation. Finally, we report the accuracy of our model on the *IITB<sub>cur</sub>* dataset -  $P : 77.4\%$ ,  $R : 54.3\%$ ,  $F_1 : 63.8\%$ .

#### 5.6 Performance Evaluation

While our model allows for efficient inference and learning, graph construction itself is an expensive operation. For a document with  $|E_d|$  candidate entities, the graph construction complexity is  $O(|E_d|^2)$ . For documents in the *Wiki<sub>cur</sub>* set with 190 candidate entities on an average, the average graph construction time was about 57 seconds. For the relatively larger documents in the *IITB<sub>cur</sub>* dataset, the average graph construction time was around 1.5 minutes. The performance could be improved by (a) pre-computing the entity-entity features for all entities in the knowledge base (b) dividing input document into chunks and performing graph construction and inference in parallel.

The running time for inference (Figure 5) shows a slightly quadratic behavior in the number of candidate entities  $|E_d|$  of a document. Inference on most documents runs in under 0.5 seconds. On the relatively sparser *Inlink+Outlink* graphs (Refer to Figure 6), training is much faster than the more dense *Category* graphs. The faster training happens without trading off much on accuracy as can be seen in Table 2. For our experiments, the model was retrained at time  $t$  using all the available train-

<sup>5</sup>due to erroneous annotations or newer Wikipedia dump

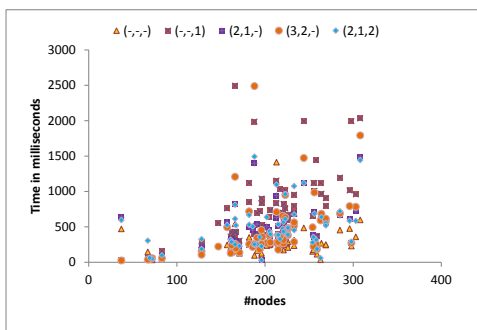


Figure 5: Running time for inference on *Wiki<sub>cur</sub>*

ing data. While this might be acceptable for offline training, online systems might benefit from faster incremental training approaches.

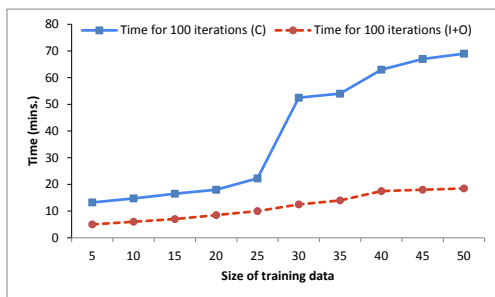


Figure 6: Scalability of training

## 6 Challenges with data curation

Data curation is a tedious and challenging task. Its inherent ambiguity often introduces annotator bias leading to either incomplete or ambiguous annotations in the curated data.

1. There might be cases when two or more entities are correct as attachments for a mention. *E.g.*, mention ‘Barack Obama’ can be tagged as *Barack Obama* or *President of United States*, and both might seem correct in the context that it appeared. A ‘one entity per mention’ assumption makes it impossible to honor such cases.
2. Human annotators often limit their attention to the candidate entities retrieved by the spotter and very rarely search the catalog for any missed candidates. This results in a lot of

missed annotations and often many mentions getting no attachments (NA).

3. Annotators also seem biased towards entity names that match with the mention text. However, this is often not true. *E.g.* a mention of ‘cone snail’ disambiguates to *Conidae* and *Conus*.
4. Wikipedia contains many disambiguation pages that often show up in the candidate set for a mention. Tagging a mention with a disambiguation page seems to beat the very purpose of a disambiguation system. Ideally, the mention should be annotated with one of the entities on the disambiguation page or NA if none of them is semantically right.

Table 4 shows some of these cases from the *IITB<sub>part</sub>* dataset. It is cases like these that we attempted to correct in coming up with the curated *IITB<sub>cur</sub>* dataset.

## 7 Conclusion

We presented an approach to jointly train the node and edge features of a collective disambiguation model for the purpose of entity linking. Our system leverages active learning to bring down labeling effort. Experiments show that the model benefits from training and improves with the availability of more labeled data. We consistently performed better than many other systems on various datasets. It also scales reasonably well and with suggested tweaks can be used for large scale document annotation.

## References

- Indrajit Bhattacharya and Lise Getoor. 2006. A latent dirichlet model for unsupervised entity resolution. In *SIAM INTERNATIONAL CONFERENCE ON DATA MINING*.
- Y. Boykov and V. Kolmogorov. 2004. An experimental comparison of min-cut/max-flow algorithms for energy minimization in vision. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 26(9):1124–1137.
- Razvan Bunescu and M Pasca. 2006. Using encyclopedic knowledge for named entity disambiguation. *Proceedings of EACL*, pages 9–16.
- Soumen Chakrabarti, Kriti Puniyani, and Sujatha Das. 2006. Optimizing scoring functions and indexes for proximity search in type-annotated corpora. In



Ground Mention	Ground Entity	Predicted Entity	Remarks
lifestyle	Lifestyle	NA	Disambiguation page attachment
harsh reality	NA	Reality	Incomplete data
effort	NA	Energy	Incomplete data
self discipline	Discipline	self → Self, discipline → Discipline	Overlapping mentions
god	God (male deity)	God	Multiple correct entities
intellect	Intelligence	Intellect	Multiple correct entities

Table 4: Examples of predictions on the *IITB<sub>part</sub>* highlighting the challenges in data curation

- Proceedings of the 15th international conference on World Wide Web, WWW '06*, pages 717–726, New York, NY, USA. ACM.
- Jinying Chen, Andrew Schein, Lyle Ungar, and Martha Palmer. 2006. An empirical study of the behavior of active learning for word sense disambiguation. In *Proceedings of the Main Conference on Human Language Technology Conference of the North American Chapter of the Association of Computational Linguistics, HLT-NAACL '06*, pages 120–127, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Tao Cheng, Xifeng Yan, and Kevin C. Chang. 2007. Entityrank: searching entities directly and holistically. In *VLDB '07: Proceedings of the 33rd international conference on Very large data bases*, pages 387–398. VLDB Endowment.
- Marco Cornolti, Paolo Ferragina, and Massimiliano Ciaramita. 2013. A framework for benchmarking entity-annotation systems. In *Proceedings of the 22nd international conference on World Wide Web, WWW '13*, pages 249–260, Republic and Canton of Geneva, Switzerland. International World Wide Web Conferences Steering Committee.
- Marco Cornolti, Paolo Ferragina, Massimiliano Ciaramita, Hinrich Schütze, and Stefan Rüd. 2014. The smaph system for query entity recognition and disambiguation. In *Proceedings of the First International Workshop on Entity Recognition & Disambiguation, ERD '14*, pages 25–30, New York, NY, USA. ACM.
- Silviu Cucerzan. 2007. Large-scale named entity disambiguation based on wikipedia data. In *Proceedings of EMNLP-CoNLL*, volume 6, pages 708–716.
- Stephen Dill, Nadav Eiron, David Gibson, Daniel Gruhl, R. Guha, Anant Jhingran, Tapas Kanungo, Sridhar Rajagopalan, Andrew Tomkins, John A. Tomlin, and Jason Y. Zien. 2003. Semtag and seeker: bootstrapping the semantic web via automated semantic annotation. In *Proceedings of the 12th international conference on World Wide Web, WWW '03*, pages 178–186, New York, NY, USA. ACM.
- Angela Fahrni and Michael Strube. 2012. Jointly disambiguating and clustering concepts and entities with markov logic. In *COLING*, pages 815–832. 227
- Paolo Ferragina and Ugo Scaiella. 2010. Tagme: on-the-fly annotation of short text fragments (by wikipedia entities). In *Proceedings of the 19th ACM international conference on Information and knowledge management, CIKM '10*, pages 1625–1628, New York, NY, USA. ACM.
- William A. Gale, Kenneth W. Church, and David Yarowsky. 1992. One sense per discourse. In *Proceedings of the workshop on Speech and Natural Language, HLT '91*, pages 233–237, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Xianpei Han and Le Sun. 2012. An entity-topic model for entity linking. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning, EMNLP-CoNLL '12*, pages 105–115, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Xianpei Han and Jun Zhao. 2009. Named entity disambiguation by leveraging wikipedia semantic knowledge. In *Proceedings of the 18th ACM conference on Information and knowledge management, CIKM '09*, pages 215–224, New York, NY, USA. ACM.
- Xianpei Han, Le Sun, and Jun Zhao. 2011. Collective entity linking in web text: a graph-based method. In *Proceedings of the 34th international ACM SIGIR conference on Research and development in Information Retrieval*, pages 765–774. ACM.
- Johannes Hoffart, Mohamed Amir Yosef, Ilaria Bordino, Hagen Fürstenau, Manfred Pinkal, Marc Spaniol, Bilyana Taneva, Stefan Thater, and Gerhard Weikum. 2011. Robust disambiguation of named entities in text. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing, EMNLP '11*, pages 782–792, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Gjergji Kasneci, Fabian M. Suchanek, Georgiana Ifrim, Shady Elbassuoni, Maya Ramanath, and Gerhard Weikum. 2008. Naga: harvesting, searching and ranking knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data, SIGMOD '08*, pages 1285–1288, New York, NY, USA. ACM.
- Saurabh S. Kataria, Krishnan S. Kumar, Rajeev R. Rastogi, Prithviraj Sen, and Srinivasan H. Sengamedu. 2011. Entity disambiguation with hierarchical topic

- models. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '11, pages 1037–1045, New York, NY, USA. ACM.
- Sayali Kulkarni, Amit Singh, Ganesh Ramakrishnan, and Soumen Chakrabarti. 2009. Collective annotation of Wikipedia entities in web text. *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD '09*, page 457.
- Ashish Kulkarni, Kanika Agarwal, Pararth Shah, Sunny Raj Rathod, and Ganesh Ramakrishnan. 2014. System for collective entity disambiguation. In *Proceedings of the First International Workshop on Entity Recognition & Disambiguation*, ERD '14, pages 111–118, New York, NY, USA. ACM.
- David D. Lewis and Jason Catlett. 1994. Heterogeneous uncertainty sampling for supervised learning. In *In Proceedings of the 11th International Conference on Machine Learning (ICML)*, pages 148–156. Morgan Kaufmann.
- Mingkun Li and Ishwar K. Sethi. 2006. Confidence-based active learning. *IEEE Trans. Pattern Anal. Mach. Intell.*, 28(8):1251–1261, August.
- Xiaonan Li, Chengkai Li, and Cong Yu. 2010. Entity-engine: answering entity-relationship queries using shallow semantics. In *Proceedings of the 19th ACM international conference on Information and knowledge management*, CIKM '10, pages 1925–1926, New York, NY, USA. ACM.
- Thomas Lin, Mausam, and Oren Etzioni. 2012. Entity linking at web scale. In *Proceedings of the Joint Workshop on Automatic Knowledge Base Construction and Web-scale Knowledge Extraction*, AKBC-WEKEX '12, pages 84–88, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Paul McNamee. 2009. Overview of the tac 2009 knowledge base population track.
- Rada Mihalcea and Andras Csomai. 2007. Wikify!: linking documents to encyclopedic knowledge. In *Proceedings of the sixteenth ACM conference on Conference on information and knowledge management*, pages 233–242. ACM.
- David Milne and Ian H Witten. 2008. Learning to link with wikipedia. In *Proceedings of the 17th ACM conference on Information and knowledge management*, pages 509–518. ACM.
- Lev Ratinov, Dan Roth, Doug Downey, and Mike Anderson. 2011. Local and global algorithms for disambiguation to wikipedia. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies - Volume 1*, HLT '11, pages 1375–1384, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Ben Taskar and Daphne Koller. 2001. Learning Associative Markov Networks.
- B. Taskar, V. Chatalbashev, and D. Koller. 2004. Learning associative markov networks. In *Proceedings of the twenty-first international conference on Machine learning*, page 102. ACM.
- P. Vernaza, B. Taskar, and D.D. Lee. 2008. Online, self-supervised terrain classification via discriminatively trained submodular markov random fields. In *Robotics and Automation, 2008. ICRA 2008. IEEE International Conference on*, pages 2750–2757. IEEE.
- Ben Wellner, Andrew McCallum, Fuchun Peng, and Michael Hay. 2004. An integrated, conditional model of information extraction and coreference with application to citation matching. In *Proceedings of the 20th Conference on Uncertainty in Artificial Intelligence*, UAI '04, pages 593–601, Arlington, Virginia, United States. AUAI Press.
- Michael L. Wick, Aron Culotta, Khashayar Rohanimanesh, and Andrew McCallum. 2009. An entity based model for coreference resolution. In *SDM*, pages 365–376.
- Yiping Zhou, Lan Nie, Omid Rouhani-Kalleh, Flavian Vasile, and Scott Gaffney. 2010. Resolving surface forms to wikipedia topics. In *Proceedings of the 23rd International Conference on Computational Linguistics*, COLING '10, pages 1335–1343, Stroudsburg, PA, USA. Association for Computational Linguistics.