

A Unified Probabilistic Framework for Name Disambiguation in Digital Library

Jie Tang, A.C.M. Fong, Bo Wang, and Jing Zhang

Abstract—Despite years of research, the name ambiguity problem remains largely unresolved. Outstanding issues include how to capture all information for name disambiguation in a unified approach, and how to determine the number of people K in the disambiguation process. In this paper, we formalize the problem in a unified probabilistic framework, which incorporates both attributes and relationships. Specifically, we define a disambiguation objective function for the problem and propose a two-step parameter estimation algorithm. We also investigate a dynamic approach for estimating the number of people K . Experiments show that our proposed framework significantly outperforms four baseline methods of using clustering algorithms and two other previous methods. Experiments also indicate that the number K automatically found by our method is close to the actual number.

Index Terms—Digital libraries, information search and retrieval, database applications, heterogeneous databases.

1 INTRODUCTION

DIFFERENT people have the same name. For example, in the DBLP database, there are 300 people with the name "Jie Tang". The number of people with the same name is 78.74% of the total number of people in the database. In this paper, we propose a unified probabilistic framework for name disambiguation in digital libraries. We formalize the problem in a unified probabilistic framework, which incorporates both attributes and relationships. Specifically, we define a disambiguation objective function for the problem and propose a two-step parameter estimation algorithm. We also investigate a dynamic approach for estimating the number of people K . Experiments show that our proposed framework significantly outperforms four baseline methods of using clustering algorithms and two other previous methods. Experiments also indicate that the number K automatically found by our method is close to the actual number.

1.1 Motivation

We believe that name disambiguation is a challenging problem in digital libraries. In this paper, we propose a unified probabilistic framework for name disambiguation in digital libraries. We formalize the problem in a unified probabilistic framework, which incorporates both attributes and relationships. Specifically, we define a disambiguation objective function for the problem and propose a two-step parameter estimation algorithm. We also investigate a dynamic approach for estimating the number of people K . Experiments show that our proposed framework significantly outperforms four baseline methods of using clustering algorithms and two other previous methods. Experiments also indicate that the number K automatically found by our method is close to the actual number.

- J. Tang and J. Zhang are with the Department of Computer Science and Technology, Tsinghua University, Rm 1-308, FIT Building, Beijing 100084, China. E-mail: jietang@tsinghua.edu.cn, zhangjing0544@gmail.com.
- A.C.M. Fong is with the School of Computing and Mathematical Sciences, Auckland University of Technology, AUT Tower Level 1, 2-14 Wakefield Street, Auckland 1142, New Zealand. E-mail: afong@aut.ac.nz.
- B. Wang is with the Department of Computer Science, Nanjing University of Aeronautics and Astronautics, Nanjing 210016, China. E-mail: bowang@nuaa.edu.cn.

Manuscript received 1 July 2008; revised 5 Apr. 2010; accepted 16 Nov. 2010; published online 27 Dec. 2010.

Recommended for acceptance by B.C. Ooi.

For information on obtaining reprints of this article, please send e-mail to: tkde@computer.org, and reference IEEECS Log Number TKDE-2008-07-0335. Digital Object Identifier no. 10.1109/TKDE.2011.13.

... a abe e ee ... e e f e ea ... (cf. Sec 2.1 f def f e ea ... e). T e d a ce be ee ... de de e e a f ... e a e e f e c e -ba ed a ... ea e e (e., c e a). T e d ... e e dea d a b a e., c d ca e ... a 11 a e d be a ed. ee d ffe e a ... A ed a e b e a f F.1 a a e d ba ed c e a (e d a ce) d be d f f c a c e e a fac e f a ce, a d a d ffe e e f e a ca be e f, b d ffe e de ee f c b. F e a e, ee a C A e a be ee de #3 a d #8. A e a be ee e de #3 a d #8. A e f e C A e a e ca a e de (a e) e a e a. O e c a, a ee a C a e a be ee de #3 a d #7, e a e a e a ed d ffe e a. T e c a e e e e de a a f e a e d a b a be b c de b a b e f a f e de a d e ea be ee de.

1.2 Prior Work

T e be a bee de e de e a ed d ffe e d a, a d a a e e [4], [5], [7], e b a e a c e d a b a [3], [20], a e de f ca [26], a d Ob e c d c [49]. De e a a a c e ed, e a e a b be e a a e e ed. I e e a, e e d f a e d a b a a fa ee ca e e: supervised based, unsupervised based, a d constraint based. T e e ed-ba ed a a c (e., [17]) e ea a e c f c c a f ca de f e a c a a e f e a -abe ed a da a. T e, e ea ed de ed c ea a e f e a c a e. I e e ed-ba ed a a c (e., [18], [36], [37], [49]), c e a c de a e e ed f d a e

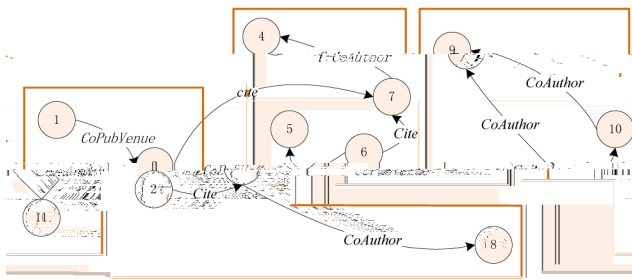


Fig. 1. An example of name disambiguation.

a d a e d ffe e a a e a ed
d ffe e a . T e c a -ba ed a a c a
7 e e c e a . T e d ffe e c e a
e - ded c a a e ed de e c e
a a d b e e d a a a (e. , [2], [51]).
F e e e e a a a c e b a e d e ,
c a /a a , a d c b a f e d ffe e
a a c e a e b e e d e d . F e a e , W a e a .
[47] d c e a e a e e -ba ed a a c e e e
e c e e c e e d a b a e a d d e e
a d e f a e e e c e a e e
e . D a e a . [11], a e d e e e d a e a c e e
c e a e e c a e e c c e c e f a e d
e e a e e . T e e d e f
e f e e c e a e a b e c (e. , a a c a b d)
e e a e d a e f a b e c a d a c e c .
McRae-S e c e a d S a d b [28] e e a a -ba ed
a a c a d a b a a e -c a e c a
e b e f -c a , c a e a . T e
a a c a a c e e a , e c b a e a e
e c a . Y e a . [50], a e d e e e d e e d a a c e
d e f e f f f a b a b b e a
e c e e a e a . M e e c e , C e e a . [8] d
c b e e d ffe e d a b a a a c e
a d e a e e e e e b e f a e ,
c c b e e e f e b a e -e e e
e e a e e e e e e e
a c c a c f e e . W a e a . [46] e a
e a e b c f a e e e e e e f
b c a e e f e c e d b e e c e e d b c . O
a d L e e [32] d e c a a b e f e a e
d a b a b e . A c e e a
b e e a d e , e d d a c e e a f a c
d a b a e d e e a :

1. S e e a , c e e , d (e. , [31], [35], [48]) f c a e d a a a , b a e d e c a c e ; e e e , d (e. , [18], [42]) a c e e d a a a , a c c d e d e a . A f e e e a c e (e. , [38], [52]) c b e e e c e f f a . F e a e , Z e a . a e c b e f a b a e d b e e a b e (e. , d e a) a d a c a c e b f c c a a b e a e e d a , e c a e f (a b e , a e) a e c e , a d b e e e a e e e e c e ' c e e a e b , d

a d a d e . T e a e c a
e a a e b e e d c a d c a
f a . A e e a e e a b e
d e a e a a b e a c e a e e
c e e f a e e c e e d a c e
e a e , a b a a c e , e c b
f e d f f e e f a a e b e
T e a e a b e c c d e , a a d d a b e
a f a e c e e b e c e
d e a d e e a c e c e e . F e e ,
[52], e e e e a d a a e c a e f e
a b e . T e f d a a e (c a b) a
e (b a) a b e a d e e c d a a e f
D B L P b b a c a d a a a a b e .
W e a e a c c e d e a b e f a
e e d f a c e a e d a b a
b e e f f e c e .

2. T e e f a c e f a e a f e e e d e d
d e e d a c c a e e a K . A
e e a c e a c a X - e a [33]
c a a a c a f d e b e K b a e d e
c e , c e a e e c a
e , d c a b e d e c a e d e a e
d a b a b e .
3. I e e , e d , e d a a a c a
e e d e a d e a ; e
b e e , e e a b e e d f f e e
e a (e. , C A a d C a) b e e e
d e . T e e f d f f e e a a a e
d f f e a c e f e a e d a b a
b e . H a a c a d e e d e e f
c b f d f f e e a a a
c a e b e .

1.3 Our Solution

H a c d c e d a e a , e e a
f e d b a b c f a e a d d e e a b e
c a e e . S e c f c a , e f a e e d a b a
b e a M a R a d F e d (M R F) [16], [24],
c e d a a a e c e e b c a a b e a d
e a . W e e e a d a c a a c f e a
e b e f e e K a d a - e a f
a a e e e a . T e e d a a c a a c e e
b e e e f a c e a e d a b a a e
e , d b e c a e e a a c a e a d a a e f e
d e e d e c e b e e a e a e . T e b e f
e d e , e f f a e a e
b e f a e d a b a a f e d f a e
a d a c e e b e e e .
T e e d f a e e e e a . O e c a
c a e a e a a f e a e c a f e a e e
f a e , e. , a f e a e b a e d e e b e a c e e
e d . T e f a e c a b e a e e d e d d e a
a e b e c a e e e a
e a a d a b a e [4].
O c b a e c d e : 1) f a
f e a e d a b a b e a f e d b a b -
c f a e ; 2) a f a e e
a a e e e a e f a e ; a d 3) a e c a
e f c a f e e f f e c e e f e e d f a e .

TABLE 1
Attributes of Each Publication p_i

Attribute	Description
$p_i.title$	title of p_i
$p_i.pubvenue$	published conference/journal of p_i
$p_i.year$	published year of p_i
$p_i.abstract$	abstract of p_i
$p_i.authors$	authors name set of p_i $\{a_i^{(0)}, a_i^{(1)}, \dots, a_i^{(u)}\}$
$p_i.references$	references of p_i

TABLE 2
Relationships between Papers

R	W	Relation Name	Description
r_1	w_1	CoPubVenue	$p_i.pubvenue = p_j.pubvenue$
r_2	w_2	CoAuthor	$\exists r, s > 0, a_i^{(r)} = a_j^{(s)}$
r_3	w_3	Citation	$p_i.cites p_j$ or $p_i.cites_{D>D_{th}} p_j$
r_4	w_4	Constraint	feedback supplied by users
r_5	w_5	τ -CoAuthor	τ -extension co-authorship ($\tau > 1$)

2 PROBLEM FORMALIZATION

2.1 Definitions

In this section, we define the problem of name disambiguation. Table 1 shows the attributes of each publication p_i . We describe the author name that we are going to disambiguate as the principle author $a_i^{(0)}$ and the rest (if any) as secondary authors.

Definition 1 (Principle Author and Secondary Author).

Each paper p_i has one or more authors $A_{p_i} = \{a_i^{(0)}, a_i^{(1)}, \dots, a_i^{(u)}\}$. We describe the author name that we are going to disambiguate as the principle author $a_i^{(0)}$ and the rest (if any) as secondary authors.

We define the following relationships between papers (Table 2). Section 2.1.1 defines the following relationships between papers:

- **CoPubVenue** (r_1) p_i and p_j share the same published venue. Formally, $p_i.pubvenue = p_j.pubvenue$. We denote this relationship as $p_i \sim_{r_1} p_j$.
- **CoAuthor** (r_2) p_i and p_j share at least one author. Formally, $A_{p_i} \cap A_{p_j} \neq \emptyset$. We denote this relationship as $p_i \sim_{r_2} p_j$.
- **Citation** (r_3) p_i cites p_j . Formally, $p_i.cites p_j$ or $p_i.cites_{D>D_{th}} p_j$. We denote this relationship as $p_i \sim_{r_3} p_j$.
- **Constraint** (r_4) p_i and p_j are related by a constraint. Formally, p_i and p_j are related by a constraint. We denote this relationship as $p_i \sim_{r_4} p_j$.
- **τ -CoAuthor** (r_5) p_i and p_j share at least τ authors. Formally, $|A_{p_i} \cap A_{p_j}| \geq \tau$. We denote this relationship as $p_i \sim_{r_5} p_j$.

The following definitions describe the relationships between papers. We define the following relationships between papers:

CoPubVenue (r_1) p_i and p_j share the same published venue. Formally, $p_i.pubvenue = p_j.pubvenue$. We denote this relationship as $p_i \sim_{r_1} p_j$.

CoAuthor (r_2) p_i and p_j share at least one author. Formally, $A_{p_i} \cap A_{p_j} \neq \emptyset$. We denote this relationship as $p_i \sim_{r_2} p_j$.

Citation (r_3) p_i cites p_j . Formally, $p_i.cites p_j$ or $p_i.cites_{D>D_{th}} p_j$. We denote this relationship as $p_i \sim_{r_3} p_j$.

Constraint (r_4) p_i and p_j are related by a constraint. Formally, p_i and p_j are related by a constraint. We denote this relationship as $p_i \sim_{r_4} p_j$.

τ -CoAuthor (r_5) p_i and p_j share at least τ authors. Formally, $|A_{p_i} \cap A_{p_j}| \geq \tau$. We denote this relationship as $p_i \sim_{r_5} p_j$.

Definition 2 (Cluster Atom).

A cluster atom is a cluster in which papers are closely connected (e.g., the similarity $K(x_i, x_j) > threshold$). Papers with similarity less than the threshold will be assigned to disjoint cluster atoms.

The following definitions describe the relationships between papers. We define the following relationships between papers:

CoPubVenue (r_1) p_i and p_j share the same published venue. Formally, $p_i.pubvenue = p_j.pubvenue$. We denote this relationship as $p_i \sim_{r_1} p_j$.

CoAuthor (r_2) p_i and p_j share at least one author. Formally, $A_{p_i} \cap A_{p_j} \neq \emptyset$. We denote this relationship as $p_i \sim_{r_2} p_j$.

Citation (r_3) p_i cites p_j . Formally, $p_i.cites p_j$ or $p_i.cites_{D>D_{th}} p_j$. We denote this relationship as $p_i \sim_{r_3} p_j$.

Constraint (r_4) p_i and p_j are related by a constraint. Formally, p_i and p_j are related by a constraint. We denote this relationship as $p_i \sim_{r_4} p_j$.

τ -CoAuthor (r_5) p_i and p_j share at least τ authors. Formally, $|A_{p_i} \cap A_{p_j}| \geq \tau$. We denote this relationship as $p_i \sim_{r_5} p_j$.

2.2 Name Disambiguation

The goal of name disambiguation is to identify the authors of a set of papers $P = \{p_1, p_2, \dots, p_n\}$. The authors of a paper p_i are $A_{p_i} = \{a_i^{(0)}, a_i^{(1)}, \dots, a_i^{(u)}\}$. We define the following relationships between papers:

$f, e, -ca, ed, f, a, e, a, [13], e, e, e, e, f, a, e, b, c, e, -ba, ed, f, a, a, d, c, e,$
 $b, ca, da, a, P, b, ca, a, d, ea, a, e, a, -$
 $f, ed, a, d, ec, ed, a, c, ea, c, de$
 $e, e, e, a, a, e, a, de, ac, ed, ea, ea, A, b, e,$
 $f, a, a, e, a, e, a, ac, ed, e, c, e, d, de, a, a,$
 $fea, e, ec, F, e, ec, e, e, d, (afe,$
 $d, f, e, a, d, e,)$, $ea, b, e, fa, a, e, a,$
 $fea, e, a, d, e, e, be, f, e, cc, e, ce, a, e,$
 $a, e, F, a, e, ca, def, e, e, b, ca, f, a, e,$
 $a, a, f, :$

Definition 3 (Publication Informative Graph). Given a set of papers $P = \{p_1, p_2, \dots, p_n\}$, let $r_k(p_i, p_j)$ be a relationship r_k between p_i and p_j . A publication informative graph is a graph $G = (P, R, V_P, W_R)$, where each $v(p_i) \in V_P$ corresponds to the feature vector of paper p_i and $w_k \in W_R$ denotes the weight of relationship r_k . Let $r_k(p_i, p_j) = 1$ iff there is a relationship r_k between p_i and p_j ; otherwise, $r_k(p_i, p_j) = 0$.

$S, e, e, e, a, e, K, e, \{y_1, \dots, y_K\}, e, a, e,$
 $a, a, d, a, b, a, e, e, n, b, ca, e, ea,$
 $e, ea, c, e, y_i, i \in [1, K], M, e, ec, f, ca, e, a, a, f,$
 $a, ed, a, b, a, ca, be, def, ed, a, :$

1. $F, a, ed, a, b, a, be, T, ef,$
 $a, a, eed, c, de, b, ca, a, b, e,$
 $fea, e, a, ca, ed, ea, c, a, e, a, d, ea,$
 $-$
 $be, ee, a, e, .$
 2. $S, e, be, a, c, ed, a, ac, Ba, ed,$
 $ef, a, ea, c, ed, a, ac,$
 $a, d, e, a, eff, ce, a, .$
 3. $Dee, e, be, f, e, e, K, G, e, a,$
 $d, a, b, a, a, (, a, f, a,$
 $), de, e, e, e, ac, a, K, .$
- $I, a, ef, e, e, a, . F,$
 $ed, ae, cea, f, a, e, e, e, ed, a, b,$
 $a, be, a, fed, fa, e, . Sec, d, e,$
 $a, de, e, , Ma, Ra, d, Fed, [16], a, e, a,$
 $a, ed, de, ea, a, da, a, H, e, e, e,$
 $b, ca, f, a, e, a, e, a, e, be,$
 $a, b, a, c, ec, ed, b, d, ffe, e, f, ea,$
 $I, cea, ef, fe, ce, (, a, a, e, e,$
 $e, a,)$, $c, a, a, a, b, a, c, e, I,$
 $add, e, a, e, be, f, e, e, K, a, a,$
 $c, a, e, a, .$

3 OUR FRAMEWORK

3.1 Basic Idea

$We, a, e, b, a, c, b, e, a, f, e, a, ed, a, b, a,$
 $be, :1) a, e, a, c, e, ed, a, e,$
 $e, a, e, a, be, (be, e, a, ea,); a, d, 2) a, e,$
 $a, ea, ed, a, e, e, a, e, a, be, f,$
 $e, a, e, a, e, ca, a, a, a,$
 $e, a, e, . A, dea, d, a, b, a, e, e,$
 $a, e, b, e, ea, b, c, e, a, a, d, a, e,$
 $ea, . T, a, a, be, , beca, e,$
 $e, c, e, e, d, ca, e, ba, a, ce, e,$
 $ece, f, f, a, .$
 $I, a, e, e, ea, fed, fa, e, ba, ed,$
 $Ma, Ra, d, Fed, [16], [24]. M, e, acc, a, e, e,$

$f, a, e, b, c, e, -ba, ed, f, a, a, d, c, e,$
 $-ba, ed, f, a, a, H, d, de, Ma, Ra, d, Fed,$
 $(HMRF) de, a, fea, e, f, c, . T, e, c, b,$
 $de, ee, f, e, e, f, f, a, a, ef, a, ed, a,$
 $e, f, e, fea, ef, c, . T, e, a, ce, f, d, ffe,$
 $e, f, ea, a, de, ed, a, e, f,$
 $c, e, d, fea, ef, c, . S, e, HMRF, de,$
 $c, de, b, e, a, e, e, f, fea, ef, c,$
 $a, da, a, e, d, ffe, e, . S, c, a, f, a, e,$
 $a, ffe, add, a, ad, a, a, e, : f,$
 $e, ed, ea, , e, ed, ea, , a, d, e,$
 $e, ed, ea, . I, a, e, e, f, c,$
 $e, ed, ea, f, a, ed, a, b, a, , b,$
 $ea, c, a, e, e, /, e, ed, f, a,$
 $e, de, Sec, d, a, a, d, de, eec, e,$
 $HMRF, de, T, e, bec, ef, c, e, HMRF, de,$
 $a, e, bab, d, b, f, d, de, a, abe, e,$
 $b, e, a, , c, ac, e, f, de, eec, a, e, .$

3.2 Hidden Markov Random Fields

$A, Ma, Ra, d, Fed, a, c, d, a, bab,$
 $d, b, f, abe, (, d, de, a, abe), a, be, e,$
 $Ma, e, [16]. Ma, eca, ca, e, f, MRF, ca, be,$
 $de, e, ed, A, H, d, de, Ma, Ra, d, Fed, a, e, be,$
 $f, e, fa, f, MRF, a, d, c, ce, de, ed, f,$
 $H, d, de, Ma, M, de, (HMM) [15]. A, HMRF,$
 $a,$
 $c, ed, f, eec, e, : a, be, abe, e, f,$
 $a, d, a, abe, X = \{x_i\}_{i=1}^n, a, d, de, fed, f, a, d,$
 $a, abe, Y = \{y_i\}_{i=1}^n, a, d, e, b, d, be, ee, ea,$
 $a, f, a, abe, e, d, de, fed, .$
 $We, f, a, e, ed, a, b, a, be, a, a, f,$
 $ea, a, a, e, d, ffe, e, c, e, . Le, e,$
 $d, de, a, abe, Y, be, e, c, e, abe, e, a, e, .$
 $E, e, d, de, a, abe, y_i, a, e, a, a, ef, e, e,$
 $\{1, \dots, K\}, c, a, e, e, de, e, f, e, c, e, . T, e,$
 $b, e, a, a, abe, X, c, e, d, a, e, , e, e, e,$
 $a, d, a, abe, x_i, e, e, a, ed, f, a, c, d, a,$
 $bab, d, b, P(x_i|y_i), de, e, ed, b, e, c, e,$
 $d, d, de, a, abe, y_i, F, e, e, a, d, a, abe,$
 $X, a, ea, ed, be, e, e, a, ed, c, d, a, de, e, de,$
 $f, e, d, de, a, abe, Y, e, ,$

$$P(X|Y) = \prod_{x_i \in X} P(x_i|y_i). \tag{1}$$

$F, . 2, e, a, ca, c, e, f, e, HMRF, f, e,$
 $e, a, e, F, . 1. We, ee, a, de, e, de, ed, e, a, e,$
 $ded, be, ee, e, d, de, a, abe, c, e, d,$
 $e, ea, F, . 1. T, e, a, e, f, ea, c, d, de,$
 $a, abe, (e, , y_1 = 1) de, e, ea, e, e, . We, d,$
 $de, e, d, ec, ea, be, ee, e, b,$
 $b, e, de, ca, a, a, e, de, e, de, ce, a,$
 $e, ea, .$
 $A, HMRF, a, eca, ca, e, f, MRF, e, bab,$
 $d, b, f, e, d, de, a, abe, be, e, Ma,$
 $e, . T, e, bab, d, b, f, e, a, e, f,$
 $y_i, f, e, be, a, a, abe, x_i, de, ed, e,$
 $c, e, abe, f, be, a, a, a, e, ea, x_i,$
 $[24]. B, ef, da, e, a, e, e, f, a, d, fed, [16],$
 $e, bab, d, b, f, e, abe, c, f, a, Y,$
 $a, ef,$

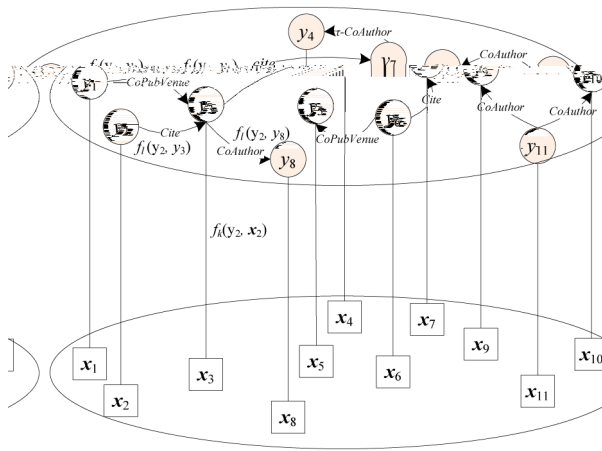


Fig. 2. Graphical representation of the HMRF model. $f_k(y_i, y_j)$ and $f_l(y_i, x_i)$ are edge feature and node feature, respectively, and will be described in the next section.

$$P(Y) = \frac{1}{Z_1} \exp\left(\sum_{(y_i, y_j) \in E, k} \lambda_k f_k(y_i, y_j)\right), \quad (2)$$

$$Z_1 = \sum_{y_i, y_j} \sum_{(y_i, y_j) \in E, k} \lambda_k f_k(y_i, y_j)$$

and $P(X|Y) = \frac{1}{Z_2} \exp\left(\sum_{x_i \in X, l} \alpha_l f_l(y_i, x_i)\right)$, $Z_2 = \sum_{y_i} \sum_{x_i \in X, l} \alpha_l f_l(y_i, x_i)$

$$P(X|Y) = \frac{1}{Z_2} \exp\left(\sum_{x_i \in X, l} \alpha_l f_l(y_i, x_i)\right), \quad (3)$$

$$Z_2 = \sum_{y_i} \sum_{x_i \in X, l} \alpha_l f_l(y_i, x_i)$$

where $f_k(y_i, y_j)$ is the edge feature between nodes y_i and y_j , and $f_l(y_i, x_i)$ is the node feature for node y_i and feature x_i . The parameters λ_k and α_l are the weights for the edge and node features, respectively. Z_1 and Z_2 are the partition functions for the author and feature layers, respectively.

The joint probability distribution is given by $P(Y, X) = P(Y)P(X|Y)$. The maximum likelihood estimation (MLE) of the parameters is performed by maximizing the log-likelihood function $L_{\max} = \log(P(Y|X)) = \log(P(Y)P(X|Y))$.

3.3 Disambiguation Objective Function

We define the disambiguation objective function as the maximum log-likelihood of the observed data X given the model parameters. The objective function is defined as $L_{\max} = \log(P(Y|X)) = \log(P(Y)P(X|Y))$.

$$L_{\max} = \log(P(Y|X)) = \log(P(Y)P(X|Y)). \quad (4)$$

By substituting (2) and (3) into (4), we have

$$L_{\max} = \log\left(\frac{1}{Z_1 Z_2} \exp\left(\sum_{(y_i, y_j) \in E, k} \lambda_k f_k(y_i, y_j) + \sum_{x_i \in X, l} \alpha_l f_l(y_i, x_i)\right)\right). \quad (5)$$

The joint probability distribution is given by $P(Y, X) = P(Y)P(X|Y)$. The maximum likelihood estimation (MLE) of the parameters is performed by maximizing the log-likelihood function $L_{\max} = \log(P(Y|X)) = \log(P(Y)P(X|Y))$.

$$f_k(y_i, y_j) = K(x_i, x_j) \sum_{r_m \in R_{ij}} [w_m r_m(x_i, x_j)]. \quad (6)$$

where $K(x_i, x_j)$ is the kernel function between features x_i and x_j , and w_m is the weight for the relationship r_m . The kernel function is defined as $K(x_i, x_j) = \exp\{-|x_i - x_j|\}$. The relationship r_m is defined as $r_m(x_i, x_j) = \exp\{-|x_i - x_j|\}$.

The node feature $f_l(y_i, x_i)$ is defined as $f_l(y_i, x_i) = K(\mu_{(i)}, x_i)$, where $\mu_{(i)}$ is the mean feature for author y_i . The kernel function is defined as $K(x_i, \mu_{(i)}) = \exp\{-|x_i - \mu_{(i)}|\}$.

$$f_l(y_i, x_i) = K(y_i, x_i) = K(\mu_{(i)}, x_i), \quad (7)$$

where $\mu_{(i)}$ is the mean feature for author y_i . The kernel function is defined as $K(x_i, \mu_{(i)}) = \exp\{-|x_i - \mu_{(i)}|\}$. The relationship r_m is defined as $r_m(x_i, x_j) = \exp\{-|x_i - x_j|\}$.

$$L_{\max} = \sum_{(x_i, x_j) \in E, k} \lambda_k K(x_i, x_j) r_k(x_i, x_j) + \sum_{x_i \in X, l} \alpha_l K(x_i, \mu_{(i)}) - \log Z, \quad (8)$$

where $Z = Z_1 Z_2$. We can maximize L_{\max} by finding the optimal parameters λ_k and α_l .

3.4 Criteria for Model Selection

We use the Bayesian Information Criterion (BIC) for model selection. The BIC is defined as $BIC = -2 \log L_{\max} + 2k$, where k is the number of parameters. We define the BIC as $BIC = -2 \log L_{\max} + 2k$.

Se fca, ef c de $K=1, \dots, e$
e e e e a e a. Te, e e a
ea e e de e e e a e c e
d be bc e. Ne, f eaç
bc e, e a a e e ea e e de e
e e. Te e a e ea e c d
a fed (e., bc e ca be). I e
ce, e ca M_h , e de c e d
e e be h . We e ef e a e a
fa fa e a e de M_h , e e h a e f 1.

n , c e.
N, a ç e e be de f M_h .
Ma ea e e ca be ed f de ec,
ç a S e e C eff ce [23], M De c
Le (MDL) [34], A a e I f a C e (AIC)
[1], a d e bab e a [22]. We ç e
BIC a e c e, beca e BIC c e f da e
a a e c e a ç a MDL a d a a
e e a a e e e c e a ç a AIC,
ç de abe be. Ba ed e e
c de a, e e a a a f e BIC ea e e
[22] a e c e.

$$BIC^v(M_h) = \log(P(M_h|P)) - \frac{|\lambda|}{2} \cdot \log(n), \quad (9)$$

e e $P(M_h|P)$ e e bab f de M_h
e e be a $P \cdot |\lambda|$ e be f a a e e
 M_h (ç ca be def ed d ffe e a, e., e
be f τ e a a e e e de M_h e
f e bab e f $P(Y)$. n e a e be. Te
ec d a a e a de c e.
I e e ce, a BIC ce a e a a -
a e e de M_h f e e da e. We e
c e f e de ec beca e ca be ea
e e ded d ffe e a. Fe a e, c e -
a c e a e K - ea [27] X -
ea [33] e a d e da a de e de a d e
e bab $P(M_h|P)$ ca be fed
 $P(P|M_h)$ acc d e Ba e a e $P(M_h|P) \propto$
 $P(P|M_h)P(M_h)$ b a e $P(M_h)$ a f.
He e, e e d a e ad a a e f de e de ce
be ee ec e e. Te, e $P(M_h)$ a
f a a e. O def (2) c de
e de e de ce a Ma f e d.

4 PARAMETER ESTIMATION

4.1 Algorithm

Te a a e e a be de e e e
a e f e a a e e $\Theta = \{\lambda_1, \lambda_2, \dots; \alpha_1, \alpha_2, \dots\}$ a d
de e e a e f a a e. Me acc a e, e
 τ e e - e d bec e f c (8)
e ec a c d a de $P(Y|X, \Theta)$.

A a e e, e ea a (cf. A 1)
f a a e e e a a c f ea e
e: Assignment f a e, a d Update f a a e e Θ .
Te ba c dea a e f a d ç e a
a a e e e Θ a d e e c a ce d f eaç c e.
Ne, e a eaç a e c e c e a d e
ca c a e, e ce d f eaç a e c e ba ed e

a e. Afe a, e da e e e f eaç
fea e f c b a τ e bec e f c.

Algorithm 1: parameter estimation

Input: $P = \{p_1, p_2, \dots, p_n\}$

Output: model parameters Θ and $Y = \{y_1, y_2, \dots, y_n\}$, where $y_i \in [1, K]$

1. Initialization

1.1 randomly initialize parameters Θ ;

1.2 for each paper x_i , choose an initial value y_i , with $y_i \in [1, K]$;

1.3 calculate each paper cluster centroid $\mu_{(i)}$;

1.4 for each paper x_i and each relationship (x_i, x_j) , calculate $f(y_i, y_j)$ and $f_k(y_i, y_j)$.

2. Assignment

2.1 assign each paper to its closest cluster centroid;

3. Update

3.1 update of each cluster centroid;

3.2 update of τ and α for each feature function.

F a a, e a d a e a e f eaç
a a e e (λ a d α). F a a f e c e
ce d, e f e a a c e e d
de f e c e a. Ba ca, a e
e a a e d be a ed d c e
a. We eed a a e de c bed fa
b a a ç e a e a a e e
a e c e ce d u. I a, e e
 γ c e a. If γ e a, e be f e e K ,
e e e γ a e ed a a a e. If
 $\gamma < K$, e a d ç e a e $(K - \gamma)$ a e a e
c e ce d. If $\gamma > K$, e e e e a e c e
a e e a e K ef. We
d ce de a e e a a e e
e a a.

Assignments. I Assignments, eaç a e x_i a ed
 $\mu_{(h)}$ a τ e $\log P(y_i|x_i)$

$$\begin{aligned} \log P(y_i|x_i) &\propto L_{x_i}(\mu_{(h)}, x_i) \\ &= \sum_{(x_i, x_j) \in E_i, R_i, k} \lambda_k K(x_i, x_j) r_k(x_i, x_j) \\ &\quad + \sum_l \alpha_l K(x_i, \mu_{(h)}) - \log Z, \end{aligned} \quad (10)$$

e e Z de ade a a fac x_i a d
ca be e ed a e ca e ab e ea e c e
f

..., e., ... d e e a d a d ... e e. H e e ,
 e a e e e e f ... c e a .
 N ... e a ... c a c a e a a a e c e (10).

T e f ... e (10) a e a ... a c b a f
 e ... a f c ... $K(x_i, \mu_{(h)})$ a d ... e e a a ...
 a f c ... $K(x_i, x_j)$, ... c a b e c a c a e d. H e e ,
 ... a c a b e ... b a a e a c ... f e a ...
 f c ... e., (Z), b e c a e ... e a ... a ... d
 a e a c e ... e a ... (Z = Z₁Z₂). A f e a ...
 ... a e b e e ... e d f a ... a e f e e c e, e.,
 b e f ... a a [30] a d c ... a e d e e c e (CD)
 [19]. W e ... e a e ... a ... a e e a ...
 f c ... a c ... a e d e e c e ... d a b a
 b e c e f c .

B a e d J e e ' e a [21], e c a b a a ... e
 b d f e e a e - e ... d (L) ... a K -
 b a c -L e b e (KL) d e e c e

$$L^{KL} = KL(q||P) = \sum_{y_i} q(y_i|x_i) \log(q(y_i|x_i)) - \sum_{y_i} q(y_i|x_i) \log(P(y_i|x_i)) = -H(q) - \langle \log(P(y_i|x_i)) \rangle_{q(y_i)}, \quad (12)$$

... e e q(y_i|x_i) ... a a ... a f ... e d b
 $P(y_i|x_i)$. $\langle \cdot \rangle_q$... e e e c a ... d e d b ... q.
 M a ... e - e ... d f ... e d a a (5) e a -
 e ... e ... e KL d e e c e (12) b e e e ... e d a a
 d b ... q^0 a d ... e e b ... d b ... e e e
 b e a a b e , q^∞ , ... e e , e f ... e c a b e c a c a e d
 b e b e a ... e c e a ... e d a b e a d
 e e c d e ... e b a b ... e e e e d e
 d b ... a ... b e a b e . A a ... e b
 d f f c ... e a b e e a ... d e e e e c d
 e . A M a ... c a M e C a (MCMC) e , d c a b e
 e d e e a e e a ... a d b ... $q^\infty(y_i|x_i)$
 ... e a ... f MCMC b e ... e e d a $q^0(y_i|x_i)$. T
 a e e e c e ... e f f c e , e c a ... e e c a e
 d e e c e a ... [19], ... c a ... a e ... e d -
 b ... b e e a G b b a ... e (... e e).
 T ... e b e c e f c ... b e c e

$$L^{KL} = KL(q^0||P) \approx KL(q^0||P) - KL(q^1||P) = \langle \log(P(y_i|x_i)) \rangle_{q^0(y_i)} - \langle \log(q^1(y_i|x_i)) \rangle_{q^1(y_i)}. \quad (13)$$

I c ... a e d e e c e e a ... e a d f ...
 $KL(q^0||q^\infty)$, e ... e , e d f f e e c e b e e e $KL(q^0||q^1)$
 a d $KL(q^1||q^\infty)$, ... e e q^1 ... e d b ... e e l - e
 e c ... c ... f ... e d a a e c (e., b e a ...) a
 a e e e a e d a f e l - e G b b a A d c a e d
 [19], ... e e l c a b e ... e a 1 ... c a e . (T a ... ,
 e c a ... c ... d e e G b b a ... e a
 ... e , e $KL(q^0||q^1)$). T e ... c e d e f e c ... c
 e d a a e c (e., q^1) f ... e d b ... q^0 d e c b e d
 A ... 2.

Algorithm 2: One-step sampling
 Input: current observation x^0 and labels y^0
 Output: sampling results of y^1 and x^1

- 1: Draw an observation x_i from the distribution of $q^0(x_i)$ ($q(x)$ can be obtained by summing over all possible labels);
- 2: Compute $P(y|x)$, the posterior probability distribution over the hidden variable given the observation x_i ;
- 3: Compute $P(y_i|\lambda)$, the probability distribution over the label y_i given labels of the other observations;
- 4: Draw a new label y_i^1 for each observation from the probability distribution $P(y_i|x)P(y_i|y_{-i})$;
- 5: Given the chosen label, compute the conditional distribution of $P(x_i|y_i)$;
- 6: Draw a new observation x_i^1 from the conditional distribution $P(x_i|y_i)$.

F a , b a e d ... e e c ... c e d d a a e c , e c a
 c a c a e (13). T e ... c a c a ... e e ... e
 d e a d . T a e ... e e f f c e , e c a ... e e
 d e e ... c e a f e d a ... [44] e a c e ... e
 a ... c e d e .

A f e ... e , d e (10), e c a c ... e e
 f e e e b e c e f c . F a , a e e d
 a ... e d e e e a d a e e a e f
 e a c a e . A a ... e f a a e e f e d e
 e e ... e e a e f e d . T e c e ... e e a e d
 a e c a e ... a ... e b e e e
 c c e e e a .

Update. I U d a e , e a c c e c e d f ... d a e d
 b ... e a ... e c e a f ... e a e c a e d

$$\mu(h) = \frac{\sum_{i:y_i=h} x_i}{\|\sum_{i:y_i=h} x_i\|_A}. \quad (14)$$

T e , b d f f e e a ... e b e c e f c ...
 e e c e a c a a e e λ_k , e , a e

$$\frac{\partial L}{\partial \lambda_k} = - \sum_{(x_i, x_j) \in E} K(x_i, x_j) r(x_i, x_j) - \frac{\partial \log Z}{\partial \lambda_k}. \quad (15)$$

W e e e ... a ... e e c d e ... a c a b e , b e c a e
 c a c a ... f Z e e d ... a ... b ... e f
 a ... e f e a c a e . A a , e a f ... e KL
 d e e c e b e c e f c (13) a d ... e e CD a
 c a c a e , e d e a e f L^{KL} ... e e c ... λ_k

$$\frac{\partial L^{KL}}{\partial \lambda_k} = \left\langle \frac{\partial \log(P(y_i|x_i))}{\partial \lambda_k} \right\rangle_{q^0(y_i)} - \left\langle \frac{\partial \log(q(y_i|x_i))}{\partial \lambda_k} \right\rangle_{q^1(y_i)} = - \sum_{(x_i, x_j) \in E} K(x_i, x_j) r(x_i, x_j) - \left\langle \frac{\partial \log(q(y_i|x_i))}{\partial \lambda_k} \right\rangle_{q^1(y_i)}. \quad (16)$$

T e f ... e ... e ... a c b a f ... e
 a ... a f c ... a d ... e e c d e ... c a b e c a c a e d
 a f e ... e l - e a ... (A ... 2).
 F a , e a c ... a a e e ... d a e d b

$$\lambda_k^{new} = \lambda_k^{old} + \Delta \frac{\partial L}{\partial \lambda_k}, \quad (17)$$

e e Δ ... e e a ... a e . W e d ... e a e f ... α .

4.2 Estimation of K

Our algorithm estimates K (see Algorithm 2) by iteratively refining a set of clusters $C = \{C_1, \dots, C_K\}$ until convergence. We start with $K=1$ and iteratively increase K until $BIC(M_2) > BIC(M_1)$. We calculate $BIC(M_2)$ for each K and choose the model with the highest BIC score.

Algorithm 3. Estimation of K

```

Input:  $P = \{p_1, p_2, \dots, p_n\}$ 
Output:  $K, Y = \{y_1, y_2, \dots, y_n\}$ , where  $y_i \in [1, K]$ 
1:  $i=0, K=1$ , that is to view  $P$  as one cluster:  $C^{(0)} = \{C_1\}$ ;
2: do {
3:   foreach cluster  $C$  in  $C^{(i)}$  {
4:     find a best two sub-clusters model  $M_2$  for  $C$ ;
5:     if  $(BIC(M_2) > BIC(M_1))$ 
6:       split cluster  $C$  into two sub clusters  $C^{(i+1)} = \{C_1, C_2\}$ ;
7:       calculate BIC score for the obtained new model;
8:   } while (existing split);
9: } while  $(BIC(M_2) > BIC(M_1))$ ;
10: choose the model as output with the highest BIC score.

```

Our algorithm estimates K by iteratively refining a set of clusters $C = \{C_1, \dots, C_K\}$ until convergence. We start with $K=1$ and iteratively increase K until $BIC(M_2) > BIC(M_1)$. We calculate $BIC(M_2)$ for each K and choose the model with the highest BIC score.

$$\sum_{i=1}^K (P(y_i) + \mu_{(i)}) + \sum_{\lambda \in \Theta} \lambda. \tag{18}$$

5 EXPERIMENTAL RESULTS

5.1 Experimental Setting

Data Sets. We use a set of 32 authors and 2,074 articles. We use the following authors: Wen Gao, Yi Li, Jie Tang, Rakesh Kumar, Bing Liu, Ajay Gupta, Dimitry Pavlov, Charles Smith, David C. Wilson, James H. Anderson, John Miller, Paul Jones, Robert Fisher, Robert Williams, Jing Zhang, Kuo Zhang, Hui Fang, Michael Wagner, Jim Smith, Wei Wang, David Jensen, David Brown, George Miller, James Johnson, Joseph Miller, Richard Taylor, Robert Moore, and William Cohen.

TABLE 3 Data Sets

Abbr. Name	#Publications	#Actual Person	Abbr. Name	#Publications	#Actual Person
Wen Gao	286	4	Jing Zhang	54	25
Yi Li	42	21	Kuo Zhang	6	2
Jie Tang	21	2	Hui Fang	15	3
Rakesh Kumar	61	5	Michael Wagner	44	1
Bing Liu	130	11	Jim Smith	33	5
Ajay Gupta	27	4	Wei Wang	306	9
Dimitry Pavlov	16	2	David Jensen	43	3
Charles Smith	7	4	David Brown	53	7
David C. Wilson	52	5	George Miller	17	2
James H. Anderson	112	2	James Johnson	17	3
John Miller	74	2	Joseph Miller	10	2
Paul Jones	13	3	Richard Taylor	93	11
Robert Fisher	105	4	Robert Moore	92	3
Robert Williams	8	2	William Cohen	110	2

Our algorithm estimates K by iteratively refining a set of clusters $C = \{C_1, \dots, C_K\}$ until convergence. We start with $K=1$ and iteratively increase K until $BIC(M_2) > BIC(M_1)$. We calculate $BIC(M_2)$ for each K and choose the model with the highest BIC score.

We use a set of 32 authors and 2,074 articles. We use the following authors: Wen Gao, Yi Li, Jie Tang, Rakesh Kumar, Bing Liu, Ajay Gupta, Dimitry Pavlov, Charles Smith, David C. Wilson, James H. Anderson, John Miller, Paul Jones, Robert Fisher, Robert Williams, Jing Zhang, Kuo Zhang, Hui Fang, Michael Wagner, Jim Smith, Wei Wang, David Jensen, David Brown, George Miller, James Johnson, Joseph Miller, Richard Taylor, Robert Moore, and William Cohen.

We use a set of 32 authors and 2,074 articles. We use the following authors: Wen Gao, Yi Li, Jie Tang, Rakesh Kumar, Bing Liu, Ajay Gupta, Dimitry Pavlov, Charles Smith, David C. Wilson, James H. Anderson, John Miller, Paul Jones, Robert Fisher, Robert Williams, Jing Zhang, Kuo Zhang, Hui Fang, Michael Wagner, Jim Smith, Wei Wang, David Jensen, David Brown, George Miller, James Johnson, Joseph Miller, Richard Taylor, Robert Moore, and William Cohen.

Experimental Design. We use a set of 32 authors and 2,074 articles. We use the following authors: Wen Gao, Yi Li, Jie Tang, Rakesh Kumar, Bing Liu, Ajay Gupta, Dimitry Pavlov, Charles Smith, David C. Wilson, James H. Anderson, John Miller, Paul Jones, Robert Fisher, Robert Williams, Jing Zhang, Kuo Zhang, Hui Fang, Michael Wagner, Jim Smith, Wei Wang, David Jensen, David Brown, George Miller, James Johnson, Joseph Miller, Richard Taylor, Robert Moore, and William Cohen.

Pairwise Precision

$$= \frac{\#PairsCorrectlyPredictedToSameAuthor}{\#TotalPairsPredictedToSameAuthor}$$

Pairwise Recall

$$= \frac{\#PairsCorrectlyPredictedToSameAuthor}{\#TotalPairsToSameAuthor}$$

$$Pairwise F_1 = \frac{2 \times Pairwise Precision \times Pairwise Recall}{Pairwise Precision + Pairwise Recall}$$

We use a set of 32 authors and 2,074 articles. We use the following authors: Wen Gao, Yi Li, Jie Tang, Rakesh Kumar, Bing Liu, Ajay Gupta, Dimitry Pavlov, Charles Smith, David C. Wilson, James H. Anderson, John Miller, Paul Jones, Robert Fisher, Robert Williams, Jing Zhang, Kuo Zhang, Hui Fang, Michael Wagner, Jim Smith, Wei Wang, David Jensen, David Brown, George Miller, James Johnson, Joseph Miller, Richard Taylor, Robert Moore, and William Cohen.

TABLE 4
Results of Name Disambiguation (Percent)

Table with 19 columns: Person Name, K-means (Prec, Rec, F1), HAC (Prec, Rec, F1), SOM (Prec, Rec, F1), SACluster (Prec, Rec, F1), CONSTRAINT (Prec, Rec, F1), and Our Approach (Fixed K) (Prec, Rec, F1). Rows include names like Cheng Chang, Jie Tang, Jing Zhang, Hui Fang, Lei Wang, and Pakchee Kummar.



fea e f eaç d; f c fe e ce, e def e a e
fea e a d e a e ec fe e ce a e; f a
e ea e e a a a e e, a, a
ea a ad def e a fea e f eaç
a ad e a e b a (dca e e ce);
ef ca, ea def e e fea e a d e
a e e a e de f, e c ed a e. I add, e
c de ed e ba e e e, d. T ef e ba ed
eaç ca a ea ec e (HAC) a f
ca ad e a eaç e e e ed a b a
a [39], e a e fea e def a def ed
ab e. T e e ba ed SAC e [52], ç e
a e de a a Kc e b b
c a a da b e f a a ca ed eaç de.
F fa c a, SAC e, e ed e a e
a b e fea e def ed a aç a d e a e
ea f a. T e d ffe e ce a SAC -
e d e d ffe e a e, e e f d ffe e ea
e c de a ea a e a e
SAC e [52].
We f e c a ed e, d e
e, d f a e d a b a : DISTINCT [49], a
c b a e, d ba ed a ea e: e

e e b e f e b e a d a d a bab ;
CONSTRAINT [51], a c a -ba ed c e a
f a e d a b a . F fa c a , 1) a
ba e e e, d a d, ec a ed e, d, e be K
f eaç a a e ea, eac a e be ;
e ef a ce e e b d f e e, d; a d
2) ed e e feedback (ea r4)
e e e (a eba e e ca e, e e feedback).

5.2 Experimental Results

5.2.1 Results

We c d c ed d a b a e e e f a e
ea ed eaç f ea a e e da a e. Tab e 4
e e . I ca be ee a e, d cea
ef eba e e e, d f a e d a b a
(+32.77% e K-Mea , +13.28% e HAC, +33.21% e
SOM, +17.57 e SAC e, a d +10.18% e CON-
STRAINT b a e a e F1 c e).
T e ba e e e, d ffe f d ad a a e :
1) e ca a e ad a a e f ea be ee
a e a d 2) e e a f ed d a ce ea e.
A, SAC e c de e ea be ee
de, c a e e ea f a a

TABLE 5
Results of Our Approach with Different Settings

Method	Precision	Recall	F1-Measure
Our Approach (Auto K)	83.01	79.54	80.05
Our Approach (w/o auto K)	90.13	88.26	88.80
Our Approach (w/o relation)	67.05	50.59	55.95

TABLE 7
Comparison with DISTINCT

Person Name	DISTINCT			Our Approach		
	Prec.	Rec.	F1	Prec.	Rec.	F1
Cheng Chang	55.07	44.19	49.03	100.00	100.00	100.00
Wen Gao	92.07	98.58	95.26	99.29	98.59	98.94
Jie Tang	79.36	93.37	85.80	100.00	100.00	100.00
Jing Zhang	100.00	75.56	86.08	83.91	100.00	91.25
Kuo Zhang	78.57	84.78	81.56	100.00	100.00	100.00
David Jensen	85.69	100.00	92.29	83.83	68.46	75.57
David Brown	69.77	74.99	72.29	89.32	91.45	90.37
David C. Wilson	87.10	90.00	88.53	94.33	67.30	78.80
Richard Taylor	68.35	63.11	65.63	94.33	79.72	86.99
Charles Smith	78.42	76.67	77.54	100.00	100.00	100.00
Hui Fang	88.60	95.00	91.69	100.00	100.00	100.00
Rakesh Kumar	92.90	96.80	94.81	99.14	96.91	98.52
Michael Wagner	72.30	75.40	73.82	85.69	82.31	83.99
Bing Liu	78.30	95.70	86.13	88.25	86.49	87.37
Jim Smith	86.30	90.40	88.30	96.37	93.80	95.08
Lei Wang	80.80	89.60	84.97	89.17	88.94	89.55
Bin Yu	68.90	77.80	73.08	95.27	72.63	82.92
Wei Wang	78.60	78.30	78.45	85.19	83.12	84.15
Ajay Gupta	98.70	92.30	95.39	97.67	96.55	97.11
Avg.	81.04	83.82	82.14	93.78	89.80	91.25

f ed d a c e f c , , , ca e c d e c b e e c e a b e e e e a e a e . O f a e d e c d e e c e a a e d e d e c e b e e e a e e e , a d e a e e e d a e a e e a f c b e e e a e . We c d c e d e e e . T e p a e a e c a e a 0.01, d c a a e e e b a a c a e a c a f c a .

Table 6 e e f a a c e a f e b e K (e b e e d b a c e e a c a b e) . We e e a e e a e d b e b a a c a e c e e a c a b e . Table 5 f e e a e a e e f a a c d f f e e e , e e / a K e e e e e f a a c a e d e f e d c e b e K a d / e a e e e e e f a a c e a (. e . , e e a e d e f e a e f c f_k(y_i, y_j) b e e) . We e e a e e a e a a a a c . W e e e a e e a e e f a c e f a a c d a (- 23.08 e c e b F_1 c e) . T c f a a d e c c a c a e d e d e c e b e e a e d e d e f a c e .

We a e d X- e a f d e b e f e e K . We a e d e e b e a l a d a b e a n , e a e e a a . We f d a X- e a f a f d e a c a b e . I a a e c e e c e Y L 2 . T e e a b e a X- e a c a a e e f e e e a b e e e a e .

TABLE 6
Result of Automatically Discovered Person Number

Person Name	Actual Number	Auto Number	Person Name	Actual Number	Auto Number
Cheng Chang	3	3	Dimitry Pavlov	2	1
Wen Gao	4	5	David Jensen	3	6
Yi Li	21	13	David Brown	7	9
Jie Tang	2	2	David C. Wilson	5	5
Gang Wu	16	12	George Miller	2	6
Jing Zhang	25	16	James H. Anderson	2	7
Kuo Zhang	2	2	James Johnson	3	3
Hui Fang	3	3	John Miller	2	5
Bin Yu	12	10	Joseph Miller	2	
Lei Wang	40	22	Paul Jones	3	
Rakesh Kumar	5	5	Richard Taylor	10	1
Michael Wagner	10	11	Robert Fisher	4	
Bing Liu	11	12	Robert Moore	3	
Jim Smith	5	5	Robert Williams	2	
Wei Wang	90	22	William Cohen	2	
Ajay Gupta	4	6	Charles Smith	4	

We c a e d a a c DISTINCT [49]. We e d e a e a e e d b [49] a d e e e f c a . We c d c e d e e e e d a a e , c a e e e f d a e d [49]. F e a e , e a e 109 a e f L e W a a d 33 a e f J S , e [49] e b e a e 55 a d 19. I a d d e d c d e e P c e e d E d e a . Table 7 e c a e e . We e e a a e a e e d c e a e f DISTINCT (+8.34% b F_1). M e e , a a c a e a d a a e a c a a c a f d e b e K , e e a DISTINCT e b e e d b e e d b e e . T e e a e d DISTINCT a d a a c a e d f f e e . DISTINCT a c d e e a - a e a d a e c f e e c e e a , a d d e d e c c d e e C A , a d C P b V e e e a , a e e e a c a b e d e d f e e a e c f e e c e a d a - a e e a .

5.2.2 Efficiency Performance

We e a a e d e e f f e c c e f a c e f a a c f e 32 a a e a e a d e c e I e C e D c e (1.6 G H) . Table 8 e C P U e e e d f a e a e d f f e e a . We e a b e a 100 a e a d e a e a e e f 100 a d a e . F a a e a e a e a e e a 1 e c d . T e a e f a a a e a c e .

TABLE 8
Comparison of Efficiency Performance (Seconds)

Person Name	k-means	k-means	HAC	SCCluster	DISTINCT	Our Approach
Wen Gao	4.8	5.1	12.9	30.4	56.0	20.3
Lei Wang	3.7	2.4	6.8	4.1	12.1	4.6
Bing Liu	1.6	1.9	4.2	5.4	1.1	5.8
Wei Wang	28.7	5.1	73.1	46.9	83.3	100.2
Robert Fisher	2.8	1.3	5.6	0.2	0.2	0.8
William Cohen	0.8	1.2	3.0	0.06	0.6	0.9
Average over 100	0.52	0.26	1.14	0.96	0.87	1.42

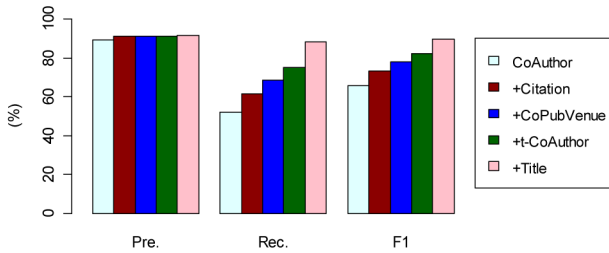


Fig. 3. Contribution of relationships.

5.2.3 Feature Contribution Analysis

We evaluate the contribution of features to the performance of the model. The features are grouped into five categories: CoAuthor, Citation, CoPubVenue, CoAuthor, and Title. The results are shown in Figure 3. The CoAuthor feature contributes the most to the performance, followed by Citation, CoPubVenue, CoAuthor, and Title.

5.2.4 Distribution Analysis

We analyze the distribution of the features. The results are shown in Figure 4. The CoAuthor feature has the highest distribution, followed by Citation, CoPubVenue, CoAuthor, and Title.

5.2.5 Application Experiments

We evaluate the performance of the model on various datasets. The results are shown in Figure 5. The CoAuthor feature has the highest performance, followed by Citation, CoPubVenue, CoAuthor, and Title.

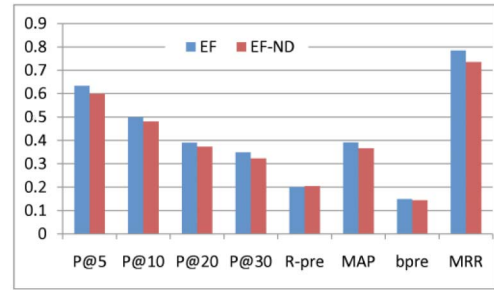


Fig. 4. Performances of expert finding.

5.3 Online System

The online system is designed to provide a user-friendly interface for expert finding. It allows users to search for experts based on various criteria and provides detailed information about each expert.

6 DISCUSSION

6.1 Connections with Previous Work

We compare our method with previous work. The results show that our method outperforms previous methods in terms of accuracy and efficiency.

$$L_{max} = \sum_{x_i \in X, l} \alpha_i K(x_i, \mu_i) - \log Z. \quad (19)$$

Our method is based on the K-means algorithm, which is a widely used clustering algorithm.

The connection with X-means is discussed in [33]. X-means is an extension of K-means that allows for variable cluster sizes.

The connection with X-means is discussed in [33]. X-means is an extension of K-means that allows for variable cluster sizes.

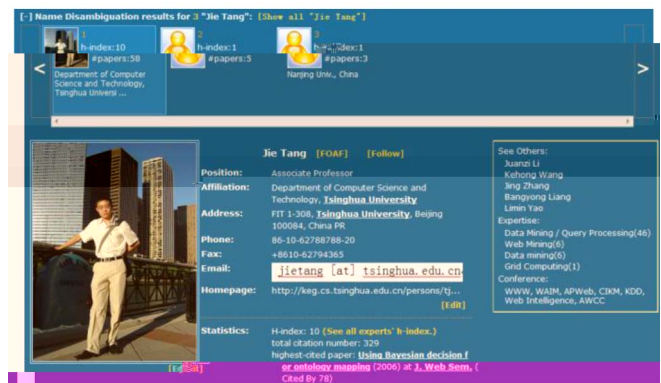


Fig. 5. Name disambiguation system (http://arnetminer.org).

e, d f a e a a X- ea f e c de e bab P(Y) f , e., de e de ce be ee daa . E ce f de eec , X- ea ef a K- ea .

Connection with the constraint-based disambiguation method:

I c a -ba ed c e , e., [2], e e ca c a de e c e ce . I e a ea ed a e d a b a ad ba ed e e [51], [41]. T e a c a c de - a dca - . M - ea a daa be ed ec e a dca - ea daa be ed d ffe e c e . We ca ada f a e a c a -ba ed c e b edef eed e e a f c .

Connection with disambiguation using spectral graph clustering:

S ec a a c e [12] a a f d b a c f ea be ee daa . K- a ec a a c e a a bee e ed f a e d a b a [18]. We ca e a e a e ed daa a f e e e a ed d ffe e c e (e., I(i ≠ j)) e bec ef c . T e f a e ca ada c e b e e ec d a f (8)

L_min = - sum_{(x_i, x_j) in E, R, k} K(x_i, x_j)r_k(x_i, x_j) + log Z. (20)

I e e ce, e bec ef c ea a e e e e e a bab e e HMRf ad f c e de e de ce be ee a e . C a e e f a e ffe e e a ad a a e :1) I ad a e , d , a - e f a e a e de e de , ca a e ad a a e f ea be ee a e . 2) T e - ed f a e ca be ea e e ded e - e - ed ea b e feedbac . 3) O f a e ca be e ed a a e e a f a e f e e a e ed e , d .

7 CONCLUSION AND FUTURE WORK

I a e, e a e e aed e be f a e d a b a . We a ef a ed e be a fed f a e ad ed a e e a ed bab - c de e be . We a edef ed a d a b a - bec ef c f e be a d a e ed a - e a a e e a a . We a e a e ed a d a ca a c f e a e be f e e K. E e e a e d ca e a e ed e, d fca ef e ba e e e, d . We a ed e e f d , cea e e (+2%) ca be ba ed.

A e e e, d be e e e a e a e e f e e f a f a e d a b a , a e a b be e e e e . M e e, a ee d , c de e LDA ca e a ed a b a .

ACKNOWLEDGMENTS

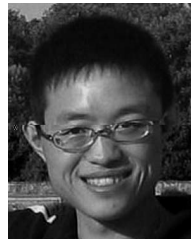
T e a d e a H C e f d e cec de f SAC e a d X a Y f d e cec de f DISTINCT f ec a e e - e . T e a a P f P Y f a abe

. Je Ta ed b e Na a Sc e ce F da f C a (N . 61073073), e C e e Na a Ke F da Re ea c (N . 60933013, N .61035004), a d a S ec a F d f FSSP.

REFERENCES

[1] H. A a e, A Ne L a e Sa ca M de Ide fca- IEEE Trans. Automatic Control, AC-19, 6, 716-723, Dec. 1974. [2] S. Ba , M. B e , a d R.J. M e , A P bab c Fa e f Se -S e ed C e , Proc. ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '04), 59-68, 2004. [3] R. Be e a d A. McCa , D a b a Web A ea - a ce f Pe e a S ca Ne , Proc. Int'l Conf. World Wide Web (WWW '05), 463-470, 2005. [4] O. Be e , H. Ga ca- M a, D. Me e a, Q. S , S.E. W a , a d J. W d , S : A Ge e c A ac E Re e , The VLDB J., 18, 255-276, 2008. [5] I. B a a c a a a d L. Ge , C ec e E Re Re a a Da a, ACM Trans. Knowledge Discovery from Data, 1, a ce 5, 2007. [6] C. B ce e a d E.M. V ee , Re e a E a a I c ee I f a , Proc. Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '04), 25-32, 2004. [7] Z. C e , D.V. Ka a , a d S. Me a, Ada e Ga ca A ac E Re , Proc. Seventh ACM/IEEE-CS Joint Conf. Digital Libraries (JCDL '07), 204-213, 2007. [8] Z. C e , D.V. Ka a , a d S. Me a, E C e A a f C b M e E Re S e , Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '09), 207-218, 2009. [9] D. C , R. Ca a a, a d A. McCa , Se - e ed C e U e Feedbac , Tec ca Re TR2003-1892, C e U , 2003. [10] D. Ca , X. He, a d J. Ha , S ec a Re e f D e a Red c , ec ca e , 2856, UIUC 2004. [11] P.T. Da , D.K. E , a d J.L. Ka a , Me d f Pec e Na ed E Mac D a C ec , Proc. ACM/IEEE-CS Joint Conf. Digital Libraries (JCDL '03), 125, 2003. [12] C. D , A T a S ec a C e , Proc. Int'l Conf. Machine Learning (ICML '04), 2004. [13] M. E e , R. Ge , B.J. Ga , Z. H , a d B. Be -M e, J C e A a f A b e Da a a d Re a Da a: T e C ec ed K-Ce e P be , Proc. SIAM Conf. Data Mining (SDM '06), 2006. [14] S. Ge a a d D. Ge a , S c a c Re a a , Gbb D b - a d e Ba e a Re a f I a e , IEEE Trans. Pattern Analysis and Machine Intelligence, PAMI-6, 6, 721-742, N . 1984. [15] Z. Ga a a a d M.I. J da , Fac a Hdde Ma M de , Machine Learning, 29, 245-273, 1997. [16] J. Ha e e a d P. C ff d, Ma Fed F e Ga , a d La ce, U b ed a c , 1971. [17] H. Ha , L. Ge , H. Z a, C. L , a d K. T , T S e ed Lea A ac ef Na e D a b a A C a , Proc. ACM/IEEE Joint Conf. Digital Libraries (JCDL '04), 296-305, 2004. [18] H. Ha , H. Z a, a d C.L. Ge e , Na e D a b a A C a U a K-Wa S ec a C e Me d, Proc. ACM/ IEEE Joint Conf. Digital Libraries (JCDL '05), 334-343, 2005. [19] G.E. H , Ta P dc f E e b M , C a e D e e ce, J. Neural Computation, 14, 1771-1800, 2002. [20] L. Ja , J. Wa , N.A , S. Wa , J. Z a , a d L. L ., GRAPE: A Ga -Ba ed Fa e f D a b a Pe e A ea a ce Web Sea c , Proc. Int'l Conf. Data Mining (ICDM '09), 199-208, 2009. [21] M.I. J da , Z. Ga a a , T. Jaa a, a d L. Sa , A I dc Va a a Me d f Ga ca M de , Learning in Graphical Models, 37, 105-161, 1999. [22] R. Ka a d L. Wa e a , A Refe ce Ba e a Te f Ne ed H ee e a d I Re a e Sc a C e , J. Am. Statistical Assoc., 90, 773-795, 1995.

- [23] L. Kaffa and P. R. ... , *Finding Groups in Data: An Introduction to Cluster Analysis*. Wiley, 1990.
- [24] R. K. Deane and J.L. S. ... , *Markov Random Fields and Their Applications*. Academic, 1980.
- [25] H. K. ... , S. Ge ... , and A. Ke ... , *Handbook of Applied Probability*, 5, 3, 577-602, 1995.
- [26] X. L. ... , P. M. ... , D. R. ... , *Defining and Testing Abnormal Names: Detecting and Evaluating Anomalies*, *Proc. 19th Nat'l Conf. Artificial Intelligence (AAAI '04)*, 419-424, 2004.
- [27] J. MacQueen, *Some Methods for Classification and Analysis of Multivariate Observations*, *Proc. Fifth Berkeley Symp. Math. Statistics and Probability*, 1967.
- [28] D.M. McRae-S. ... and N.R. S. ... , *Abnormal Names: A Case Study*, *Proc. ACM/IEEE Joint Conf. Digital Libraries (JCDL '06)*, 53-54, 2006.
- [29] E. M. ... , W.W. C. ... , and A.Y. N. ... , *Chinese Search and Named Entities in English*, *Proc. 29th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '06)*, 27-34, 2006.
- [30] K.P. M. ... , Y. We ... , and M.I. J. da ... , *Learning from Ambiguous Labels*, *Proc. Conf. Uncertainty in Artificial Intelligence (UAI '99)*, 467-475, 1999.
- [31] M.E.J. Ne ... and M. G. ... , *Finding and Evaluating Clusters in Networks*, *Physical Rev. E*, 69, 026113, 2004.
- [32] B. O. ... and D. Lee, *Scalable Named Entity Recognition*, *Proc. SIAM Int'l Conf. Data Mining (SDM '07)*, 2007.
- [33] D. Pe ... and A. M. ... , *X-Mean: Efficient K-Mean Clustering*, *Proc. Int'l Conf. Machine Learning (ICML '00)*, 2000.
- [34] J. R. ... , A. U. ... , and P. ... , *Efficient Data Mining*, *J. Annals of Statistics*, 11, 2, 416-431, 1983.
- [35] J. S. ... and J. Ma ... , *Named Entity Recognition*, *IEEE Trans. Pattern Analysis and Machine Intelligence*, 22, 8, 888-905, Aug. 2000.
- [36] L. S. ... , B. L. ... , and W. Me ... , *Large Text Mining Clustering*, *Proc. IEEE Int'l Conf. Data Eng. (ICDE '09)*, 880-891, 2009.
- [37] Y. S. ... , J. H. a. ... , I.G. C. ... , J. L. ... , and C.L. G. ... , *Efficient Text-Based Named Entity Recognition*, *Proc. ACM/IEEE Joint Conf. Digital Libraries (JCDL '07)*, 342-351, 2007.
- [38] Y. S. ... , Y. Y. ... , and J. H. a. ... , *Robust Named Entity Recognition*, *Proc. ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '09)*, 2009.
- [39] Y.F. Ta ... , M. Ka ... , and D. Lee, *Search Engine Data Analysis*, *Proc. ACM/IEEE Joint Conf. Digital Libraries (JCDL '06)*, 314-315, 2006.
- [40] J. Ta ... , J. Z. a. ... , L. Ya ... , J. L. ... , L. Z. a. ... , and Z. S. ... , *Entity and Relationship Classification*, *Proc. 14th ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '08)*, 2008.
- [41] J. Ta ... , L. Ya ... , D. Z. a. ... , and J. Z. a. ... , *Acquiring Web User Profiles*, *ACM Trans. Knowledge Discovery from Data*, 5, 4, Dec. 2010.
- [42] Y. Ta ... , R.A. Ha ... , and J.M. Pa. ... , *Efficient Graph Search*, *Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '08)*, 567-580, 2008.
- [43] J. Ve ... and E. A. ... , *Classification of Self-Organizing Maps*, *IEEE Trans. Neural Network*, 11, 3, 586-600, Mar. 2000.
- [44] M. We ... and G.E. H. ... , *Neural Networks for Feature Extraction*, *Proc. Int'l Conf. Artificial Neural Networks (ICANN '01)*, 351-357, 2001.
- [45] M. We ... and K. K. ... , *Bayesian K-Means*, *Proc. SIAM Int'l Conf. Data Mining (SDM '06)*, 472-476, 2006.
- [46] S.E. W. a. ... , D. Me ... , G. K. ... , M. T. e. b. a. d. a. d. H. G. a. c. a. -M. a. ... , *Entity Recognition*, *Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '09)*, 219-232, 2009.
- [47] S.E. W. a. ... , O. Be ... , and H. G. a. c. a. -M. a. ... , *Entity Recognition*, *The VLDB J.*, 18, 6, 1261-1277, 2009.
- [48] X. X. ... , N. Y. ... , Z. Fe ... , and T.A.J. Sc. ... , *Scalable Semantic Analysis of News*, *Proc. ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '07)*, 824-833, 2007.
- [49] X. Y. ... , J. Ha ... , and P.S. Y. ... , *Obec Detection: Detecting Obvious Information*, *Proc. Int'l Conf. Data Eng. (ICDE '07)*, 1242-1246, 2007.
- [50] H. Y. ... , W. K. ... , V. Ha ... , and J. W. b. ... , *Large Scale Classification of Academic Databases*, *ACM Trans. Information Systems*, 24, 3, 380-404, 2006.
- [51] D. Z. a. ... , J. Ta ... , J. L. ... , and K. Wa ... , *Acquiring and Evaluating Named Entities*, *Proc. ACM Conf. Information and Knowledge Management (CIKM '07)*, 1019-1022, 2007.
- [52] Y. Z. ... , H. C. e. ... , and J.X. Y. ... , *Graph-Based Semantic Analysis*, *Proc. VLDB Endowment*, 2, 1, 718-729, 2009.



Jie Tang is an associate professor at Tsinghua University. His research interests are social network analysis, data mining, and semantic web.



A.C.M. Fong is a professor in the School of Computing and Mathematical Sciences, Auckland University of Technology. He has published widely in the areas of data mining and communications.



Bo Wang is currently working toward the PhD degree from Nanjing University of Aeronautics and Astronautics. His research interests include transfer learning and information network analysis.



Jing Zhang received the MS degree from Tsinghua University in 2008. Her research interests include information retrieval and text mining.

► For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.