

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 e e a a e de ca a e e ea
d.I e a ed a e 300 c a e
a e a e ed b e a 114 e e (a
ab 78.74 e ce) e U ed S a e (:// a e .
aba .c / a e_ a e .). I a a ca
c a ce fc ea e a a e e ad f a
e a , e e e a e a e ed a e de fe
e e e f a . Na e a b ea
e a f e e e ed f a .
T de e e e e f e be , e a e
e a ed 100 e a e e b ca da a ad
f d, f e a e, e e a e 54 a e a ed b
25 d ffe e J Z a e DBLP da aba e. A , ee
de a ed Y L a e ad a ed f e f
a ' ab.

1.1 Motivation

We be b a e be a e a e
da f a ea- d e (://a e e .)
[40]. I e , e e ac e ea ce f e
f e e ba d e a e e b ca da af e
da aba e c a DBLP, ACM D a L b a , C e See ,
ad SCI. I e e a , e e ab a e e a e
a b be . F . 1 a f ed e a e. I
F . 1, eac de de e a a e (e ed). Eac
d ec ed ed e de e a e a be ee a e

- 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 ff c ac 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 , be ef
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 ee 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 ea a ce d a b a [3], [20], a e
de f ca [26], a d Ob ec d c [49]. De e a
a ac 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, unsuper-
vised based, a d constraint based. T e e ed-ba ed
a ac (e ., [17]) e ea a ec f c c a f ca
de f eac a a e f e a -abe ed
a da a. T e , e ea ed de ed ed c
e a a e f eac a e. I e e ed-
ba ed a ac (e ., [18], [36], [37], [49]), c e
a c de a e e ed f d a e

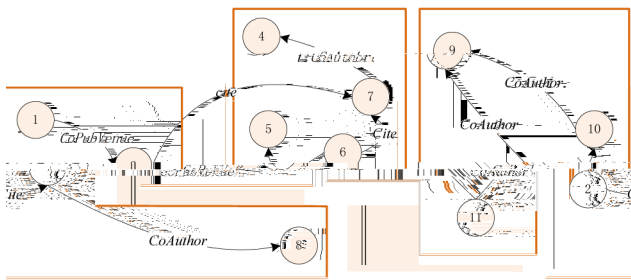


Fig. 1. An example of name disambiguation.

a , a d a e d ffe e a a e a ed
d ffe e a . T e c a -ba ed a ac a
e e e c e a . T e d ffe e ce a
e - ded c a a e ed de e c e
a a d be e da a a (e. , [2], [51]).
F e e, e e a e a ac e ba ed e,
c a /a a , a d c b a f e d ffe e
a ac e a e bee d ed. F e a e, W a e a .
[47] d ce a e a e e -ba ed a ac e e
e c e e ce e da aba e a d de e
a de f a e e e ce a e e
e . Da e a . [11] a e de e da e ac e e
c e a e ca e e cc e ce f a ed
e e a e e . T e e de f
efe e ce a e a bec (e. , a a c a b d)
e e a ed a e f a bec ad ac ec .
McRae-S e ce a d S adb [28] e e a a -ba ed
a ac a d a b a a e -ca e c a
e b e f -c a , c a e a . T e
a ac ca ac e e a ec b a e a e
eca . Y e a . [50] a e de e ed e ed a ac e
de f e f f f a b abbe a
ec e e a e a . M e e ce , C e e a . [8] d
c b e e d ffe e d a b a a ac e
a d e a e e e e e b e f a e ,
c c b e e e f e ba e - e e e
e e a e e e e e e e
acc ac fe e . W a e a . [46] e a
e a eb c f a e e e e e e f
b c a e ef e ced b e e ce ed b c . O
a d Lee [32] d e ca ab e f e a e
d a b a b e . A c e a
bee ade, e e d d ac e e a fac
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 da a a ba ed e ca c e; e e e d (e. , [18], [42]) a c e e da a a acc d de a . A fe e ea ce (e. , [38], [52]) c b e e e ce f f a . F e a e, Z e a . a e c b e f a ba ed b e e a b e (.e., de a) a d a ca c e b f c c a a b e a e ed a e c a e f (a b e, a e) a e ce, a d b e e e a e a e e e ce' c e e a e b d

- a d a de . T e a e c a
ea a e b e e d ca d ca
f a . A e e a e e abe
de a e a a b e a ce a e e
c e e f a e e ce e d a ce
ea e, a ba a ce ec b
f e d ffe e f a a e be
T e a e abe c c de a add a b e
a f a ec e bec e
de ade e ac e c e e . F e,
[52], e e e e a da a e c a e fe
a b e . T e f da a e (ca b) a
e (b a) a b e a d e ec d da a e f
DBLP b b a ca da a a a b e .
W e a e a c ce de a b e f a
e ed f ac e a e d a b a
be effec e .
2. T e e f a ce fa eaf e e ed e d
de e d acc a e e a K. A
e e a c e a c a X- ea [33]
ca a a ca f d e be K ba ed e
c e , cea e e c a
e d ca be d ec a ed e a e
d a b a b e .
 3. I e e d , e da a a c a
e e de a d e a ; e
be e , e e a be e d ffe e
ea (e. , C A a d C a) be ee
de . T e e f d ffe e ea a a e
d ffe e a ce f e a e d a b a
be . H a a ca de e de ee f
c b f d ffe e ea a
c a e b e .

1.3 Our Solution

Ha c d c ed a e a , e e a
fed bab c f a e adde e ab e
c a e e . S ec f ca , ef a e e d a b a
be a Ma Ra d Fed (MRF) [16], [24],
c e da a a e c e e b ca a b e a d
ea . W e e ad a ca ac f e a -
e be f e e K a da - e a f
a a e e e a . T e ed a ac ca ac e e
be e ef a ce a e d a b a a e
e d beca e ea ac a e ad a a e f e -
de e de ce be ee a e a e . T e be f
ed e, ef f a e a e
be f a e d a b a a fed f a e
a d ac e e be e e .
T e ed f a e e e e a . O e ca
c a e a e a a fea e ca fea e e
f a e , e. , a fea e ba ed e eb ea c e e
ed. T e f a e ca be a e e ded dea
a e be c a e e a
ea a da aba e [4].
O c b a e c de: 1) f a a
f e a e d a b a b e a fed bab -
c f a e ; 2) a f a a e e
a a e e e a e f a e ; a d 3) a e ca
e f ca f e effec e e f e ed 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 $p_i.cites p_{j-1}$
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. 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.

- **CoPubVenue (r_1)** r_1 is a relationship between two papers p_1 and p_2 . $r_1(p_1, p_2)$ holds if and only if $p_1.pubvenue = p_2.pubvenue$. We denote this relationship as r_1 .
- **CoAuthor (r_2)** r_2 is a relationship between two papers p_1 and p_2 . $r_2(p_1, p_2)$ holds if and only if $\exists r, s > 0, a_i^{(r)} = a_j^{(s)}$. We denote this relationship as r_2 .
- **Citation (r_3)** r_3 is a relationship between two papers p_1 and p_2 . $r_3(p_1, p_2)$ holds if and only if $p_1.cites p_2$ or $p_1.cites p_{2-1}$. We denote this relationship as r_3 .
- **Constraint (r_4)** r_4 is a relationship between a paper p_i and a user u . $r_4(p_i, u)$ holds if and only if u has provided feedback for p_i . We denote this relationship as r_4 .
- **τ -CoAuthor (r_5)** r_5 is a relationship between two papers p_1 and p_2 . $r_5(p_1, p_2)$ holds if and only if τ -extension co-authorship ($\tau > 1$) exists between p_1 and p_2 . We denote this relationship as r_5 .

The following definitions describe the relationships between papers. Table 2 shows the relationships between papers. We define the following relationships between papers (Table 2). Section 2.1.1.

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.

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

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 (Table 2). Section 2.1.1.

$f \in \{-ca, ed, f, a, e, a\}$ [13] $e \in \{e, e, e, e, f, a, a, b, c, d, e, f, g, h, i, j, k, l, m, n, o, p, q, r, s, t, u, v, w, x, y, z\}$
 $b \in \{ca, da, a, P, b, ca, ad, ea, ae, a, -\}$
 $f \in \{ed, a, d, ec, ed, a, c, eac, de\}$
 $e \in \{e, e, a, a, e, a, deac, ed, ea, ea, ., A, b, e\}$
 $f \in \{a, a, e, a, e, a, ac, ed, ec, e, d, de, a, a\}$
 $fea \in \{e, ec, ., F, e, ec, , e, e, d, (afe, e, f, ea, a, de, ed, a, e, f, c, e, d, fea, ef, c, ., S, e, HMRF, de, c, de, b, e, a, e, e, f, fea, ef, c, a, da, a, e, d, ffe, e, e, ., S, c, a, f, a, e, a, ffe, add, a, ad, a, a, e, :, f, , e, ed, ea, , e, ed, ea, , a, d, e, -\}$
 $d \in \{f, e, a, d, e,)\}$ $ea \in \{b, e, fa, a, e, a\}$
 $fea \in \{e, a, d, e, e, be, f, e, cc, e, ce, a, e, a, e, ., F, a, , e, ca, def, e, e, b, ca, f, a, e, a, a, f, :\}$

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, $r_k(p_i, p_j) = 0$.

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

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

3 OUR FRAMEWORK

3.1 Basic Idea

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

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

3.2 Hidden Markov Random Fields

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

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

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

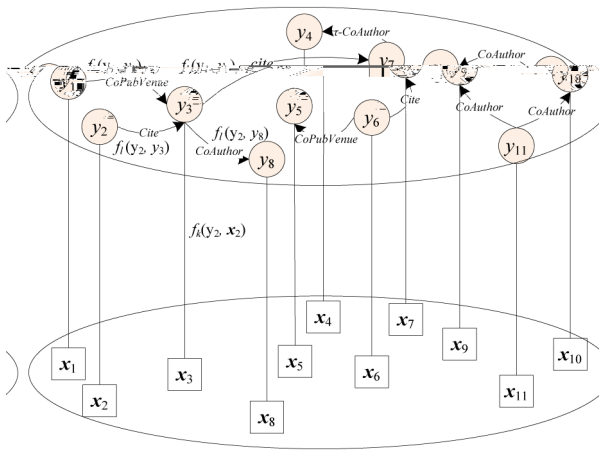


Fig. 2. Graphical representation of the HMRF model. $f(y_i, y_j)$ and $f(y_i, x_i)$ are edge feature and node feature, respectively, and will be described in the next section.

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

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

and

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

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

where $f_k(y_i, y_j)$ is the edge feature function, and $f_l(y_i, x_i)$ is the node feature function. λ_k and α_l are the weights of the edge and node features, respectively. Z_1 and Z_2 are the partition functions.

The joint probability distribution is defined as $P(Y, X) = P(Y)P(X|Y)$.

3.3 Disambiguation Objective Function

We define the maximum likelihood estimation (MLE) objective function as follows:

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

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

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

where $f_k(y_i, y_j)$ and $f_l(y_i, x_i)$ are the edge and node features, respectively. The edge features are defined as follows:

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

where $K(x_i, x_j)$ is the kernel function, and R_{ij} is the set of relationships between x_i and x_j .

The kernel function is defined as follows:

The kernel function is defined as follows:

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

where $\mu_{(i)}$ is the mean vector of the node x_i .

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

where $Z = Z_1 Z_2$.

3.4 Criteria for Model Selection

We use the Bayesian Information Criterion (BIC) for model selection.

Section, effective of $K=1$, a, e e
 e e e e a e a. Te, e e a
 ea e e de e e e e a e c e
 d be bc e. Ne, f eac
 bc e, e a a e e ea e e de e e
 e e. Te ea 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 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 c e e be de f M_h .
 Ma ea e e ca be ed f de e ec,
 c 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 c e
 BIC a e c e, beca e BIC c e f da e -
 a a e c e a c a MDL a d a a
 e e a a e e c e a c a AIC,
 c 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 (c ca be def ed d ffe e a, e., e
 be f γ 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 c e a a e -
 a e e de M_h f e e da e e. We e
 c e f e de e ec beca e ca be ea
 e e ded d ffe e a. F e a e, c e -
 a c e a e K - ea [27] X -
 ea [33] e a d e da a a de e de a d e
 e bab $P(M_h|P)$ ca be f ed
 $P(P|M_h)$ 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.
 H e e, e e d a e ad a a e f de e de ce
 be ee ec e e. T, 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 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 ea e fa a e. M e acc a e, e
 γ e e - e d bec e f c (8)
 e ec ac 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 e a e
 e : *Assignment* f a e, a d *Update* f a a e e Θ .
 Te ba c dea a e f a d c e a
 a a e e e Θ a d e ec a ce d f eac c e.
 Ne, ea eac a e c e c e a d e
 ca c a e e ce d f eac a e - c e ba ed e

a e. Af e a, e da e e e f eac
 fea e f c b a γ 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 $f_i(y_i, y_j)$ 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 the weight for each feature function.

F a γ a, e a d a e a e f eac
 a a e e (λ a d α). F a γ 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 a
 e a a e d be a ed d c e
 a. We eed a a e e de c bed fa
 b a a c 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 c e a e $(K - \gamma)$ a e a e
 c e ce d. If $\gamma > K$, e e ea 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*, eac a e x_i a ed
 $\mu_{(h)}$ a γ e $\log P(y_i|x_i)$

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

e e Z de ade a a γ a fac x_i a d
 ca be e ed a e ca e ab e ea e c e
 f ac 33 21 6251 0 0 a 21 3-491.[(f 6-629.5(d ffe ())TJ]/F3 1

$f_c, e., d e e a d a d e e. H e e,$
 $e a e e e e f c e a .$
 $N, e a c a c a e a a a e c e (10).$
 $T e f e (10) a e a a c b a f$
 $e a f c K(x_i, \mu_{(i)}) a d e e a a -$
 $a f c K(x_i, x_j), c c a b e c a c a e d. H e e,$
 $a c a b e b a a e a c f e a$
 $f c, e., (Z), b e c e e a a \gamma a d$
 $a e a c e e a (Z = Z_1 Z_2). A f e a -$
 $a e b e e e d f a a e f e e c e, e.,$
 $b e f a a [30] a d c a e d e e c e (CD)$
 $[19]. W e e a e a a e e a$
 $f c a c a e d e e c e d a b a$
 $b e c e f c .$
 $B a e d J e e' e a [21], e c a b a a e$
 $b d f e e a e - e d (L) a K -$
 $b a c - L e b e (KL) d e e c e$

$$\begin{aligned}
 L^{KL} &= KL(q \| P) \\
 &= \sum_{y_i} q(y_i | x_i) \log(q(y_i | x_i)) - \sum_{y_i} q(y_i | x_i) \log(P(y_i | x_i)) \\
 &= -H(q) - \langle \log(P(y_i | x_i)) \rangle_{q(y_i)},
 \end{aligned}
 \tag{12}$$

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

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

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

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 $P(y_i | x_i)$, the posterior probability distribution over the label variable given the observation x_i ;
- 3: Compute $q^1(y_i | x_i)$, the probability distribution over the label variable given labels of its neighboring observations;
- 4: Draw a new label y_i^1 for each observation from the probability distribution $P(y_i | x_i)P(y_i | y_{-i})$;
- 5: Given the chosen label, compute the conditional distribution of $P(x_i | y_i)$;
- 6: Draw each feature of the new observation x_i^1 from the conditional distribution $P(x_i | y_i)$.

$F a, b a e d e e c c e d d a a e c, e c a$
 $c a c a e (13). T e c a c a e e e e$
 $d e a d. T a e e e f f c e, e c a e e$
 $d e e c e a f e d a [44] e a c e e$
 $a c e d e.$
 $A f e e d e (10), e c a c e e e$
 $f e e b e c e f c. F a, a e e d$
 $a e d e e a d a e e a e f$
 $e a c a e. A a e f a a e e f e d e$
 $e e e e a e f e d. T e c e e e a e d$
 $a e c a e a e b e e e$
 $c c e e e a.$

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

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

$T e, b d f f e e a e b e c e f c$
 $e e c e a c a a e e \lambda_k, e a e$

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

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

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

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

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

$e e \Delta e e a a e. W e d e a e f \alpha.$

4.2 Estimation of K

Our algorithm estimates the number of clusters K by iteratively splitting clusters until the BIC score is maximized. The algorithm starts with $K=1$ and iteratively splits clusters until the BIC score is maximized. The algorithm starts with $K=1$ and iteratively splits clusters until the BIC score is maximized.

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 $y_i \in [1, K]$

- $i=0, K=1$, that is to view P as one cluster: $C^{(0)}=\{C_1\}$;
- do {
- foreach cluster C in $C^{(i)}$ {
- find a best two sub-clusters model M_2 for C ;
- if($BIC(M_2) > BIC(M_1)$)
- split cluster C into two sub clusters $C^{(i+1)}=\{C_1, C_2\}$;
- calculate BIC score for the obtained new model;
- }while(existing split);
- }while(BIC score of the model as output, until the BIC score is maximized).

Our algorithm estimates the number of clusters K by iteratively splitting clusters until the BIC score is maximized. The algorithm starts with $K=1$ and iteratively splits clusters until the BIC score is maximized.

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

5 EXPERIMENTAL RESULTS

5.1 Experimental Setting

Data Sets. We use a real-world dataset from [40]. We consider a dataset of 2,074 authors and their publications. The dataset is divided into two parts: a training set and a test set. The training set consists of 25 authors and their publications, and the test set consists of 40 authors and their publications. The authors in the training set are 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, and Robert Williams. The authors in the test set are 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	12
Bing Liu	130	11	Jim Smith	33	5
Ajay Gupta	27	4	Wei Wang	306	90
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	10
Robert Fisher	105	4	Robert Moore	92	3
Robert Williams	8	2	William Cohen	110	2

Our algorithm estimates the number of clusters K by iteratively splitting clusters until the BIC score is maximized. The algorithm starts with $K=1$ and iteratively splits clusters until the BIC score is maximized.

We use a real-world dataset from [40]. We consider a dataset of 2,074 authors and their publications. The dataset is divided into two parts: a training set and a test set. The training set consists of 25 authors and their publications, and the test set consists of 40 authors and their publications.

The authors in the training set are 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, and Robert Williams. The authors in the test set are 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 real-world dataset from [40]. We consider a dataset of 2,074 authors and their publications. The dataset is divided into two parts: a training set and a test set. The training set consists of 25 authors and their publications, and the test set consists of 40 authors and their publications.

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 real-world dataset from [40]. We consider a dataset of 2,074 authors and their publications. The dataset is divided into two parts: a training set and a test set. The training set consists of 25 authors and their publications, and the test set consists of 40 authors and their publications.

1. //a e e. /d a b a .

TABLE 4
Results of Name Disambiguation (Percent)

Person Name	K-means			HAC			SOM			SACluster			CONSTRAINT			Our Approach (Fixed <i>K</i>)		
	Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1
Cheng Chang	89.47	68.00	77.27	100.0	100.0	100.0	76.30	65.42	70.44	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
Jie Tang	95.38	72.09	82.12	100.0	100.0	100.0	84.92	70.65	77.13	90.14	82.04	85.90	100.0	100.0	100.0	100.0	100.0	100.0
Geny Wu	72.84	72.04	72.33	93.54	93.54	93.54	74.79	73.72	74.26	72.66	72.66	72.66	78.72	78.32	77.86	98.36	83.05	81.62
Jing Zhang	7.88	26.03	12.10	85.00	69.86	76.69	38.76	64.23	48.35	72.00	86.75	78.69	83.91	100.0	91.75	83.91	100.0	100.0
Kuo Zhang	60.00	60.00	60.00	100.0	100.0	100.0	82.50	70.20	75.85	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
Hui Fang	60.87	90.32	72.73	100.0	100.0	100.0	40.60	80.60	54.00	92.21	54.20	68.27	100.0	100.0	100.0	100.0	100.0	100.0
Lei Wang	11.98	21.87	15.48	68.45	41.12	51.38	21.52	57.34	31.29	44.40	75.59	55.94	91.58	92.59	92.08	88.14	96.91	98.01
Rakesh Kumar	68.82	91.28	78.47	63.36	92.41	75.18	62.83	90.17	74.06	80.98	82.43	81.70	92.37	90.18	95.65	96.14	96.91	98.01

fea e f eac 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 f a e d a b a . F fa c a , 1) a
 ea a a d def e a fea e f eac ba e e e d a d ec a ed e d, e be *K*
 a a d e a e b a (dca e e ce); f eac a a e ea eac a e be ; ,
 ef ca , ea def e e fea e a d e e ef ace e e b df e e d; a d
 a e ea e de f e c ed a e. I add , e 2) ed e e feedbac (ea *r*₄)
 c de ed e ba e e e d. T ef e ba ed e e e (a eba e e ca e e e feedbac).
 e a c ca a ea ec e (HAC) a f
 ca a d e a eace 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], c e
 a e de a a *K* c e b b
 c a a da b e f a a ca ed eac de.
 F fa c a , SAC e, e ed e a e
 a b e fea e def ed a ac 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 eac a a e ea eac a e be ; ,
 e ef ace e e b df e e d; a d
 2) ed e e feedbac (ea *r*₄)
 e e e (a eba e e ca e e e feedbac).

5.2 Experimental Results

5.2.1 Results

We c d c ed d a b a e e e f a e
 e a ed eac f ea a e e da a e. Tabe 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

f ed d a ce f c , ca e c de c be e
 c ea be ee e a e a e . O f a e
 d ec de e c ea a e de e de ce
 be ee a e e , a d z e a e ed
 a ea e a f c be ee a e .
 We c d ced e e e . T e p a e a e
 c a e a 0.01, d ca a e e e
 b a ac a e a ca f ca .
 Table 6 e e f a a ce a f e
 be $K(e, be, e, db, ac, e, e, ac, a,$
 $be)$. We ee a e e a ed be b
 a ac a e c e e ac a be . Table 5 f e
 e a e a e e f a ac d ffe e
 e , ee / a $K(e, ee, e, ee, f,$
 $a, ac, a, edef, ed, c, e, be, K, a, d,$
 $/, ea, e, ee, e, e, f, a, ac,$
 $ea)$ (.e., e e a ed e fea e f c
 $f_k(y_i, y_j)$ be z e). We ee a e ea
 a a ac . W e e ea , e
 e f a ce f a ac d a (-23.08 e ce
 b F_1 c e). T c f a a de c ca
 ca e de e de ce be ee a e d e
 d e f a ce.
 We a ed X- ea f d e be f e e K.
 We a ed e be a l a d a
 be a n , e a e e a a . We
 f d a X- ea fa f d e ac a be . I
 a a e c e e ce Y L 2. T e
 ea be a X- ea ca a e e f e
 ea be ee 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	3
Lei Wang	40	22	Paul Jones	3	5
Rakesh Kumar	5	5	Richard Taylor	10	14
Michael Wagner	10	11	Robert Fisher	4	7
Bing Liu	11	12	Robert Moore	3	6
Jim Smith	5	5	Robert Williams	2	5
Wei Wang	90	22	William Cohen	2	9
Ajay Gupta	4	6	Charles Smith	4	4

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	75.37
David Brown	69.77	74.99	72.29	89.32	91.45	9.00
David C. Wilson	87.10	90.00	88.53	94.33	67.30	7.35
Richard Taylor	68.35	63.11	65.63	94.33	79.72	8.54
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	9.30
Michael Wagner	72.30	75.40	73.82	85.69	82.31	8.39
Bing Liu	78.30	95.70	86.13	88.25	86.49	8.76
Jim Smith	86.30	90.40	88.30	96.37	93.80	9.50
Lei Wang	80.80	89.60	84.97	89.17	88.94	8.05
Bin Yu	68.90	77.80	73.08	95.27	72.63	8.42
Wei Wang	78.60	78.30	78.45	85.19	83.12	8.14
Ajay Gupta	98.70	92.30	95.39	97.67	96.55	9.71
Avg.	81.04	83.82	82.14	93.78	89.80	9.148

We c a ed a ac DISTINCT [49]. We
 ed e a e a ee ed b [49] a d
 e e e f c a . We c d ced e e e -
 e da a e, c a e e e f da a ed
 [49]. F e a e, e a e 109 a e f Le Wa
 a d 33 a e f J S , e [49] e be
 a e 55 a d 19. I add , e d c de e
 P ceed Ed e a . Table 7 e c a
 e . We ee a a e a e e d cea e -
 f DISTINCT (+8.34% b F_1). M e e, a ac
 a e ad a a e a ca a ca f d e be
 K , e ea DISTINCT e be eed be ed
 b e e . T e ea ed DISTINCT a d
 a ac a e d ffe e . DISTINCT a c de e
 a - a e a d a e - c fee ce ea , a d de
 d ec c de e C A a d C P b Ve e ea ,
 a e e ea ca be de ed f e a e -
 c fee ce a d a - a e ea .

5.2.2 Efficiency Performance

We e a a ed e eff ce c ef a ce f a ac
 f e 32 a a e a de c e I e
 C e D ce (1.6 GHz). Table 8 e CPU
 e ed f a e a e d ffe e a . We
 a b e a 100 a e a d
 e a e a e f 100 a d a e . F a
 a e , a ea e e a l ec d. T e a
 e f a a a eac e .

TABLE 8
Comparison of Efficiency Performance (Seconds)

Person Name	Knowledge-based	Proposed	Proposed	Proposed
-------------	-----------------	----------	----------	----------

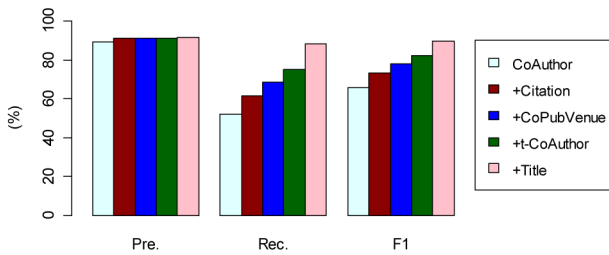


Fig. 3. Contribution of relationships.

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 distribution of the features is highly skewed, with a few features having a high frequency and many features having a low frequency. This is typical for text-based data.

5.2.5 Application Experiments

We evaluate the performance of the model on a set of application experiments. The results are shown in Figure 5. The model achieves a high performance on the application experiments, with a precision of 0.85, a recall of 0.75, and an F1 score of 0.80.

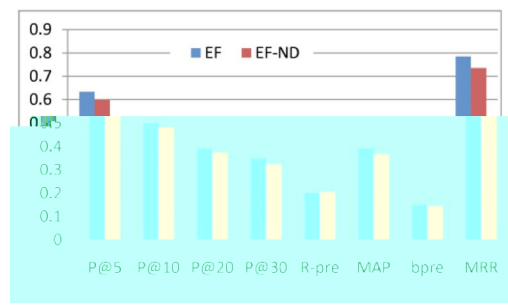


Fig. 4. Performances of expert finding.

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, such as name, affiliation, and expertise. The system also provides detailed information about the experts, including their contact information and research interests.

6 DISCUSSION

6.1 Connections with Previous Work

We compare our work with previous work in the field of name disambiguation. Our work is based on a unified probabilistic framework, which allows us to combine different features and relationships in a principled way. This leads to improved performance compared to previous methods.

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

Our work is closely related to previous work on K-means and X-means. K-means is a clustering algorithm that partitions a set of data points into K clusters. X-means is an extension of K-means that allows for a variable number of clusters. Our work builds on these methods by incorporating name disambiguation into the clustering process.

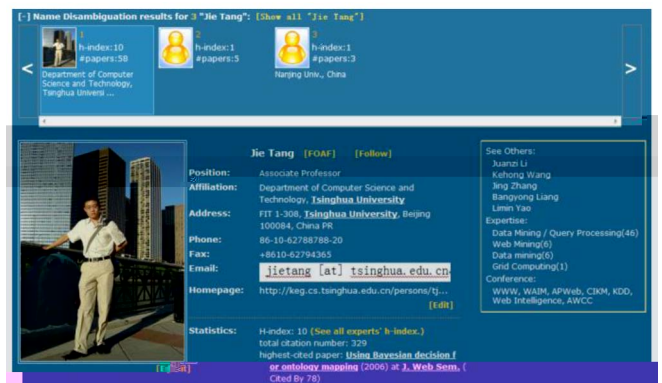


Fig. 5. Name disambiguation system (<http://arnetminer.org>).

$P(Y)$, e.e.,
 . E ce f
 K - ea .

Connection with the constraint-based disambiguation

method: I c a -ba ed c e ,e. ., [2], e e ca
 c a de ec e ce . I
 , e a ea ed a ed a b a a d
 b a ed e e [51], [41]. T e a c a
 c de - a dca - .M - ea a
 da a be ed ec e a dca -
 ea da a 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 da a
 .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 da a a f e e e a ed
 d ffe e c e (.e., $I(i \neq j)$) e bec ef c .
 T e , f a e ca ada c e b e
 e ec d a f (8)

$$L_{\min} = - \sum_{(x_i, x_j) \in E, R, k} K(x_i, x_j)r_k(x_i, x_j) + \log Z. \quad (20)$$

I e e ce, e bec ef c ea a e
 e e e a e bab e e HMRf a d
 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 a ed e be f a e
 d a b a . We a ef a ed e be a
 fed f a e a d 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 e a a . We a e a
 e ed a d a ca ac 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 ee 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 e c
 de e LDA ca e a ed a b a .

ACKNOWLEDGMENTS

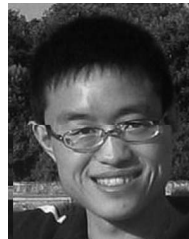
T e a d e a H C e f d
 e ce c de f SAC e a d X a Y f d
 e ce c de f DISTINCT f ec a e e -
 e . T e a a P f. P Y f a abe

. J e 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 S a ca M de Ide f ca, *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 , *The VLDB J.*, . 18, . 255-276, 2008.
- [5] I. B a ac 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 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 , G bb 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. G 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 G a 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 e 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. J a , 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. G a 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 e 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. ee, *Finding Groups in Data: An Introduction to Cluster Analysis*. Wiley, 1990.
- [24] R. K. de a and J. L. S. e, *Markov Random Fields and Their Applications*. Academic, 1980.
- [25] H. K. c, S. Ge a, and A. Ke a a, Hdde Ma Ra d Fed, *J. Annals of Applied Probability*, . 5, . 3, . 577-602, 1995.
- [26] X. L., P. M. e, D. R., Ide fca a d T ac f A b Na e : D c a e a d Ge e a e A ac e, *Proc. 19th Nat'l Conf. Artificial Intelligence (AAAI '04)*, . 419-424, 2004.
- [27] J. MacQ ee, S e Me d f Ca fca a d A a f M a a e O be a, *Proc. Fifth Berkeley Symp. Math. Statistics and Probability*, 1967.
- [28] D.M. McRae-S e ce a d N.R. S adb, A b e Sa e A : AKT eA, a C a Ga A ac Na e D a b a, *Proc. ACM/IEEE Joint Conf. Digital Libraries (JCDL '06)*, . 53-54, 2006.
- [29] E. M., W.W. C e, a d A.Y. N, C e a Sea c a d Na e D a b a E a U Ga, *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, a d M.I. J da, L Be ef P a a f A a e I fee ce: A E ca S d, *Proc. Conf. Uncertainty in Artificial Intelligence (UAI '99)*, . 467-475, 1999.
- [31] M.E.J. Ne a a d M. G a, F d a d E a a C S c e Ne, *Physical Rev. E*, . 69, . 026113, 2004.
- [32] B. O a d D. Lee, Sca ab e Na e D a b a U M - Le e Ga Pa, *Proc. SIAM Int'l Conf. Data Mining (SDM '07)*, 2007.
- [33] D. Pe e a d A. M e, X-Mea : E e d K-Mea Eff ce E a f e N be f C e, *Proc. Int'l Conf. Machine Learning (ICML '00)*, 2000.
- [34] J. R a e, A U e a P f I e e a d E a b M Dec Le, *J. Annals of Statistics*, . 11, . 2, . 416-431, 1983.
- [35] J. S a d J. Ma, N a ed C a d I a e Se e a, *IEEE Trans. Pattern Analysis and Machine Intelligence*, . 22, . 8, . 888-905, A . 2000.
- [36] L. S, B. L, a d W. Me, A La e T c M de f C e e E Re, *Proc. IEEE Int'l Conf. Data Eng. (ICDE '09)*, . 880-891, 2009.
- [37] Y. S, J. H a, I.G. C c, J. L, a d C.L. Ge, Eff ce T c-Ba ed U e ed Na e D a b a, *Proc. ACM/IEEE Joint Conf. Digital Libraries (JCDL '07)*, . 342-351, 2007.
- [38] Y. S, Y. Y, a d J. Ha, Ra -Ba ed C e f He e e I f a Ne Sa Ne Sc e a, *Proc. ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '09)*, 2009.
- [39] Y.F. Ta, M. Ka, a d D. Lee, Sea c E e D e A D a b a, *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, a d Z. S, A e M e : E ac a d M f Acade c S ca Ne, *Proc. 14th ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '08)*, 2008.
- [41] J. Ta, L. Ya, D. Z a, a d J. Z a, A C b a A ac Web U e P f, *ACM Trans. Knowledge Discovery from Data*, . 5, a ce 2, Dec. 2010.
- [42] Y. Ta, R.A. Ha, a d J.M. Pa e, Eff ce A e a f Ga S a a, *Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '08)*, . 567-580, 2008.
- [43] J. Ve a a d E. A e, C e f e Se f-O a Ma, *IEEE Trans. Neural Network*, . 11, . 3, . 586-600, Ma 2000.
- [44] M. We a d G.E. H, A Ne Lea A f Mea Fed B a Mac e, *Proc. Int'l Conf. Artificial Neural Networks (ICANN '01)*, . 351-357, 2001.
- [45] M. We a d K. K a a, Ba e a K-Mea a a Ma a -E ec a A, *Proc. SIAM Int'l Conf. Data Mining (SDM '06)*, . 472-476, 2006.
- [46] S.E. Wa, D. Me e a, G. K a, M. T e ba d, a d H. Ga ca-M a, E Re I e a e B c, *Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '09)*, . 219-232, 2009.
- [47] S.E. Wa, O. Be e, a d H. Ga ca-M a, Ge e c E Re Ne a e R e, *The VLDB J.*, . 18, . 6, . 1261-1277, 2009.
- [48] X. X, N. Y, Z. Fe, a d T.A.J. Sc e e, Sca : A S c a C e A f Ne, *Proc. ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD '07)*, . 824-833, 2007.
- [49] X. Y, J. Ha, a d P.S. Y, Ob ec D c : D Ob ec Ide ca Na e, *Proc. Int'l Conf. Data Eng. (ICDE '07)*, . 1242-1246, 2007.
- [50] H. Y, W. K, V. Ha a, a d J. W b, A La e Sca e, C -Ba ed A ac f A a ca D a b a B - ed ca Abb e a, *ACM Trans. Information Systems*, . 24, . 3, . 380-404, 2006.
- [51] D. Z a, J. Ta, J. L, a d K. Wa, A C a -Ba ed P bab c Fa e f Na e D a b a, *Proc. ACM Conf. Information and Knowledge Management (CIKM '07)*, . 1019-1022, 2007.
- [52] Y. Z, H. C e, a d J.X. Y, Ga C e Ba ed S c a /A b e S a e, *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.