

Incorporating World Knowledge to Document Clustering via Heterogeneous Information Networks

[Extended Version]

Chenguang Wang[†], Yangqiu Song[‡], Ahmed El-Kishky[‡], Dan Roth[‡], Ming Zhang[†], Jiawei Han[‡]

[†]School of EECS, Peking University

[‡]Department of Computer Science, University of Illinois at Urbana-Champaign

wangchenguang@pku.edu.cn, {yqsong, elkishk2, danr, hanj}@illinois.edu, mzhang@net.pku.edu.cn

ABSTRACT

One of the key obstacles in making learning protocols realistic in applications is the need to supervise them, a costly process that often requires hiring domain experts. We consider the framework to use the world knowledge as indirect supervision. World knowledge is general-purpose knowledge, which is not designed for any specific domain. Then the key challenges are how to adapt the world knowledge to domains and how to represent it for learning. In this paper, we provide an example of using world knowledge for domain dependent document clustering. We provide three ways to specify the world knowledge to domains by resolving the ambiguity of the entities and their types, and represent the data with world knowledge as a heterogeneous information network. Then we propose a clustering algorithm that can cluster multiple types and incorporate the sub-type information as constraints. In the experiments, we use two existing knowledge bases as our sources of world knowledge. One is Freebase, which is collaboratively collected knowledge about entities and their organizations. The other is YAGO2, which is a knowledge base automatically extracted from Wikipedia and maps the knowledge to the linguistic knowledge base, WordNet. Experimental results on two text benchmark datasets (20newsgroups and RCV1) show that incorporating world knowledge as indirect supervision can significantly outperform the state-of-the-art clustering algorithms as well as clustering algorithms enhanced with world knowledge features.

1. INTRODUCTION

Machine learning algorithms have become pervasive in multiple domains and have started to impact applications. Nonetheless, a key obstacle in making learning protocols realistic in applications is the need to supervise them, a costly process that often requires hiring domain experts. In the past decades, machine learning community has elaborated to reduce the labeling work done by human for supervised machine learning algorithms or to improve unsupervised learning with only minimum supervision. For example, semi-supervised learning [8] is proposed to use only partially labeled data and a lot of unlabeled data to perform learning with the

hope that it can perform as good as fully supervised learning. Transfer learning [34] uses the labeled data from other relevant domains to help the learning task in the target domain. However, there are still many cases that neither semi-supervised learning nor transfer learning can help. For example, in the era of big data, we can have a lot textual information from different Web sites, e.g., blogs, forums, mailing lists, etc.. It is impossible to ask human to annotate all the required tasks. It is also difficult to find relevant labeled domains. Some domains can be very specific and really need the domain experts to perform annotation, e.g., the medical domain publication classification. Therefore, we should consider a more general approach to further reduce the labeling cost for learning tasks in diverse domains.

Fortunately, nowadays we possess an abundance of general-purpose knowledge bases, e.g., Cyc project [25], Wikipedia, Freebase [6], KnowItAll [12], TextRunner [2], WikiTaxonomy [35], DBpedia [1], YAGO [40], NELL [7] and Knowledge Vault [11]. We call these knowledge bases world knowledge [14], because they are universal knowledge that are either collaboratively annotated by human labelers or automatically extracted from big data. When world knowledge is annotated or extracted, it is not collected for any specific domain. However, we believe the facts in world knowledge bases are very useful and often of high quality. Therefore we consider using them as supervision for many machine learning problems. People have found it useful to use world knowledge as distant supervision for entity and relation extraction [32]. This is a direct use of the facts in world knowledge bases, where the entities in the knowledge bases are matched in the context regardless the ambiguity. A more interesting question is *can we use the world knowledge to “supervise” more machine learning algorithms or applications?* Particularly, if we can use world knowledge as indirect supervision, then we can extend the knowledge about entities and relations to more generic text analytics problems, e.g., categorization and information retrieval.

Thus, we consider a general world knowledge enabled machine learning framework, that can incorporate world knowledge into machine learning algorithms. As mentioned, world knowledge is not designed for any specific domain. For example, when we want to cluster the documents about entertainment or sports, then the world knowledge about names of celebrities and athletes may help while the terms used in science and technology may not be very useful. Thus, a key issue is how we should adapt world knowledge to the domain specific tasks. Another problem is when we have the world knowledge, how we can represent it for the domain dependent tasks. For example, because most of the knowledge bases use a linked network to organize the knowledge, to adapt the world knowledge to domains, we should consider how to use the linked data. Although traditional machine learning algorithms using world

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

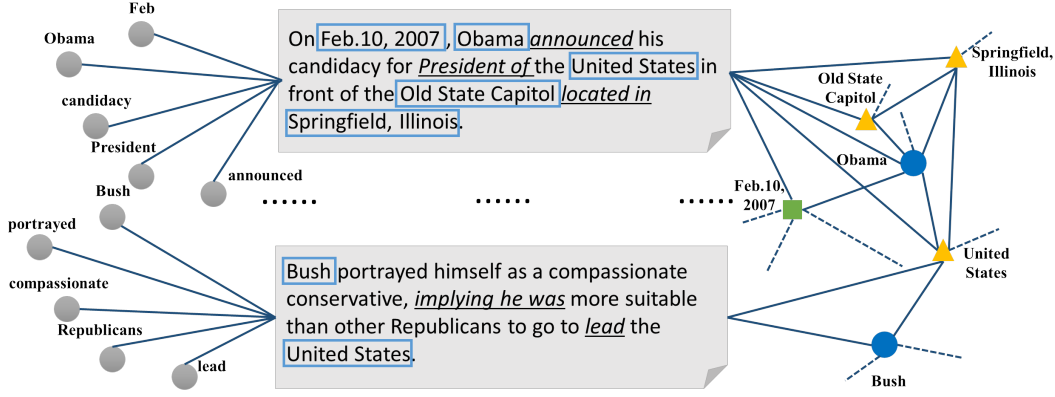


Figure 1: Heterogeneous information network example. The network \mathcal{G} contains five entity types: document, word, date, person and location, which are represented with gray rectangle, gray round, green square, blue round, and yellow triangle, respectively.

knowledge just treat world knowledge as “flat” features in addition to the original text data [14, 29], the structure of the knowledge provides rich information about the connections of entities and relations. Therefore, we should also carefully consider the best way to represent the world knowledge for machine learning algorithms.

In this paper, we illustrate a framework of machine learning with world knowledge using a document clustering problem. We select two knowledge bases, i.e., Freebase, YAGO2, as the sources of world knowledge. Freebase [6] is a collaboratively collected knowledge base about entities and their organizations. YAGO2 [40] is a knowledge base automatically extracted from Wikipedia and maps the knowledge to the linguistic knowledge base, WordNet [13]. To adapt the world knowledge to domain specific tasks, we first use semantic parsing to ground any text to the knowledge bases [5]. Then we propose to use frequency, document frequency, and conceptualization [38] based semantic filters to resolve the ambiguity problem when adapting world knowledge to the domain tasks. After that, we have the documents as well as the extracted entities and their relations. Since the knowledge bases provide the entity types, the resulting data naturally form a heterogeneous information network (HIN) [17]. We show an example of such HIN in Figure 1. The specified world knowledge, such as named entities (“Bush”, “Obama”) and their types (*Person*), as well as the documents and the words form the HIN. We then formulate the document clustering problem as an HIN partitioning problem, and provide a new algorithm to better perform clustering by incorporating the rich structural information as constraints in the HIN. For example, the HIN builds a link (a must-link constraint) between “Obama” of sub-type *Politician* in one document and “Bush” of sub-type *Politician* in another document. Such link and type information could be very useful if the target clustering domain is “Politics.”

The main contributions of this work can be highlighted as follows:

- We study a novel problem of supervising machine learning algorithms with world knowledge.
- We propose to use semantic parsing and semantic filtering to specify world knowledge to the domain dependent documents, and develop a new constrained HIN clustering algorithm to make better use of the structural information from the world knowledge for document clustering task.
- We conduct experiments on two benchmark datasets (20news-groups and RCV1) to evaluate the clustering algorithm using

HIN, compared with the state-of-the-art document clustering algorithms and clustering with “flat” world knowledge features. We show that our approach can be 13.3% better than the semi-supervised clustering algorithm incorporating 250K constraints which are generated by ground-truth labels.

2. TO LEARN WITH WORLD KNOWLEDGE

In this section, we discuss how we enable world knowledge to indirectly “supervise” machines, instead of just using world knowledge as additional features. In general, performing machine learning with world knowledge, we should follow four steps: (1) Knowledge acquisition. (2) Domain adaptation. (3) Data and knowledge representation. (4) Learning. Since we assume the world knowledge is given, we skip step one in this study. Then given the world knowledge, we should consider to adapt it to specific domains. Since the knowledge can be ambiguous without context, we should consider to use domain dependent data to find the best knowledge to use. For example, when a text mentioning “apple,” it can refer to a company or a fruit. In the knowledge base, we have both. Therefore, we should choose the right one to use. Then given the filtered knowledge we have as well as the domain dependent data, we use a better representation which considers the structure information of the linked knowledge rather than just considering the knowledge as flat features. After we have the representation, we can design a learning algorithm for domain dependent task.

The above four steps are general, which means they may apply to many applications. In this section, we introduce when we are given a document clustering problem, how we choose the right knowledge to use and how we represent it given the knowledge. Then in the next section, we will introduce the learning algorithm to perform better document clustering given the representation.

2.1 World Knowledge Specification

In this subsection, we propose a world knowledge specification approach to generate specified world knowledge given a set of domain dependent documents. We first use semantic parsing to ground any text to the knowledge base, then provide three semantic filtering approaches to avoid ambiguity of the extracted information.

2.1.1 Semantic Parsing

Semantic parsing is the task of mapping a piece of natural language text to a formal meaning representation [33]. This can support question answering by querying a knowledge base [23]. To our

best knowledge, most previous semantic parsing algorithms or tools developed are for small scale problems but with complicated logical forms, until Berant et al. [5] develop a system that can handle very large scale knowledge bases such as Freebase. They use the developed system to solve question answering problem with Freebase. In their work, they formulate their problem to match answers to the questions, which is a supervised learning process. Similar to them, we are also working with very large scale world knowledge bases, but unlike them, we do not match question and answers. Our task is to ground any text to the knowledge base entities and their relationships in the prescribed logical form. Therefore, our problem is a fully unsupervised problem.

We first introduce the problem formulation and then introduce how we perform unsupervised semantic parsing. Let \mathcal{E} be a set of entities and \mathcal{R} be a set of relations in the knowledge base. Then the knowledge base \mathcal{K} consists of triplets in the form of (e_1, r, e_2) , where $e_1, e_2 \in \mathcal{E}$ and $r \in \mathcal{R}$. We follow [5] to use a simple version of Lambda Dependency-Based Compositional Semantics (λ -DCS) [28] as the logic language. From each sentence in the document, we can parse four possible λ -DCS logic forms [5]: (1) Unary: an entity e is a unary logic form (e.g., *Obama*); (2) Binary: a relation r is a binary logic form (e.g., *PresidentofCountry*); (3) Join: $r.e$ is a unary logic form, denoting a join, where r is a binary and e is a unary (e.g., *PresidentofCountry.Obama*); (4) Intersection: $e_1 \cap e_2$ ($e_1, e_2 \in \mathcal{E}$) denotes set intersection, where e_1 and e_2 are both unaries (e.g., *Location.Olympics* \cap *PresidentofCountry.Obama*).

In simpler terms, semantic parsing can be understood as the following process. First, given “Obama is the president of United States of America,” it maps the entities, as well as the relation phrases in the text to knowledge base. So “Obama” and “United States of America” are mapped to knowledge base, resulting in two unary logic forms *People.BarackObama* and *Country.USA*, where *People* and *Country* are the type information in Freebase. The relation phrase “president” is mapped to a binary logic form *PresidentofCountry*. Notice that, the mapping process skips the words “is” and “of.” The mapping dictionary is constructed by aligning a large text corpus to the knowledge base. A phrase and a knowledge base entity or relation align if they co-occur with many of the same entities. We select two knowledge bases, i.e., Freebase and YAGO2. For Freebase, we just use the mapping already existing in the released tool shown in [5]. For YAGO2, we follow [5] and download a subset of ClueWeb09¹ to find the new mapping for YAGO2 entities and relations. Second, it uses some rules (i.e., grammar) to combine the basic logic forms to generate the restricted four logic forms above, and rank the results. For the example shown in this paragraph, *People.BarackObama* \wedge *President.USA* is generated to represent its semantic meaning. Notice that, *President.USA* is generated by joining the unary *Country.USA* with the binary *PresidentofCountry*.

When there are more than one candidate semantic meanings for a sentence, in [5], they learn the ranks based on the annotated question-answer pairs. For our task, this annotation is not available. Therefore, instead of ranking or enumerating all the possible logic forms (which is found to be not feasible in limited time), we constrain the entities to be the maximum length spanning phrases recognized by a state-of-the-art named entity recognition tool [36]. We then perform the two steps introduced above by using the maximum length spanning noun phrase as entities, and use the phrase between them in the text as relation phrase. From the multiple resulting meaning representations of the given text, we propose to use the following three semantic filtering methods to resolve the ambiguity prob-

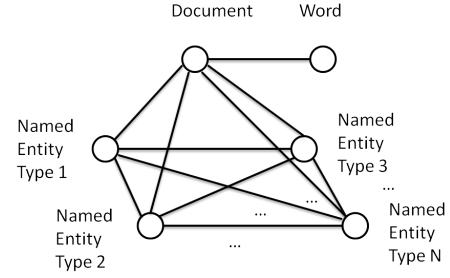


Figure 2: Heterogeneous information network schema. The specified knowledge is represented in the form of heterogeneous information network. The schema contains multiple entity types: document \mathcal{D} , word \mathcal{W} , named entities $\{\mathcal{E}^i\}_{i=1}^T$, and the relation types connecting the entity types.

lem.

2.1.2 Semantic Filtering

For each sentence in the given document, the output of semantic parsing is a set of logic forms that represent the semantic meaning. However, the extracted entities can be ambiguous. For example, “apple” may be associated with type *Company* or *Fruit*. Therefore, we should filter out the noisy entities and their types to ensure that the knowledge we have is good enough as indirect supervision for document clustering. We assume that in the domain specific tasks, given the context, the entities seldom have multiple meanings. Thus, we propose the following three approaches to select the best knowledge to use for further learning process.

Frequency based semantic filter (FBSF). We use the frequency of a type for an entity as the criterion to decide whether the entity should be extracted for the domain specific task in a sentence. Here we assume the most frequent type of an entity from all the sentences of the document is the correct semantic meaning.

Document frequency based semantic filter (DFBSF). Similar to the frequency based method, we use the document frequency (DF) of a type of an entity as the criterion to find the most likely semantic meaning. Here we assume that if an entity appears in multiple documents with the same type, then the type should be the correct semantic meaning.

Conceptualization based semantic filter (CBSF). Motivated by the approaches of conceptualization [38] and entity disambiguation [27], we represent each entity with a feature vector of entity types, and use standard Kmeans to cluster the entities. Then in each cluster, we use the intersection operation to find the most likely entity type for the entities in the cluster. In this case, different entities can be used to disambiguate each other. Here we assume that the type that can best fit the context is the correct semantic meaning.

2.2 World Knowledge Representation

The output of semantic parsing and semantic filtering is then the document associated with the entities, which are further associated with the types (or concepts, categories, the names can be different for different knowledge bases). For example, in Freebase, we select the top level named entity categories (i.e., domains) as the types, e.g., *Person*, *Location*, and *Organization*. In addition to the named entities, we also regard the document and word as two types. Then we use an HIN to represent the data we get after semantic parsing and semantic filtering.

DEFINITION 1. A **heterogeneous information network (HIN)** is a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with an entity type mapping $\phi: \mathcal{V} \rightarrow \mathcal{A}$ and

¹<http://www.lemurproject.org/clueweb09.php/>

a relation type mapping $\psi: \mathcal{E} \rightarrow \mathcal{R}$, where \mathcal{V} denotes the entity set and \mathcal{E} denotes the link set, \mathcal{A} denotes the entity type set and \mathcal{R} denotes the relation type set, and the number of entity types $|\mathcal{A}| > 1$ or the number of relation types $|\mathcal{R}| > 1$.

The network schema provides a high-level description of a given heterogeneous information network.

DEFINITION 2. Given an HIN $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with the entity type mapping $\phi: \mathcal{V} \rightarrow \mathcal{A}$ and the relation type mapping $\psi: \mathcal{E} \rightarrow \mathcal{R}$, the **network schema** for network \mathcal{G} , denoted as $\mathcal{T}_{\mathcal{G}} = (\mathcal{A}, \mathcal{R})$, is a graph with nodes as entity types from \mathcal{A} and edges as relation types from \mathcal{R} .

Then for our world knowledge dependent network, we use the network schema shown in Figure 2 to represent the data. The network contains multiple entity types: **document** \mathcal{D} , **word** \mathcal{W} , **named entities** $\{\mathcal{E}^t\}_{t=1}^T$, and a few relation types connecting the entity types. We denote the document set as $\mathcal{D} = \{d_1, d_2, \dots, d_M\}$, where M is the size of \mathcal{D} , the word set as $\mathcal{W} = \{w_1, w_2, \dots, w_N\}$, where N is the size of \mathcal{W} , and the entity set as $\mathcal{E}^t = \{e_1^t, e_2^t, \dots, e_{V_t}^t\}$, where V_t is the size of \mathcal{E}^t . We have $t = 1, \dots, T$ where T is the total number of named entity types we find in the knowledge base. Note that if there are no named entities, then the network reduces to a bipartite graph containing only documents and words.

3. DOCUMENT CLUSTERING WITH WORLD KNOWLEDGE

In this subsection, we present our clustering algorithm using HIN, constructed from domain dependent documents and the world knowledge. Given the HIN, it is natural to perform HIN partitioning to obtain the document clusters. In addition to the HIN itself, let us revisit the structural information in a typical world knowledge base, e.g., Freebase. In the world knowledge base, the named entities are often organized in a hierarchy of categories. Although there are additional category information for each entity, we only use the top level named entity types as the entity types in HIN. For example, “Barack Obama” is a person, where person is the top level category. In addition, he is the president of the “United States,” a politician, a celebrity, etc.. Another example is that “Google” is a software company, plus it has a CEO. This shows that the entities can have some attributes. We choose to use top level entity types for the HIN schema since then we will have a dense graph for each pairwise nodes in the network schema. The fine-grained named entity sub-types or the attributes are also very useful to identify the topics or the clusters of the documents. Therefore, in this section, we introduce how we incorporate the fine-grained level of named entity types as constraints in the HIN clustering algorithm.

3.1 Constrained Clustering Modeling

To formulate the clustering algorithm for the domain dependent documents, we denote latent label sets of the documents as $\mathcal{L}_d = \{l_{d_1}, l_{d_2}, \dots, l_{d_M}\}$. We also denote $\mathcal{L}_w = \{l_{w_1}, l_{w_2}, \dots, l_{w_N}\}$ for words, and $\mathcal{L}_{e^t} = \{l_{e_1^t}, l_{e_2^t}, \dots, l_{e_{V_t}^t}\}$ for the t^{th} named entities set. In general, we follow the framework of information-theoretic co-clustering (ITCC) [10] and constrained ITCC [37] to formulate our approach. Instead of only performing on the bipartite graph, we need to handle multi-type relational data, as well as more complicated constraints.

The original ITCC uses a variational function to approximate the joint probability of documents and words, which is:

$$q(d_m, w_i) = p(\hat{d}_{k_d}, \hat{w}_{k_w}) p(d_m | \hat{d}_{k_d}) p(w_i | \hat{w}_{k_w}), \quad (1)$$

where \hat{d}_{k_d} and \hat{w}_{k_w} are cluster indicators to formulate the conditional probability, and k_d and k_w are the corresponding cluster indices. $q(d_m, w_i)$ is used to approximate $p(d_m, w_i)$ by minimizing the Kullback-Leibler (KL) divergence:

$$\begin{aligned} & D_{KL}(p(\mathcal{D}, \mathcal{W}) || q(\mathcal{D}, \mathcal{W})) \\ &= D_{KL}(p(\mathcal{D}, \mathcal{W}, \hat{\mathcal{D}}, \hat{\mathcal{W}}) || q(\mathcal{D}, \mathcal{W}, \hat{\mathcal{D}}, \hat{\mathcal{W}})) \\ &= \sum_{k_d}^{K_d} \sum_{d_m: l_{d_m} = k_d} p(d_m) D_{KL}(p(\mathcal{W} | d_m) || p(\mathcal{W} | \hat{d}_{k_d})) \\ &= \sum_{k_w}^{K_w} \sum_{w_i: l_{w_i} = k_w} p(w_i) D_{KL}(p(\mathcal{D} | w_i) || p(\mathcal{D} | \hat{w}_{k_w})) \end{aligned} \quad (2)$$

where $\hat{\mathcal{D}}$ and $\hat{\mathcal{W}}$ are the cluster sets, $p(\mathcal{W} | \hat{d}_{k_d})$ denotes a multinomial distribution based on the probabilities

$$p(\mathcal{W} | \hat{d}_{k_d}) = (p(w_1 | \hat{d}_{k_d}), \dots, p(w_N | \hat{d}_{k_d}))^T,$$

$$p(w_i | \hat{d}_{k_d}) = p(w_i | \hat{w}_{k_w}) p(\hat{w}_{k_w} | \hat{d}_{k_d}) \text{ and,}$$

$$p(w_i | \hat{w}_{k_w}) = p(w_i) / p(l_{w_i} = \hat{w}_{k_w}).$$

Symmetrically, we have

$$p(\mathcal{D} | \hat{w}_{k_w}) = (p(d_1 | \hat{w}_{k_w}), \dots, p(d_M | \hat{w}_{k_w}))^T,$$

$$p(d_i | \hat{w}_{k_w}) = p(d_i | \hat{d}_{k_d}) p(\hat{d}_{k_d} | \hat{w}_{k_w}) \text{ and,}$$

$$p(d_i | \hat{d}_{k_d}) = p(d_i) / p(l_{d_i} = \hat{d}_{k_d}).$$

Moreover, $p(\hat{w}_{k_w} | \hat{d}_{k_d})$ and $p(\hat{d}_{k_d} | \hat{w}_{k_w})$ are computed based on the joint probability $q(\hat{d}_{k_d}, \hat{w}_{k_w}) = \sum_{l_{d_m} = k_d} \sum_{l_{w_i} = k_w} p(d_m, w_i)$.

Our problem in HIN is that each edge in the network schema can be represented as a bipartite graph. The co-occurrence frequency of the semantic parsing results is used as the edge weights in the bipartite graph. Motivated by ITCC, according to the network schema shown in Figure 2, our problem of HIN clustering is formulated as

$$\begin{aligned} \mathcal{J}_{\text{HINC}} &= D_{KL}(p(\mathcal{D}, \mathcal{W}) || q(\mathcal{D}, \mathcal{W})) \\ &+ \sum_{t=1}^T D_{KL}(p(\mathcal{D}, \mathcal{E}^t) || q(\mathcal{D}, \mathcal{E}^t)) \\ &+ \sum_{t=1}^T \sum_{s=1}^T D_{KL}(p(\mathcal{E}^t, \mathcal{E}^s) || q(\mathcal{E}^t, \mathcal{E}^s)). \end{aligned} \quad (3)$$

where all the probabilities can be defined similar to the document-word bipartite graph. We omit the detailed definitions due to the space limitation. A summary of the notations is shown in Table 1.

Table 1: Notations for clustering algorithm. The indicators are used for the probability representation, while the indices are used as ids for the clusters.

Meaning	Document	Word	Entity
Cluster Index	k_d	k_w	k_{e^t}
Cluster Indicator	\hat{d}_{k_d}	\hat{w}_{k_w}	$\hat{e}_{k_{e^t}}^t$
Data Indicator	d_m	w_i	e_i^t
Data Indicator Set	\mathcal{D}	\mathcal{W}	\mathcal{E}^t
Label	l_{d_m}	l_{w_i}	$l_{e_i^t}$
Label Indicator Set	\mathcal{L}_d	\mathcal{L}_w	\mathcal{L}_{e^t}

To incorporate the side information of the fine-grained named entity sub-types or the attributes as indirect supervision for document clustering, we define the constraints for the named entities we find after semantic parsing. We take the t^{th} entity label set \mathcal{E}^t as an example, and use must-links and cannot-links as the constraints. We denote the must-link set associated with e_i^t as $\mathcal{M}_{e_i^t}$, and the cannot-link set as $\mathcal{C}_{e_i^t}$. For must-links, the cost function is defined as

$$\begin{aligned} & V_{\mathcal{M}}(e_{i_1}^t, e_{i_2}^t \in \mathcal{M}_{e_{i_1}^t}) \\ &= w_{\mathcal{M}} D_{KL}(p(\mathcal{D} | e_{i_1}^t) || p(\mathcal{D} | e_{i_2}^t)) \cdot \mathcal{I}_{l_{e_{i_1}^t} \neq l_{e_{i_2}^t}}, \end{aligned} \quad (4)$$

Input: HIN defined on documents \mathcal{D} , words \mathcal{W} , and entities $\mathcal{E}^t, t = 1, \dots, T$; Set maxIter and $\text{max}\delta$.

while $\text{iter} < \text{maxIter}$ and $\delta > \text{max}\delta$ **do**

D Label Update: minimize Eq. (7) w.r.t. \mathcal{L}_d .

D Model Update: update $q(d_m, w_i)$ and $q(d_m, e_i^t)$.

for $t = 1, \dots, T$ **do**

\mathcal{E}^t Label Update: minimize Eq. (9) w.r.t. \mathcal{L}_{e^t} .

\mathcal{E}^t Model Update: update $q(d_m, e_i^t)$ and $q(e_j^s, e_i^t)$.

end for

D Label Update: minimize Eq. (7) w.r.t. \mathcal{L}_d .

D Model Update: update $q(d_m, w_i)$ and $q(d_m, e_i^t)$.

W Label Update: minimize Eq. (8) w.r.t. \mathcal{L}_w .

W Model Update: update $q(d_m, w_i)$.

 Compute cost change δ using Eq. (6).

end while

Algorithm 1: Alternating Optimization for CHINC.

where $w_{\mathcal{M}}$ is the weight for must-links, and $p(\mathcal{D}|e_{i_1}^t)$ denotes a multinomial distribution based on the probabilities $(p(d_1|e_{i_1}^t), \dots, p(d_M|e_{i_1}^t))^T$, and $\mathcal{I}_{true} = 1, \mathcal{I}_{false} = 0$. The above must-link cost function means that if the label of $e_{i_1}^t$ is not equal to the label of $e_{i_2}^t$, then we should take into account the cost function of how dissimilar the two entities $e_{i_1}^t$ and $e_{i_2}^t$ are. The dissimilarity is computed based on the probability of document \mathcal{D} given the entities $e_{i_1}^t$ and $e_{i_2}^t$ as Eq. (4). The more dissimilar the two entities are, the larger cost is imposed.

For cannot-links, the cost function is defined as

$$V_C(e_{i_1}^t, e_{i_2}^t \in \mathcal{C}_{e_{i_1}^t}) = w_C(D_{max}^t - D_{KL}(p(\mathcal{D}|e_{i_1}^t)||p(\mathcal{D}|e_{i_2}^t))) \cdot \mathcal{I}_{e_{i_1}^t \neq e_{i_2}^t}, \quad (5)$$

where w_C is the weight for cannot-links, and D_{max}^t is the maximum value for all the $D_{KL}(p(\mathcal{D}|e_{i_1}^t)||p(\mathcal{D}|e_{i_2}^t))$. The cannot-link cost function means that if the label of $e_{i_1}^t$ is equal to the label of $e_{i_2}^t$, then we should take into account the cost function of how similar they are.

Integrating the constraints for $\mathcal{L}_{e^1}, \dots, \mathcal{L}_{e^T}$ to Eq. (3), the objective function of constrained HIN clustering is:

$$\begin{aligned} \mathcal{J}_{CHINC} &= D_{KL}(p(\mathcal{D}, \mathcal{W})||q(\mathcal{D}, \mathcal{W})) \\ &+ \sum_{t=1}^T D_{KL}(p(\mathcal{D}, \mathcal{E}^t)||q(\mathcal{D}, \mathcal{E}^t)) \\ &+ \sum_{t=1}^T \sum_{s=1}^T D_{KL}(p(\mathcal{E}^t, \mathcal{E}^s)||q(\mathcal{E}^t, \mathcal{E}^s)) \\ &+ \sum_{t=1}^T \sum_{e_{i_1}^t=1}^{V_t} \sum_{e_{i_2}^t \in \mathcal{M}_{e_{i_1}^t}} V_{\mathcal{M}}(e_{i_1}^t, e_{i_2}^t \in \mathcal{M}_{e_{i_1}^t}) \\ &+ \sum_{t=1}^T \sum_{e_{i_1}^t=1}^{V_t} \sum_{e_{i_2}^t \in \mathcal{C}_{e_{i_1}^t}} V_C(e_{i_1}^t, e_{i_2}^t \in \mathcal{C}_{e_{i_1}^t}). \end{aligned} \quad (6)$$

From this objective function we can see that, the must-links and cannot-links are imposed to the entities that the semantic parsing detects. Since the task is document clustering, the sub-types of entities serve as indirect supervision because they cannot directly affect the cluster labels of the documents. However, the constraints can affect the labels of entities, and then the labels of entities can be transferred to the document side to affect the labels of documents.

3.2 Alternating Optimization

Since global optimization of all the latent labels as well as the approximate function $q(\cdot, \cdot)$ is intractable, we perform an alternating optimization shown in Algorithm 1. We iterate the process to optimize the labels of documents, words, and entities. Meanwhile, we update the function $q(\cdot, \cdot)$ for the corresponding types.

For example, to find label l_{d_m} of document d_m , we have:

$$l_{d_m} = \arg \min_{l_{d_m}=k_d} D_{KL}(p(\mathcal{W}|d_m)||p(\mathcal{W}|\hat{d}_{k_d})) + \sum_{t=1}^T D_{KL}(p(\mathcal{E}^t|d_m)||p(\mathcal{E}^t|\hat{d}_{k_d})) \quad (7)$$

To find label l_{w_i} of word w_i , we have:

$$l_{w_i} = \arg \min_{l_{w_i}=k_w} D_{KL}(p(\mathcal{D}|w_i)||p(\mathcal{D}|\hat{w}_{k_w})) \quad (8)$$

To find the label $l_{e_i^t}$, we use the iterated conditional mode (ICM) algorithm [4] to iteratively assign a label to the entity. We update one label $l_{e_i^t}$ at a time, and keep all the other labels fixed:

$$\begin{aligned} l_{e_i^t} &= \arg \min_{l_{e_i^t}=k_{e^t}} D_{KL}(p(\mathcal{D}|e_i^t)||p(\mathcal{D}|\hat{e}_{k_{e^t}}^t)) \\ &+ \sum_{s=1}^T D_{KL}(p(\mathcal{E}^s|e_i^t)||p(\mathcal{E}^s|\hat{e}_{k_{e^t}}^t)) \\ &+ \sum_{e_{i'}^t \in \mathcal{M}_{e_i^t};} w_{\mathcal{M}} D_{KL}(p(\mathcal{D}|e_i^t)||p(\mathcal{D}|e_{i'}^t)) \\ &\quad \mathcal{I}_{l_{e_i^t} \neq l_{e_{i'}^t}} \\ &+ \sum_{e_{i'}^t \in \mathcal{C}_{e_i^t};} w_C (D_{max}^t - D_{KL}(p(\mathcal{D}|e_i^t)||p(\mathcal{D}|e_{i'}^t))) \\ &\quad \mathcal{I}_{l_{e_i^t} = l_{e_{i'}^t}} \end{aligned} \quad (9)$$

To transfer the original objective function (6) to Eq. (9), we should follow Eq. (2) where we replace the document and word notations to the entity notations. To understand why Eq. (2) holds, we suggest to refer to the original ITCC for detailed derivation [10].

Then, with the labels $\mathcal{L}_d, \mathcal{L}_{e^t}$ and \mathcal{L}_w fixed, we update the model function $q(d_m, w_i)$, $q(d_m, e_i^t)$, and $q(e_j^s, e_i^t)$. The update of q is not influenced by the must-links and cannot-links. Thus we can modify them the same as ITCC [10] and only show the update of $q(d_m, e_i^t)$ here:

$$q(\hat{d}_{k_d}, \hat{e}_{k_{e^t}}^t) = \sum_{l_{d_m}=k_d} \sum_{l_{e_i^t}=k_{e^t}} p(d_m, e_i^t); \quad (10)$$

$$q(d_m|\hat{d}_{k_d}) = \frac{q(d_m)}{q(l_{d_m}=k_d)} \quad [q(d_m|\hat{d}_{k_d}) = 0 \text{ if } l_{d_m} \neq k_d]; \quad (11)$$

$$q(e_i^t|\hat{e}_{k_{e^t}}^t) = \frac{q(e_i^t)}{q(l_{e_i^t}=k_{e^t})} \quad [q(e_i^t|\hat{e}_{k_{e^t}}^t) = 0 \text{ if } l_{e_i^t} \neq k_{e^t}]; \quad (12)$$

where $q(d_m) = \sum_{e_i^t} p(d_m, e_i^t)$, $q(e_i^t) = \sum_{d_m} p(d_m, e_i^t)$, $q(\hat{d}_{k_d}) = \sum_{k_{e^t}} p(\hat{d}_{k_d}, \hat{e}_{k_{e^t}}^t)$ and $q(\hat{e}_{k_{e^t}}^t) = \sum_{k_d} p(\hat{d}_{k_d}, \hat{e}_{k_{e^t}}^t)$.

Algorithm 1 summarizes the main steps in the procedure. The objective function (6) with our alternating update monotonically decreases to a local optimum. This is because the ICM algorithm decreases the non-negative objective function (6) to a local optimum given a fixed q function. Then the update of q is monotonically decreasing as guaranteed by the theorem proven in [37].

The time complexity of Algorithm 1 is $O(n_{\mathcal{D}, \mathcal{W}} \cdot (K_d + K_w) + \sum_{t=1}^T n_{\mathcal{D}, \mathcal{E}^t} \cdot (K_d + K_{e^t}) + \sum_{t=1}^T \sum_{s=1}^T (n_{\mathcal{E}^t, \mathcal{E}^s} + (n_c * \text{iter}_{ICM}))) \cdot (K_{e^t} + K_{e^s})) \cdot \text{iter}_{AO}$, where $n_{\cdot, \cdot}$ is the total number of non-zero elements in the corresponding co-occurrence matrix, n_c is the number of constraints, iter_{ICM} is the number of ICM iterations, K_d, K_w and K_{e^t} are the number of document clusters, word clusters and entity clusters of type t , and iter_{AO} is the number of the alternating optimization iterations.

4. EXPERIMENTS

In this section, we show the experimental results to demonstrate the effectiveness and efficiency of our approach on document clustering with world knowledge as indirect supervision.

4.1 Datasets

We use the following two benchmark datasets to evaluate domain dependent document clustering. For both datasets we assume the numbers of document clusters are given.

20Newsgroups (20NG): The 20newsgroups dataset [24] contains about 20,000 newsgroups documents evenly distributed across 20 newsgroups.² We use all the 20 groups as 20 classes.

RCV1: The RCV1 dataset is a dataset containing manually labeled newswire stories from Reuter Ltd [26]. The news documents are categorized with respect to three controlled vocabularies: industries, topics and regions. There are 103 categories including all nodes except for root in the hierarchy. The maximum depth is four, and 82 nodes are leaves. We select top categories MCAT, CCAT and ECAT in one portion of the test partition to form three clustering tasks. The three clustering tasks are summarized in Table 2. We use the original source of this data, and use the leaf categories in each task as the ground-truth classes.

Table 2: RCV1 dataset statistics. #(Categories) is the number of all categories; #(Leaf Categories) is the number of leaf categories; #(Documents) is the number of documents.

	#(Categories)	#(Leaf Categories)	#(Documents)
MCAT	9	7	44,033
CCAT	31	26	47,494
ECAT	23	18	19,813

4.2 World Knowledge Bases

Then we introduce the knowledge bases we use.

Freebase: Freebase³ is a publicly available knowledge base consisting of entities and relations collaboratively collected by its community members. Now, it contains over 2 billions relation expressions between 40 millions entities. We convert a logical form generated by our unsupervised semantic parser of the world knowledge specification approach introduced in Section 2.1 into a SPARQL query and execute it on our copy of Freebase using the Virtuoso engine.

YAGO2: YAGO2⁴ is also a semantic knowledge base, derived from Wikipedia, WordNet and GeoNames. Currently, YAGO2 has knowledge of more than 10 million entities (like persons, organizations, cities, etc.) and contains more than 120 million facts about these entities. Similar to Freebase, we also convert a logical form into a SPARQL query and execute it on our copy of YAGO2 using the Virtuoso engine.

In Table 3, we show some statistics about Freebase and YAGO2.

Note that in most knowledge bases, such as Freebase and YAGO2, entities types are often organized in a hierarchical manner. For example, *Politician* is a sub-type of *Person*. *University* is a sub-type of *Organization*. All the types or attributes share a common root, called *Object*. Figure 3 depicts an example of hierarchy of types. In general, we use the highest level under the root object as the entity types (e.g., *Person*) as specified world knowledge incorporated in the HIN, and the direct children (e.g., *Politician*) as entity constraints that will be introduced later. In the following experiments, we select *Person*, *Organization*, and *Location* as the three entity types in the HIN, because they are popular in both Freebase and YAGO2.

²<http://qwone.com/~jason/20Newsgroups/>

³<https://developers.google.com/freebase/>

⁴<http://www.mpi-inf.mpg.de/departments/databases-and-information-systems/research/yago-naga/yago/>

Table 3: Statistics of Freebase and YAGO2. #(Entity Types) is the number of entity types; #(Entity Instances) is the number of entity instances; #(Relation Types) is the number of relation types; #(Relation Instances) is the number of relation instances.

Name	Freebase	YAGO2
#(Entity Types)	1,500	350,000
#(Entity Instances)	40 millions	10 millions
#(Relation Types)	35,000	100
#(Relation Instances)	2 billions	120 millions

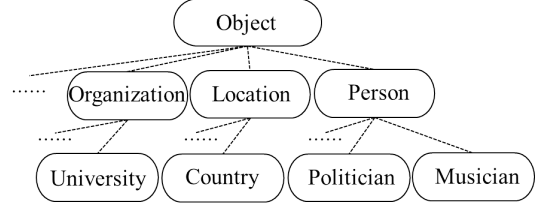


Figure 3: Hierarchy of entity types.

4.3 Effectiveness of World Knowledge Specification

Before applying the specified world knowledge to downstream text analytics tasks, such as document clustering in our case, we need to evaluate whether our world knowledge specification approach could produce the correct specified world knowledge.

In order to test the effectiveness of our world knowledge specification approach, we first sample 200 documents from 20 newsgroups, i.e., 10 documents from each category. Second, we split the documents into sentences using Stanford CoreNLP. After post-processing, 3,232 sentences are generated for human evaluation. Third, we use our world knowledge specification approach in Section 2.1 with three different semantic filtering modules to generate the specified world knowledge for each sentence, which consists of relation triplets in the form of (e_1, r, e_2) with the type information. Afterwards, we ask three annotators to label the specified world knowledge according two criterion: (1) whether the boundaries of e_1 and e_2 are correctly recognized or not; (2) whether the entity type of e_1 and e_2 are correct or not. The label equals to 1 if both (1) and (2) are satisfied. Otherwise, the label equals to 0. We check the mutual agreement of the human annotation, which is around 91.3% accuracy.

Table 4: Precision of different semantic filtering results. FBSF represents frequency based semantic filter; DFBSF represents document frequency based semantic filter; CBSF represents conceptualization based semantic filter.

Semantic Filter	FBSF	DFBSF	CBSF
Precision	0.751	0.890	0.916

We then test the precision of three different specified world knowledge generated by the corresponding semantic filtering method. The results are shown in Table 4. From the results we can see that, CBSF outperforms the other two ways to generate the correct semantic meaning. The main reason is that, conceptualization based method is able to use the context information to help judge the real semantic of the text rather than only taking the statistics of the data into account. Here we only care about precision because we wish to use world knowledge as indirect supervision. The recall will not

Table 5: Error analysis of specified world knowledge generated by the world knowledge specification approach with three different semantic filters. FBSF represents frequency based semantic filter; DFBSF represents document frequency based semantic filter; CBSF represents conceptualization based semantic filter.

Type of error	Example sentence	Number and percentage of errors		
		FBSF (805)	DFBSF (359)	CBSF (272)
Entity Recognition	“Einstein ’s theory of relativity explained mercury ’s motion.”	179 (22.2%)	129 (35.9%)	105 (38.6%)
Entity Disambiguation	“Bill said all this to make the point that Christianity is eminently.”	537 (66.7%)	182 (50.7%)	130 (47.8%)
Subordinate Clause	“Bruce S. Winters, worked at United States Technologies Research Center, bought a Ford.”	89 (11.1%)	48 (13.4%)	37 (13.6%)

Table 6: Performance of different clustering algorithms on 20NG and RCV1 data. CHINC is our proposed method. BOW, FB (Freebase), or YG (YAGO2) represent bag of word features, the entities generated by our world knowledge specification approach based on Freebase or YAGO2, respectively. We compared all the numbers of HINC and CHINC with CITCC, which is the strongest baseline. CITCC uses 250K constraints generated based on ground-truth labels of documents.

Features Data	Kmeans			ITCC			CITCC	HINC		CHINC	
	BOW	BOW +FB	BOW +YG	BOW	BOW +FB	BOW +YG		FB	YG	FB	YG
20NG	0.429	0.447	0.437	0.501	0.525	0.513	0.569	0.571 (+0.4%)	0.541 (−4.9%)	0.631 (+10.9%)	0.600 (+5.5%)
MCAT	0.549	0.575	0.559	0.604	0.630	0.619	0.652	0.645 (−1.1%)	0.625 (−4.1%)	0.698 (+7.1%)	0.685 (+5.1%)
CCAT	0.403	0.419	0.410	0.472	0.494	0.481	0.535	0.542 (+1.3%)	0.515 (−3.7%)	0.606 (+13.3%)	0.574 (+7.3%)
ECAT	0.417	0.436	0.424	0.493	0.516	0.505	0.562	0.561 (−0.2%)	0.530 (−5.7%)	0.624 (+11.0%)	0.588 (+4.6%)

be very important.

Error Analysis

To further investigate what triggers the errors in our semantic parsing and semantic filtering pipelines, we analyze the cause of errors for the incorrect specified world knowledge. As shown in Table 5, we categorize the errors as follows:

Entity Recognition: In semantic parsing, entities can be extracted incorrectly. Long entities are composed of multiple simple entities. For example, “Einstein ’s theory of relativity” may be extracted as “Einstein” and “theory of relativity.” Paraphrasing and misspelling entities cause their textual expressions to deviate from any knowledge base entries. Idiomatic expressions are incorrectly picked up as entities. Using a larger mapping from text to knowledge base phrases, or paraphrasing techniques will help avoid some errors. However, this is out of the scope of this paper.

Entity Disambiguation: Selecting an incorrect entity out of multiple matching candidates causes this error, e.g., “Bill” in our example sentence can be “Bill Clinton” or “Bill Gates.” Primarily due to two reasons: first, entity disambiguation is a tough research problem in NLP community. Second, the type information of relations are not sufficient to further prune out mismatching entities during semantic filtering process. Notice that, entity disambiguation is the major cause of the errors. By using CBSF, the number of incorrect entities caused by disambiguation can be dramatically reduced.

Subordinate Clause: Semantic parsing sometimes produces wrong relation phrases in the subordinate clauses. For example, in the example wrong case shown in Table 5, it takes the relation phrase “worked at” meaning the working place of “Bruce S. Winters,” ignores the phrase “bought,” which could be more informative for the target clustering domain. This could be resolved by adding more concrete rules in the semantic parsing grammar.

In the following experiments, we use the world knowledge specification approach with CBSF, because it performs the best among the three semantic filtering methods.

4.4 Clustering Result

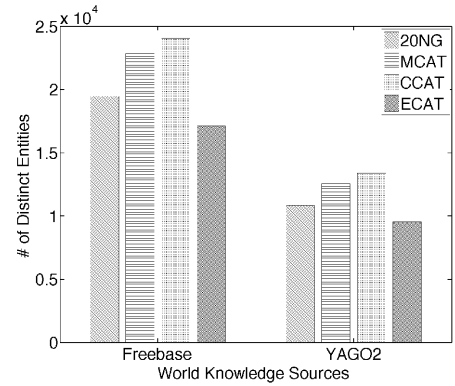


Figure 4: Statistics of the number of entities in different document datasets with different world knowledge sources.

In this experiment, we compare the performance of our model, constrained heterogeneous information network clustering (CHINC), with several representative clustering algorithms such as Kmeans, ITCC [10] and CITCC [37]. The parameters used in CHINC to control the constraints are w_M and w_C . We set them following the rules tested in [37]. We also denote our algorithm without constraints as HINC. “FB” and “YG” represent two different world knowledge sources, Freebase and YAGO2, respectively. We re-implement all the above clustering algorithms. Notice that, for CITCC, we follow [37] to generate and add constraints for documents and words. We also use the specified world knowledge as features to enhance the Kmeans and ITCC. The feature settings are defined as below:

- BOW: Traditional bag-of-words model with the tf-idf weighting mechanism.
- BOW+FB: BOW integrated with additional features from entities in specified world knowledge of Freebase.
- BOW+YG: BOW integrated with additional features from

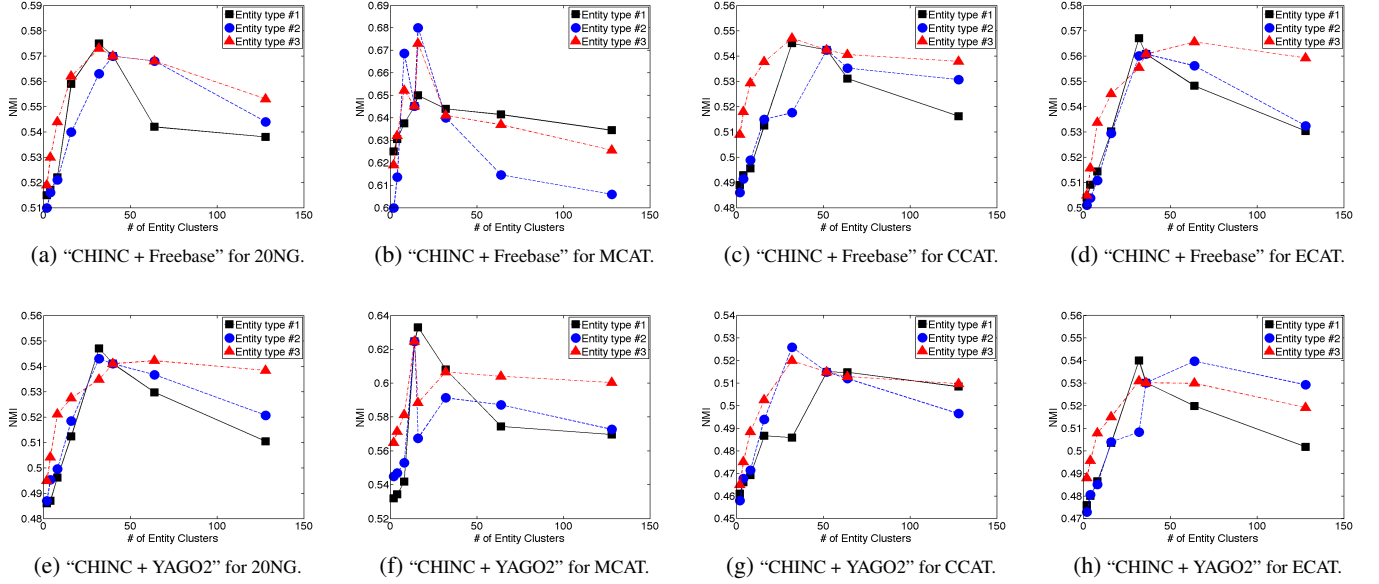


Figure 5: Effect of number of entity clusters of each entity type on document clustering on different dataset and world knowledge source combinations.

entities in specified world knowledge of YAGO2.

We employ the widely-used normalized mutual information (NMI) [39] as the evaluation measure. The NMI score is 1 if the clustering results match the category labels perfectly and 0 if the clusters are obtained from a random partition. In general, the larger the scores are, the better the clustering results are.

In Table 6, we show the performance of all the clustering algorithms with different experimental settings. The NMI is the average NMI of five random trials per experiment setting. Overall, among all the methods we test, CHINC consistently performs the best among all the clustering methods. We can see that HINC+FB and HINC+YG perform better than ITCC with BOW+FB or BOW+YG features, respectively. This means that by using the structural information provided by the world knowledge, we can further improve the clustering results. In addition, the algorithms with Freebase consistently outperform the ones with YAGO2, since Freebase has much more facts compared with YAGO2 as shown in Table 3; besides, one can see in Figure 4 that Freebase could consistently specify more entities than YAGO2 does from all of the document datasets. CITCC is the strongest baseline clustering algorithm, because it uses the ground-truth constraints derived from category labels based on the human knowledge. We use 250K constraints to perform CITCC. As shown in Table 6, HINC performs competitive with the CITCC. CHINC significantly outperforms CITCC. This shows that by automatically using world knowledge, it has the potential to perform better than the algorithm with the specific domain knowledge.

4.4.1 Analysis of Number of Entity Clusters

We also evaluate the effect of varying the number of entity clusters of each entity type in CHINC on the document clustering task. Figure 5a shows the results of clustering with different numbers of entity clusters of each entity type on “CHINC + Freebase” for the 20NG dataset. The number of entity clusters varies from 2 to 128. The default number of iterations is set as 20, which will be discussed in Section 4.4.2. When testing the effect of the number of entity clusters of one entity type, the numbers of entity clusters

of the other two entity types are fixed as twice as the number of document clusters, which are 40 and 40 in 20NG, respectively. It is shown that for this dataset, more entity clusters may not result in improved document clustering results when a sufficient number of entity clusters is reached. For example, as shown in Figure 5a, after reaching 32, the NMI scores of CHINC actually decrease when the numbers of entity clusters further increase.

One can also find the effects of the numbers of entity clusters on the clustering performance with the other document dataset and knowledge base combinations in Figure 5b—5h. From the results, we can conclude that, there exist certain values of the number of entity clusters leading to the best clustering performance. Similar to the results on “CHINC + Freebase” for 20NG dataset, in the rest of the experiments, we fix the number of entity clusters of each entity type to be twice the number of document clusters.

4.4.2 Analysis of Number of Iterations in Alternating Optimization

We herein evaluate the impact of the number of iterations of the alternating optimization (Algorithm 1) on CHINC in relation to the execution time of the optimization algorithm as well as the clustering performance. We increase the number of iterations from 1 to 80. For example, for each number of iterations, we run CHINC five trials, and the average execution time and NMI are summarized in Figure 6a. From the result, one can conclude that the larger number of iterations is, the improvement on clustering performance is more significant, but the improvement will drop and become stable. The reason is that, along with the increase of the number of iterations, the alternating optimization algorithm comes to convergence. However, the execution time still increase in a nearly linear manner. For example, as shown in Figure 6a, after reaching 20, the performance stays stable. Thus, we set the number of iterations as 20 in the remaining experiments with the consideration of both performance and efficiency. As shown in Figure 6b—6h, we set the number of iterations as 20 when conducting experiments on the other combinations of document datasets and world knowledge bases.

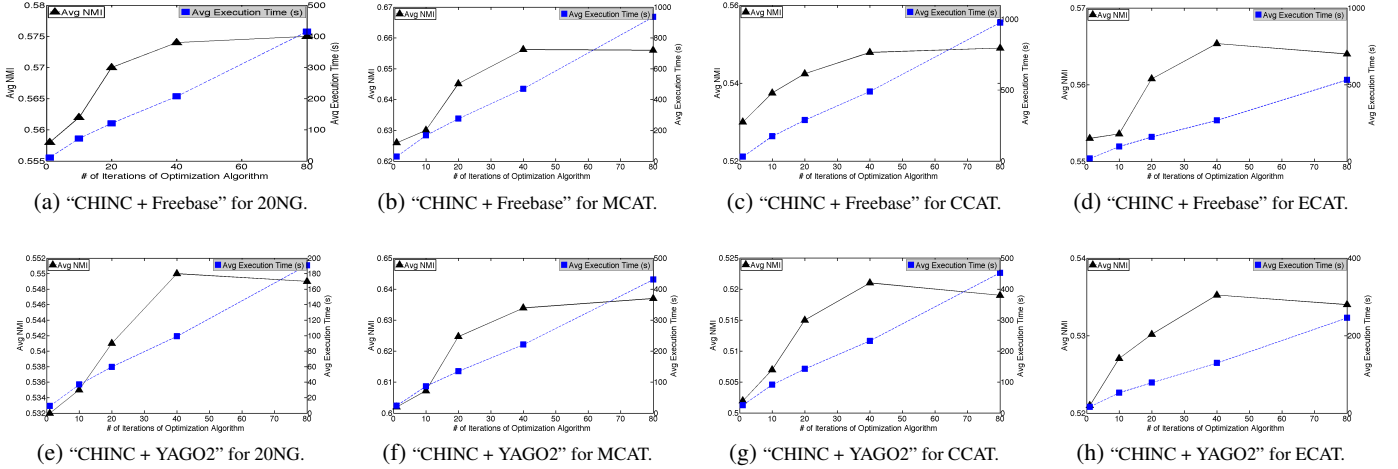


Figure 6: Analysis of # of iterations in alternating optimization algorithm on different dataset and world knowledge source combinations. Left y -axis: average NMI; Right y -axis: average execution time (s).

4.4.3 Analysis of Specified World Knowledge based Constraints

Rather than using human knowledge as constraints, we instead use the specified world knowledge automatically generated by our approach as constraints in CHINC. Based on the specified world knowledge, it is straightforward to design constraints for entities.

Entity constraints. (1) **Must-links.** If two entities belong to the same entity sub-type, we add a must-link. (2) **Cannot-links.** If two entities belong to different entity sub-types, we add a cannot-link. For example, the entity sub-types of “Obama” and “United States” are *Politician* and *Country* respectively. In this case, we add a cannot-link to them.

We then test the performance of our proposed CHINC by using the specified world knowledge as constraints described above. We show the experiments on all of the different combinations of datasets and world knowledge sources in Figure 7. Each x -axis represents the number of entity type constraints used in each experiment, and y -axis is the average NMI of five random trials. For example, the constraints derived from entity type #1, #2, and #3 are eventually added to CHINC as shown in Figure 7a, Figure 7b and Figure 7c respectively, when using Freebase as world knowledge and testing on 20NG dataset. We can see that CHINC outperforms the best clustering algorithm with the human knowledge as shown in Table 6 (CITCC: 0.569) with even no constraints (HINC: 0.571). By adding more and more constraints, the clustering result of CHINC is significantly better. So CHINC is able to use information in world knowledge specified in the HIN, and the entity sub-type information can be transferred to the document side. The results show the power of modeling data as heterogeneous information networks, as well as the high quality of constraints derived from world knowledge.

Moreover, we also test the scalability of CHINC algorithm by adding more constraints into the algorithm. The results are also shown in Figure 7. By increasing the number of constraints, we find that the average execution time of five trials increases linearly, and the clustering performance measured by NMI is increasing as mentioned before. For example, Figure 7c shows the effects of the constraints of all the three entity types on the clustering performance as well as the execution time, when Freebase is used as world knowledge and CHINC is tested on 20NG dataset. After the number of

constraints reach 50M, the increase of performance drops and stays stable. At this point, the execution time is around 1.2M (s). In Figure 8, one can see the similar results on the other combinations of document datasets and knowledge bases. We also find that the average execution time of our algorithm with Freebase as world knowledge source is greater than that with YAGO2. As shown in Figure 4, the reason is that each document datasets with Freebase could be specified much more entities than that with YAGO2. From the results, we can see that our algorithm is scalable to use the large scale specified world knowledge as constraints, and cluster large amounts of documents.

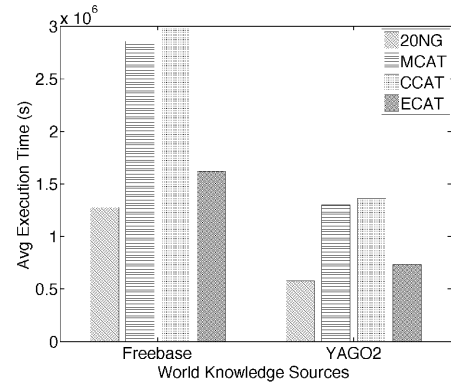


Figure 8: Analysis of the efficiency of our algorithm on different document datasets with different world knowledge sources.

5. RELATED WORK

In this section, we review the related work on document clustering, machine learning with world knowledge, and heterogeneous information network.

5.1 Document Clustering

Document clustering has been studied for many years. We can use traditional one-dimensional clustering algorithms (e.g., Kmeans) to cluster the documents. If we treat the document and corresponding words as a bipartite graph, we can use co-clustering algorithms [10]

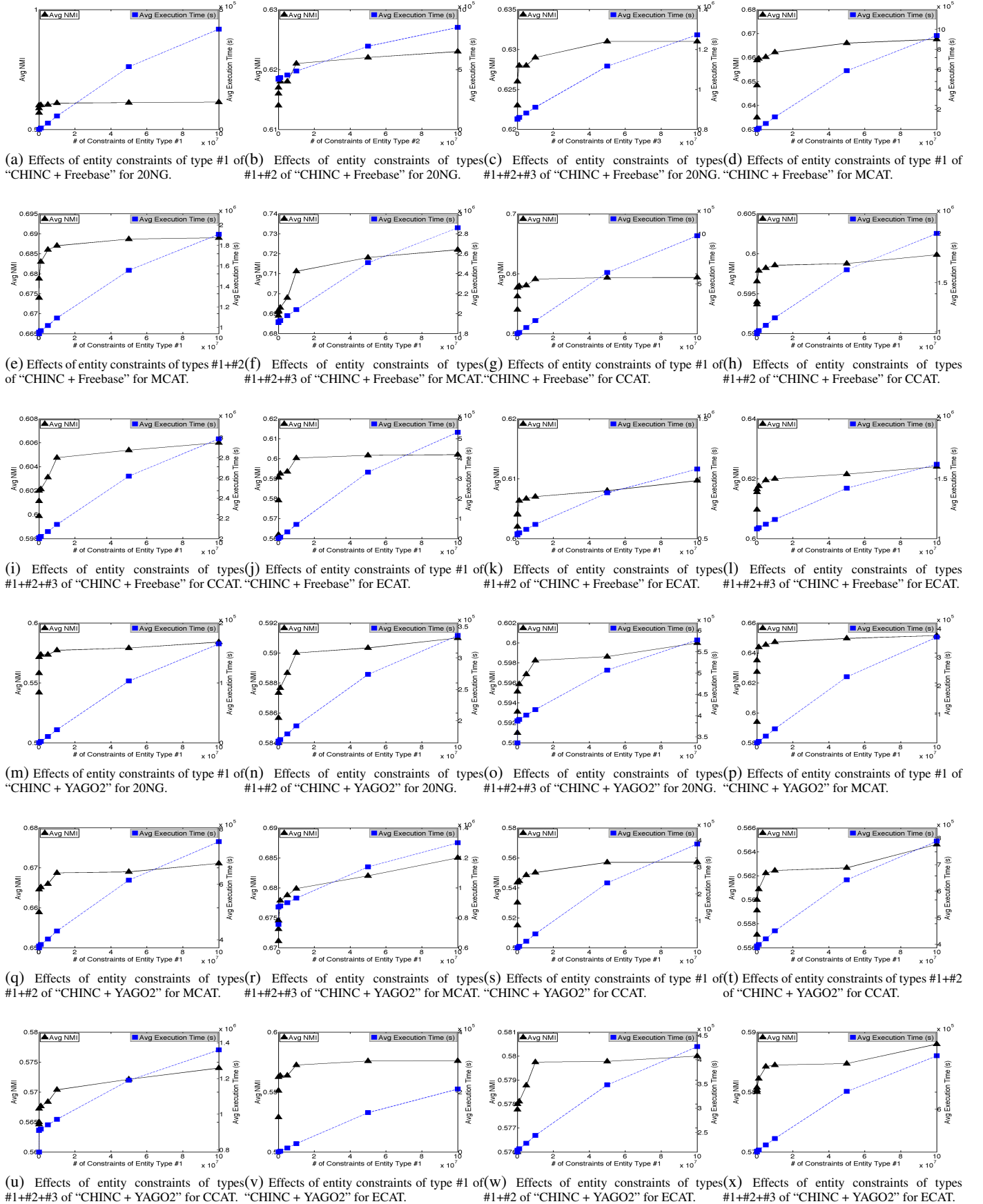


Figure 7: Effects of entity constraints on different dataset and world knowledge source combinations. Left y -axis: average NMI; Right y -axis: average execution time (s).

to cluster the documents. Moreover, with the help of labeled seed documents, semi-supervised clustering can be used [3]. When the seeds are not available, we can use side information as constraints to guide clustering algorithms [4]. When the supervision from target domain is not available, we can also perform transfer learning to transfer the labeled information from other domains to the target domain [9, 43]. All the above clustering algorithms with supervision need domain or relevant domain knowledge. When there are diverse domains and the supervision is needed, they will still be very costly to ask a lot of different domain experts to label.

5.2 Machine Learning with World Knowledge

Most of the existing usage of world knowledge is to enrich the features beyond bag-of-words representation of documents. For example, by using the linguistic knowledge base WordNet to resolve synonyms and introduce WordNet concepts, the quality of document clustering can be improved [18]. The first paper using the term “world knowledge” [14] extends the bag-of-words features with the categories in Open Directory Project (ODP), and shows that it can help improve text classification with additional knowledge. Following this, by mapping the text to the semantic space provided by Wikipedia pages, it has been proven to be useful for short text classification [15, 16] and clustering [19, 20, 21]. Liu et al. [29] also use another knowledge base of taxonomy, Probase, to enrich the features of ads keywords to build a new taxonomy of domain dependent keyword set. All of the above approaches just consider to use world knowledge as a source of features. However, the knowledge in the knowledge bases indeed has annotations of types, categories, etc.. Thus, it can be more effective to consider this information as “supervision” to supervise other machine learning algorithms and tasks.

Distant supervision uses the knowledge of entities and their relationships from world knowledge bases, e.g., Freebase, as supervision for the task of entity and relation extraction [32]. It only considers to extract more entities and relations from new text. Thus, the application of direct supervision is limited.

Song et al. [37] consider using fully unsupervised method to generate constraints of words using an external general-purpose knowledge base, WordNet. This can be regarded as an initial attempt to use general knowledge as indirect supervision to help clustering. However, the knowledge from WordNet is mostly linguistically related. It lacks of the information about named entities and their types. Moreover, their approach is still a simple application of constrained co-clustering, where it misses the rich structural information in the knowledge base.

5.3 Heterogeneous Information Network

A heterogeneous information network (HIN) is defined as a graph of multi-typed entities and relations [17]. Different from traditional graphs, HIN incorporates the type information which can be useful to identify the semantic meaning of the paths in the graph [41]. This is a good property to perform graph search and matching. Original HINs are developed for the applications of scientific publication network analysis [41, 42]. Then social network analysis also leverages this representation for user similarity and link prediction [22, 44, 45]. Seamlessly, we can see that the knowledge in world knowledge bases, e.g., Freebase and YAGO2, can be naturally represented as an HIN, since the entities and relations in the knowledge base are all typed. We introduce this representation to knowledge based analysis, and show that it can be very useful for our document clustering task. Note that there is also a series of methods called multi-type relational data clustering [30, 31]. While they require the data to be structural beforehand (e.g., providing in-

formation of authors, co-authors, etc.), our method only needs the input of raw documents. In addition to the multi-type relational information, we also incorporate the type information provided by the knowledge base as constraints to further improve the clustering results.

6. CONCLUSION

In this paper, we study a novel problem of machine learning with world knowledge. Particularly, we take document clustering as an example and show how to use world knowledge as indirect supervision to improve the clustering results. To use the world knowledge, we show how to adapt the world knowledge to domain dependent tasks by using semantic parsing and semantic filtering. Then we represent the data as a heterogeneous information network, and use a constrained network clustering algorithm to obtain the document clusters. We demonstrate the effectiveness and efficiency of our approach on two real datasets along with two popular knowledge bases. In the future, we plan to use world knowledge to help more text mining and text analytics tasks, such as text classification and information retrieval.

7. REFERENCES

- [1] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, and Z. Ives. *Dbpedia: A nucleus for a web of open data*. Springer, 2007.
- [2] M. Banko, M. J. Cafarella, S. Soderland, M. Broadhead, and O. Etzioni. Open information extraction from the web. In *IJCAI*, pages 2670–2676, 2007.
- [3] S. Basu, A. Banerjee, and R. J. Mooney. Semi-supervised clustering by seeding. In *ICML*, pages 27–34, 2002.
- [4] S. Basu, M. Bilenko, and R. J. Mooney. A probabilistic framework for semi-supervised clustering. In *KDD*, pages 59–68, 2004.
- [5] J. Berant, A. Chou, R. Frostig, and P. Liang. Semantic parsing on freebase from question-answer pairs. In *EMNLP*, pages 1533–1544, 2013.
- [6] K. D. Bollacker, C. Evans, P. Paritosh, T. Sturge, and J. Taylor. Freebase: a collaboratively created graph database for structuring human knowledge. In *SIGMOD*, pages 1247–1250, 2008.
- [7] A. Carlson, J. Betteridge, B. Kisiel, B. Settles, E. R. Hruschka Jr, and T. M. Mitchell. Toward an architecture for never-ending language learning. In *AAAI*, volume 5, page 3, 2010.
- [8] O. Chapelle, B. Schölkopf, and A. Zien, editors. *Semi-Supervised Learning*. MIT Press, 2006.
- [9] W. Dai, G.-R. Xue, Q. Yang, and Y. Yu. Co-clustering based classification for out-of-domain documents. In *KDD*, pages 210–219, 2007.
- [10] I. S. Dhillon, S. Mallela, and D. S. Modha. Information-theoretic co-clustering. In *KDD*, pages 89–98, 2003.
- [11] X. Dong, E. Gabrilovich, G. Heitz, W. Horn, N. Lao, K. Murphy, T. Strohmann, S. Sun, and W. Zhang. Knowledge vault: A web-scale approach to probabilistic knowledge fusion. In *KDD*, pages 601–610, 2014.
- [12] O. Etzioni, M. Cafarella, and D. Downey. Webscale information extraction in knowitall (preliminary results). In *WWW*, pages 100–110, 2004.
- [13] C. Fellbaum, editor. *WordNet: an electronic lexical database*. MIT Press, 1998.

- [14] E. Gabrilovich and S. Markovitch. Feature generation for text categorization using world knowledge. In *IJCAI*, pages 1048–1053, 2005.
- [15] E. Gabrilovich and S. Markovitch. Overcoming the brittleness bottleneck using wikipedia: Enhancing text categorization with encyclopedic knowledge. In *AAAI*, pages 1301–1306, 2006.
- [16] E. Gabrilovich and S. Markovitch. Computing semantic relatedness using wikipedia-based explicit semantic analysis. In *IJCAI*, pages 1606–1611, 2007.
- [17] J. Han, Y. Sun, X. Yan, and P. S. Yu. Mining knowledge from databases: An information network analysis approach. In *SIGMOD*, pages 1251–1252, 2010.
- [18] A. Hotho, S. Staab, and G. Stumme. Ontologies improve text document clustering. In *ICDM*, pages 541–544, 2003.
- [19] J. Hu, L. Fang, Y. Cao, H.-J. Zeng, H. Li, Q. Yang, and Z. Chen. Enhancing text clustering by leveraging Wikipedia semantics. In *SIGIR*, pages 179–186, 2008.
- [20] X. Hu, N. Sun, C. Zhang, and T.-S. Chua. Exploiting internal and external semantics for the clustering of short texts using world knowledge. In *CIKM*, pages 919–928, 2009.
- [21] X. Hu, X. Zhang, C. Lu, E. K. Park, and X. Zhou. Exploiting wikipedia as external knowledge for document clustering. In *KDD*, pages 389–396, 2009.
- [22] X. Kong, J. Zhang, and P. S. Yu. Inferring anchor links across multiple heterogeneous social networks. In *CIKM*, pages 179–188, 2013.
- [23] T. Kwiatkowski, L. S. Zettlemoyer, S. Goldwater, and M. Steedman. Lexical generalization in CCG grammar induction for semantic parsing. In *EMNLP*, pages 1512–1523, 2011.
- [24] K. Lang. Newswelder: Learning to filter netnews. In *ICML*, pages 331–339, 1995.
- [25] D. B. Lenat and R. V. Guha. *Building Large Knowledge-Based Systems: Representation and Inference in the Cyc Project*. Addison-Wesley, 1989.
- [26] D. D. Lewis, Y. Yang, T. G. Rose, and F. Li. Rcv1: A new benchmark collection for text categorization research. *JMLR*, 5:361–397, 2004.
- [27] Y. Li, C. Wang, F. Han, J. Han, D. Roth, and X. Yan. Mining evidences for named entity disambiguation. In *KDD*, pages 1070–1078, 2013.
- [28] P. Liang. Lambda dependency-based compositional semantics. *arXiv*, 2013.
- [29] X. Liu, Y. Song, S. Liu, and H. Wang. Automatic taxonomy construction from keywords. In *KDD*, pages 1433–1441, 2012.
- [30] B. Long, Z. M. Zhang, X. Wú, and P. S. Yu. Spectral clustering for multi-type relational data. In *ICML*, pages 585–592, 2006.
- [31] B. Long, Z. M. Zhang, and P. S. Yu. A probabilistic framework for relational clustering. In *KDD*, pages 470–479, 2007.
- [32] M. Mintz, S. Bills, R. Snow, and D. Jurafsky. Distant supervision for relation extraction without labeled data. In *ACL/AFNLP*, pages 1003–1011, 2009.
- [33] R. J. Mooney. Learning for semantic parsing. In *CICLing*, pages 311–324, 2007.
- [34] S. J. Pan and Q. Yang. A survey on transfer learning. *IEEE TKDE*, 22(10):1345–1359, 2010.
- [35] S. P. Ponzetto and M. Strube. Deriving a large-scale taxonomy from wikipedia. In *AAAI*, pages 1440–1445, 2007.
- [36] L. Ratinov and D. Roth. Design challenges and misconceptions in named entity recognition. In *CoNLL*, pages 147–155, 2009.
- [37] Y. Song, S. Pan, S. Liu, F. Wei, M. Zhou, and W. Qian. Constrained text coclustering with supervised and unsupervised constraints. *IEEE TKDE*, 25(6):1227–1239, 2013.
- [38] Y. Song, H. Wang, Z. Wang, H. Li, and W. Chen. Short text conceptualization using a probabilistic knowledgebase. In *IJCAI*, pages 2330–2336, 2011.
- [39] A. Strehl and J. Ghosh. Cluster ensembles—a knowledge reuse framework for combining multiple partitions. *JMLR*, 3:583–617, 2003.
- [40] F. M. Suchanek, G. Kasneci, and G. Weikum. Yago: a core of semantic knowledge. In *WWW*, pages 697–706, 2007.
- [41] Y. Sun, J. Han, X. Yan, P. S. Yu, and T. Wu. Pathsims: Meta path-based top-k similarity search in heterogeneous information networks. *PVLDB*, pages 992–1003, 2011.
- [42] Y. Sun, B. Norick, J. Han, X. Yan, P. S. Yu, and X. Yu. Integrating meta-path selection with user-guided object clustering in heterogeneous information networks. In *KDD*, pages 1348–1356, 2012.
- [43] Z. Wang, Y. Song, and C. Zhang. Knowledge transfer on hybrid graph. In *IJCAI*, pages 1291–1296, 2009.
- [44] J. Zhang, X. Kong, and P. S. Yu. Predicting social links for new users across aligned heterogeneous social networks. In *ICDM*, pages 1289–1294, 2013.
- [45] J. Zhang, X. Kong, and P. S. Yu. Transferring heterogeneous links across location-based social networks. In *WSDM*, pages 303–312, 2014.