

A Unified Framework for Document Clustering with Dual Supervision

Yeming Hu
Dalhousie University
Faculty of Computer Science
6050 University Avenue
Halifax, Canada
yeming@cs.dal.ca

Evangelos E. Milios
Dalhousie University
Faculty of Computer Science
6050 University Avenue
Halifax, Canada
eem@cs.dal.ca

James Blustein
Dalhousie University
Faculty of Computer Science
and School of Information
Management
jamie@cs.dal.ca

ABSTRACT

Semi-supervised clustering algorithms for general problems use a small amount of labeled instances or pairwise instance constraints to aid the unsupervised clustering. However, user supervision can also be provided in alternative forms for document clustering, such as labeling a feature by associating it with a document or a cluster. Besides labeled documents, this paper also explores labeled features to generate cluster seeds to seed the unsupervised clustering. In this paper, we present a unified framework in which one can use both labeled documents and features in terms of seeding clusters and refine this information using intermediate clusters. We introduce two methods of using labeled features to generate cluster seeds. Experimental results on several real-world data sets demonstrate that constraining the clustering by both documents and features seeding can significantly improve document clustering performance over random seeding and document only seeding. We also demonstrate that the clustering performance can be improved even with only a fraction of clusters being seeded compared to unsupervised clustering.¹

Categories and Subject Descriptors

H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*Clustering*; I.5.4 [Pattern Recognition]: Application—*Text Processing*

General Terms

Algorithm, Document Clustering, Features

Keywords

User Supervision, Feature Supervision, Seeding, Text Cloud

1. INTRODUCTION

Traditional document clustering is an unsupervised categorization that partitions a given document collection into clusters so that topically similar documents are placed into the same clusters. However, given the same document collection, different users may want to organize it in their own point of view instead of a universal one, which is addressed to some extent by incorporating document supervision [3]. In this paper, we have two types of user supervision, namely,

¹This work is based on an earlier work: SAC '12 Proceedings of the 2012 ACM Symposium on Applied Computing, Copyright 2012 ACM 978-1-4503-0857-1/12/03. <http://doi.acm.org/10.1145/2245276.2245306>

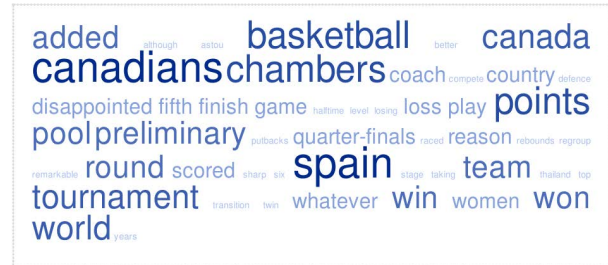


Figure 1: Text Cloud of a Document about Canadian Basketball

document supervision and *feature supervision* for document clustering. *Document Supervision* involves labeling documents, i.e., assigning a document to a cluster. *Feature Supervision* involves labeling features, i.e., associating a feature with a document if that feature describes the topic of that document.

Most prior semi-supervised clustering algorithms use user supervision in the form of document supervision such as labeled instances [3] or instance pairwise constraints [26] for general clustering problems. However, user supervision can also be provided in alternative forms such as labeling features (words) for document clustering in addition to labeling instances (documents). Since this paper focuses on document clustering, we may use *instance* and *document*, *feature* and *word* interchangeably. Labeling documents and words can be performed at the same time, with *little additional effort* for labeling words, if an appropriate document visualization is used, such as text clouds [19]. While the user assigns a document to a cluster based on the document's text cloud, the words appearing in the text cloud can also be labeled by being clicked or highlighted.

Example 1. Consider a collection of news articles about international sports. While the user labels the document displayed as text cloud (Fig. 1) to a cluster, the words associating the document with the specific cluster can also be labeled by being clicked or highlighted. In one scenario, the document (Fig. 1) can be labeled to cluster “Canada”, in which the words “Canada”, “Canadians” should be labeled (associated) with the document. In another scenario, the document would be labeled to cluster “Basketball”, in which the words “basketball”, “points” should be associated with the document.

Example 2. Assume we have two papers and one talks about programming languages while the other is about software debugging. One human user can assign them into the same cluster “software engineering” while another one would like to put them into two clusters, i.e., “languages” and “debugging”. Clearly, the keywords (features) assigned for the two cases will be different too.

Therefore, different labeled words reflect different organizations and the user forms his point of view based on the perception of the words in the text clouds. By using the text cloud for labeling documents, the user can not only label documents to seed the clustering but also label the words discriminating among clusters. It has been argued that document supervision and feature supervision are complementary rather than completely redundant and this joint use has been called *dual supervision* [1].

In this paper, we assume that the user labels a document by reading its content. At the same time, the user can label a word by indicating (e.g. highlighting) whether it is associated with the document or the specific cluster. The text cloud could be used to visualize the document content and enhance the labeling. We extend two methods incorporating the labeled features from document classification to document clustering, namely, feature-vote-model [9] which uses labeled features to vote for cluster label of an unlabeled documents, and feature-generative-model [21] which uses labeled features to infer a multinomial generative model. In (semi-supervised) document classification, labeled documents and features are required for each category. However, knowledge of the relevant categories is incomplete in many domains. Semi-supervised document clustering can group documents into partial clusters with labeled documents and features, as well as extend and modify the existing set of clusters to reflect other topical groupings in document collection [3]. In this paper, we propose a clustering model built from both the labeled documents and the labeled features can be used to guide the clustering process. At the same time, the knowledge from the labeled documents and features will be refined by intermediate clusters in an iterative manner. To this end, we present a unified framework which combines knowledge from labeled documents, labeled features, and unlabeled documents by an iterative clustering process. Finally, we demonstrate the effectiveness of the framework on several real-word data sets.

The rest of this paper is organized as follows. Related work on semi-supervised clustering and feature supervision is discussed in Section 2. In Section 3, we introduce the models to incorporate the labeled features and present the unified framework to combine knowledge from labeled documents, labeled features and intermediate clusters. The details of the experimental results on several real-world text datasets are presented and discussed in Section 4. We conclude this paper and discuss the future work in Section 5.

2. RELATED WORK

Existing semi-supervised clustering techniques, employing user supervision in the form of instance-level constraints, are generally grouped into four categories. First, constraints are used to modify the loss function [4, 18, 25]. Second, cluster seeds derived from the constraints initialize the cluster

centers [3]. Third, constraints are employed to learn adaptive distance metrics using metric learning techniques [2, 6]. Finally, the original high-dimensional feature space can be projected into low-dimensional feature subspaces guided by constraints [24]. However, alternative forms of user supervision exist when we apply semi-supervised clustering algorithms to group documents. In this paper, we explore words labeled by being associated with a document when the document is assigned to a cluster.

Liu et al. [20] propose to ask the user to label features for each class and use the set of features labeled for each class to label a set of documents for training classifiers. Druck et al. [9] use labeled features for each class to constrain the probabilistic model estimation on unlabeled instances instead of creating pseudo-instances as done in other approaches. Sindhwani and Melville [23] present a novel semi-supervised sentiment prediction method which use both labeled documents and features for each class to train the classifiers. Raghavan et al. [22] make use of feature feedback in the active learning with support vector machine by up-weighting the accepted features. Unlike the above classification methods which require labeled documents and/or features for each class, our framework can deal with partial clusters with labeled documents and/or features. In addition, it explores the unlabeled documents to refine the prior knowledge provided by the user. Huang and Mitchell [17] propose a generative probabilistic framework to incorporate various types of user feedback including feedback on features. In their work, the user needs to assign a feature to an intermediate cluster, which requires the user browse the intermediate clusters and understand them. In our framework, the user associates the features with documents through text clouds, which is much easier and more convenient than understanding intermediate clusters. Hu et al. [11] propose an interactive framework for feature selection for document clustering, in which the user only indicates whether a feature is suitable for clustering. However, they ask the user to label features from a standalone ranked list of features, which requires extra effort for labeling. In addition, they did not explore the usefulness of integrating labeling documents and features together or compare feature supervision with document supervision for clustering.

3. METHODOLOGY

In this section, we first briefly describe basic K Means algorithm and then present a unified framework to combine the document supervision, feature supervision, and unlabeled documents.

3.1 Background

K Means [5] is a clustering algorithm based on iterative assignments of data points to clusters and partitions a dataset into K clusters so that the average squared distance between the data points and the closest cluster centers are locally minimized. For a dataset with data points $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$, $x_i \in \mathbb{R}^d$, K Means algorithm generates K clusters $\{\mathcal{X}_l\}_{l=1}^K$ of \mathcal{X} so that the objective function

$$J = \sum_{l=1}^K \sum_{x_i \in \mathcal{X}_l} \|x_i - \mu_l\|^2 \quad (1)$$

is locally minimized and $\{\mu_1, \mu_2, \dots, \mu_K\}$ represents the cen-

ters of the K clusters.

3.2 Algorithms

In this section, we first introduce document supervision and feature supervision in the form of document seeding and feature seeding separately. Then, we present two methods to model feature seeding. At the end, we describe a unified framework to incorporate both document seeding and feature seeding into the K Means algorithm, namely, *DualSeededK Means*.

3.2.1 Document Seeding

Given a dataset \mathcal{X} , as previously described, K Means can partition it into K clusters $\{\mathcal{X}_i\}_{i=1}^K$. Then, we can define the document seed set $\mathcal{D}^L \subseteq \mathcal{X}$ as the following subset of data points: for each $x_i \in \mathcal{D}^L$, the user provides the cluster \mathcal{X}_i to which it belongs. We assume that there is at least one data point x_i for each cluster \mathcal{X}_i . Note that there is a K -disjoint partitioning $\{\mathcal{D}_i^L\}_{i=1}^K$ of the seed set \mathcal{D}^L such that all $x_i \in \mathcal{D}_i^L$ belong to \mathcal{X}_i according to the supervision. We define the centers of the document seed set $\{\mathcal{D}_i^L\}_{i=1}^K$ as $\{\mu_i^d\}_{i=1}^K$:

$$\mu_i^d = \frac{\sum_{x_i \in \mathcal{D}_i^L} x_i}{|\mathcal{D}_i^L|} \quad (2)$$

Those seed centers can be used to both initialize the clustering algorithms and guide the clustering process.

3.2.2 Feature Seeding

Similar to document seed set \mathcal{D}^L , we can define the feature seed set \mathcal{W}^L as the following subset of features: for each $w_i \in \mathcal{W}^L$, the user indirectly associates it with the cluster \mathcal{X}_i through document $x_j \in \mathcal{X}_i$ in which w_i occurs and is labeled from. We assume that each cluster has a topic and at least one feature is associated with it. Note that there does not exist a K -disjoint partitioning $\{\mathcal{W}_i^L\}_{i=1}^K$ of the feature seed set because one feature can be associated with multiple clusters. We define the centers of the feature seed set $\{\mathcal{W}_i^L\}_{i=1}^K$ as $\{\mu_i^w\}_{i=1}^K$, which can be derived from either feature-vote-model (see Section 3.2.4 for details) or feature-generative-model (see Section 3.2.5 for details). And then those seed centers can be used to both initialize the clustering algorithms and guide the clustering process.

3.2.3 Feature supervision

A document d can be considered as a list of words in the order in which the words appear in the document, i.e., $\langle w_1, w_2, \dots, w_{|d|} \rangle$, where $|d|$ is the length of the document in terms of the number of words. Note that w_i might be the same as w_j where $i \neq j$, $1 \leq i \leq |d|$ and $1 \leq j \leq |d|$. To label a document, we assume that the user needs to read at least a fraction of the document content, i.e., $\langle w_s, w_2, \dots, w_e \rangle$, where $1 \leq s \leq |d|$ and $s \leq e \leq |d|$. While reading a document, the user is assumed to be able to recognize useful words for clustering. The useful words should describe the cluster topic of the document from which they are labeled. The fraction of the document content could be displayed as a text cloud and the user could label words by highlighting them through double-clicking on the text clouds. The user labels a feature if it is a good description of the topic of a cluster and discriminates the cluster from others. Then, a labeled feature is associated with a cluster indirectly through the labeled documents from which it is

labeled. After a cluster being created, additional features can be associated with a cluster directly by being assigned into the cluster.

3.2.4 Feature-Vote-Model

In this method, we use the labeled features in the feature seed set to vote on cluster labels for the unlabeled documents. A similar approach was introduced for document classification [9, 27]. For each labeled feature w in a document x , it contributes one vote for each of its cluster labels (could be associated with multiple clusters). Then, we normalize the vote totals to get a probabilistic distribution over the cluster labels for each document, i.e., $\{P_{li}\}$ for document x_i and cluster \mathcal{X}_i . Assume document x_i contains n_{il} labeled features for cluster \mathcal{X}_i , we define

$$P_{li} = \frac{n_{il}}{\sum_{k=1}^K n_{ik}} \quad (3)$$

With this soft labeled documents, we can derive the center of μ_i^w from the feature seed set as:

$$\mu_i^w = \sum_{x_i \in \mathcal{X}} P_{li} x_i \quad (4)$$

where x_i is the vector of TFIDF values of the features selected for clustering.

3.2.5 Feature-Generative-Model

This model was introduced for binary sentiment analysis [21] and we extend it for document clustering with multiple clusters. In this method, we generate each cluster center from the feature seed set directly. We choose to represent the cluster center as a multinomial distribution which generates documents for the corresponding cluster. Without losing generality, we derive the cluster center for cluster \mathcal{X}_i and *words* and *features* are used interchangeably. We define the following notations to aid our derivations:

- \mathcal{V} – set of words used for clustering, including both labeled and unlabeled words
- $\mathcal{P}_{\mathcal{X}_i}$ – set of words labeled for cluster \mathcal{X}_i
- $\mathcal{N}_{\mathcal{X}_i}$ – set of words labeled for the other clusters
- \mathcal{U} – set of unlabeled words used for clustering
- m – size of vocabulary, i.e. $|\mathcal{V}|$
- $p_{\mathcal{X}_i}$ – number of words labeled for cluster \mathcal{X}_i , i.e. $|\mathcal{P}_{\mathcal{X}_i}|$
- $n_{\mathcal{X}_i}$ – number of words labeled for the other clusters, i.e. $|\mathcal{N}_{\mathcal{X}_i}|$

In order to derive the multinomial distribution for cluster center of \mathcal{X}_i , we assume the following properties about the relationships between words and clusters.

Property 1: All words in $\mathcal{P}_{\mathcal{X}_i}$ are equally likely to occur in a document from cluster \mathcal{X}_i .

$$P(w_i|\mathcal{X}_i) = P(w_j|\mathcal{X}_i), \forall w_i, w_j \in \mathcal{P}_{\mathcal{X}_i} \quad (5)$$

We refer to the probability of any word in $\mathcal{P}_{\mathcal{X}_i}$ appearing in a document from cluster \mathcal{X}_i simply as $P(w_p|\mathcal{X}_i)$.

Property 2: All words in $\mathcal{N}_{\mathcal{X}_i}$ are equally likely to occur in a document from cluster \mathcal{X}_i .

$$P(w_i|\mathcal{X}_i) = P(w_j|\mathcal{X}_i), \forall w_i, w_j \in \mathcal{N}_{\mathcal{X}_i} \quad (6)$$

We refer to the probability of any word in $\mathcal{N}_{\mathcal{X}_i}$ appearing in a document from cluster \mathcal{X}_i simply as $P(w_n|\mathcal{X}_i)$.

Property 3: The unlabeled words are treated equally in each cluster.

$$P(w_i|\mathcal{X}_i) = P(w_j|\mathcal{X}_i), \forall w_i, w_j \in \mathcal{U} \quad (7)$$

We refer to the probability of any word in \mathcal{U} appearing in a document from cluster \mathcal{X}_i simply as $P(w_u|\mathcal{X}_i)$.

Property 4: A document from cluster \mathcal{X}_i is more likely to contain a word from $\mathcal{P}_{\mathcal{X}_i}$ than a word from $\mathcal{N}_{\mathcal{X}_i}$

$$P(w_p|\mathcal{X}_i) = r \times P(w_n|\mathcal{X}_i) \quad (8)$$

where r is referred to as polarity level, which measures how much more likely a word in $\mathcal{P}_{\mathcal{X}_i}$ occurs in a document from cluster \mathcal{X}_i compared with a word in $\mathcal{N}_{\mathcal{X}_i}$. Since a word in $\mathcal{P}_{\mathcal{X}_i}$ is more likely occurs in a document from cluster \mathcal{X}_i , we have $0 < 1/r \leq 1$.

Property 5: The multinomial probability distribution learned from labeled features for each cluster is constrained by summing to one.

$$\sum_i^m P(w_i|\mathcal{X}_i) = 1 \quad (9)$$

We use property 5 as constraints to derive the appropriate probability distribution based on labeled features. By Eq. 9 it follows that

$$pP(w_p|\mathcal{P}_{\mathcal{X}_i}) + nP(w_n|\mathcal{P}_{\mathcal{X}_i}) + (m-p-n)P(w_u|\mathcal{P}_{\mathcal{X}_i}) = 1 \quad (10)$$

which gives us the following inequality using Eq. 8,

$$\begin{aligned} pP(w_p|\mathcal{X}_i) + nP(w_n|\mathcal{X}_i) &\leq 1 \\ \Rightarrow pP(w_p|\mathcal{X}_i) + n\frac{P(w_p|\mathcal{X}_i)}{r} &\leq 1 \end{aligned}$$

Since $0 < 1/r \leq 1$, it follows that,

$$P(w_p|\mathcal{X}_i) \leq \frac{1}{p+n}$$

By assigning the maximum probability mass to the known words, $P(w_p|\mathcal{X}_i)$ is set to the maximum value possible, i.e.

$$P(w_p|\mathcal{X}_i) = \frac{1}{p+n} \quad (11)$$

Now, it follows from Eq. 8,

$$P(w_n|\mathcal{X}_i) = \frac{1}{p+n} \times \frac{1}{r} \quad (12)$$

Now, solving Eq. 10, we can have the probabilities for the unlabeled words:

$$P(w_u|\mathcal{X}_i) = \frac{n(1-1/r)}{(p+n)(m-p-n)} \quad (13)$$

Finally, we use Eqs. 11, 12 and 13 to derive the center μ_i^w of cluster \mathcal{X}_i . The cluster center μ_i^w is defined as a vector, whose elements are the probabilities of words in \mathcal{V} given the cluster \mathcal{X}_i , namely,

$$\mu_i^w = (P(w_1|\mathcal{X}_i), P(w_2|\mathcal{X}_i), \dots, P(w_m|\mathcal{X}_i)) \quad (14)$$

where $w_i \in \mathcal{V}$ and $m = |\mathcal{V}|$ as previously defined.

In our experiments, we set $r = 100$ based on previous experimental results [21].

3.2.6 Combining Multiple Centers

Opinion pool is a general approach to combine information from multiple sources, such as the centers derived from document seed set and feature seed set in our document clustering problem. Particularly, we use *linear opinion pool* approach to aggregate multiple centers. which was used to

Algorithm 1 DualSeededK Means

Input: Set of data points \mathcal{X} , the document seed set $\mathcal{D}^L = \cup_{l=1}^K \mathcal{D}_l^L$, the feature seed set $\mathcal{W}^L = \cup_{l=1}^K \mathcal{W}_l^L$

Output: K clusters $\{\mathcal{X}_i\}_{i=1}^K$

Method:

- 1: Compute $\{\mu_l^d\}$ from $\{\mathcal{D}_l^L\}$ using Eq. 2
 - 2: Compute $\{\mu_l^w\}$ from $\{\mathcal{W}_l^L\}$ using Eq. 4 or Eq. 14
 - 3: initialize: $\mu_l^{(0)} = \alpha_d \mu_l^d + \alpha_w \mu_l^w$, for $l = 1, \dots, K; t \leftarrow 0$
 - 4: **repeat**
 - 5: **for all** $x_i \in \mathcal{X}$ **do**
 - 6: Assign x_i to the closest cluster $\mathcal{X}_l^{(t+1)}$ based on $\{\mu_l^t\}$ and get $\{\mathcal{X}_l^{(t+1)}\}_{l=1}^K$
 - 7: **end for**
 - 8: Update intermediate cluster centers:
 $\mu_l^c \leftarrow \frac{1}{|\mathcal{X}_l^{(t+1)}|} \sum_{x \in \mathcal{X}_l^{(t+1)}} x$
 - 9: Update cluster centers:
 $u_l^{(t+1)} \leftarrow \alpha_d \mu_l^d + \alpha_w \mu_l^w + \alpha_c \mu_l^c$
 - 10: $t \leftarrow t + 1$
 - 11: **until convergence**
-

combine probability distributions for text classification [21]. In this approach, the aggregated (pooling) center is defined as

$$\mu_l = \sum_{s=1}^S \alpha_s \mu_l^s \quad (15)$$

where S is the number of sources we have.

In addition, we compute the weights α 's of individual sources based on their error in labeling the document seed set. In particular, we use the same weighting scheme as [21]:

$$\alpha_s = \log \frac{1 - err_s}{err_s} \quad (16)$$

where err_s is the classification error of the source s when the derived centers based on the information provided by the source s are used to classify the documents in the document seed set. All α_s 's are normalized to one.

3.2.7 Dual Semi-supervised KMeans

In *DualSeededK Means*, both the document seeds and feature seeds are used to initialize the K Means algorithm through derived cluster centers. To this end, the center of the l^{th} cluster is initialized with the pooling center derived from both μ_l^d and μ_l^w (Eq. 15) before the clustering starts. During the clustering, the cluster centers are refined using the information contained in the intermediate clusters. This information is expressed in the form of intermediate cluster centers μ_l^c

$$\mu_l^c = \frac{\sum_{x_i \in \mathcal{X}_l^c} x_i}{|\mathcal{X}_l^c|} \quad (17)$$

where \mathcal{X}_l^c is the l^{th} intermediate cluster. Then, we can incorporate μ_l^c to the *DualSeededK Means* algorithm using the *linear opinion pool* technique (Eq. 15). The algorithm is described in details in Alg. 1. Note that *DualSeededK Means* can be specialized to *DocumentSeededK Means* when feature seed set is empty and *FeatureSeededK Means* when document seed set is empty.

Table 1: Six Datasets from the 20-newsgroups, Webkb, Industry Sectors and Reuters21578

Dataset	Description	Categories included	Category Doc.	Total Doc.
news-similar-3-100 (D1)	The 20-Newsgroup data set consists of 20 different Usenet newsgroups, each of which has approximately 1000 newsgroup messages.	3:comp.graphics,comp.os.ms-windows.misc,comp.windows.x	100	300
news-multi-7-100 (D2)		7:alt.atheism,comp.sys.mac.hardware, misc.forsale,rec.sport.hockey,sci.crypt, talk.politics.guns,soc.religion.christian	100	700
news-multi-10-100 (D3)		10:alt.atheism,comp.sys.mac.hardware,misc.forsale,rec.autos,rec.sport.hockey,sci.crypt,sci.med, sci.electronics, sci.space, talk.politics.guns	100	1000
webkb-sfcp-4-250 (D4)	webpages from different universities	4:student, faculty, course, project	250	1000
sector-multi-10-100 (D5)	webpages from different industrial sectors	10:basic.materials,capital.goods,consumer.cyclical, oil.and.gas.integrated, investment.services, biotechnology.and.drugs, hotels.and.motels, communications.equipment, railroad, water.utilities	100 (railroad-95)	995
reuters-multi-10-100 (D6)	news articles from Reuters21578. We use the top 10 most frequent categories, documents of which does not have multiple labels.	10:acq, coffee, crude, earn, gold, interest, money-fx, ship, sugar, trade	100 (gold-90)	990

3.3 Oracles

Most research involving labeling documents simulates human input by a document oracle that uses the underlying class labels of documents in the dataset [1, 3, 4, 6, 16, 15, 18, 24]. However, in the case of features, we do not have a gold-standard set of feature labels. Ideally, we should have a human expert in the loop labeling the selected features. However, such a manual process is not feasible for repetitive large-scale experiments. Therefore, we construct a feature oracle similar to the method described by [9, 10, 11, 13, 14]. Using the document labels, the oracle computes the χ^2 value of each feature with cluster/class label, and accept a feature if the χ^2 value is above a threshold β . In this paper, the β value is the mean of the top f most predictive features, where $f = 100K$, namely, 100 times the number of clusters. If accepted, the feature oracle labels a feature with the cluster in which it occurs the most and any other clusters in which the feature occurs at least half of the most occurrences.

4. EXPERIMENTAL RESULTS

4.1 Datasets

We conducted our experiments on several real-word datasets of different sizes and also consisting of different types of text documents. We derive three datasets of different sizes from the 20-Newsgroup corpus² and three more datasets from webkb³, industry sector⁴, and reuters21578⁵ separately. The datasets are different from each other in terms of sizes of the datasets and types of documents, i.e., webpages, newsgroup messages, etc. The descriptions and details of the datasets are summarized in Table 1.

We pre-processed each document by tokenizing the text into bags-of-words⁶. Then, we removed the stop words and stemmed all the remaining words. Next, we selected the top 2000

²<http://people.csail.mit.edu/jrennie/20Newsgroups/>

³<http://www.cs.cmu.edu/webkb>

⁴<http://www.cs.umass.edu/mccallum/data.html>

⁵<http://kdd.ics.uci.edu>

⁶A word is defined as a sequence of alphabetic characters delimited by non-alphabetic characters.

words using mutual information between words and documents [7]. Finally, a feature vector for each document is constructed with TFIDF weighting and then normalized.

4.2 Evaluation Measures

In this paper, we employed normalized mutual information (NMI) [8] as the clustering evaluation measure. NMI measures the share information between the cluster assignments S and class labels L of documents. It is defined as:

$$NMI(S, L) = \frac{I(S, L)}{(H(S) + H(L))/2} \quad (18)$$

where $I(S, L)$, $H(S)$, and $H(L)$ denote the mutual information between S and L , the entropy of S , and the entropy of L respectively. The range of NMI values is 0 to 1.

4.3 Analysis of Results

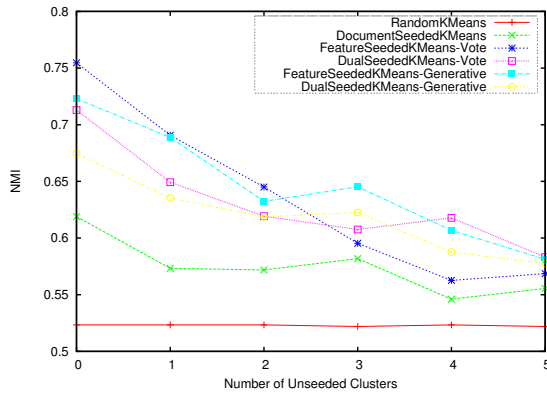
First, we have two sets of comparisons in our experiments. The first set of comparisons is designed to see whether the user provided information can be refined by the intermediate clusters, i.e., clustering models incorporating unlabeled documents categorize documents better than classification models which only use labeled information. The clustering and classification models are defined as:

- *DualSeededK Means*, or its specialized algorithms when one of the seed set is empty, i.e., *DocumentSeededK Means* and *FeatureSeededK Means*. Note that *FeatureSeededK Means* has two variants, namely, *Feature-Vote-Model* and *Feature-Generative-Model* to derive cluster centers.
- *SupervisedK Means*, which performs clustering by assigning documents to nearest cluster centers inferred from either document seed set or feature seed set or both. It can be achieved by running the *DualSeededK Means* or its specialized cases, i.e., *DocumentSeededK Means* and *FeatureSeededK Means*, with only one iteration. Correspondingly, we have *DualSupervisedK Means*, *DocumentSupervisedK Means*, and *FeatureSupervisedK Means*.

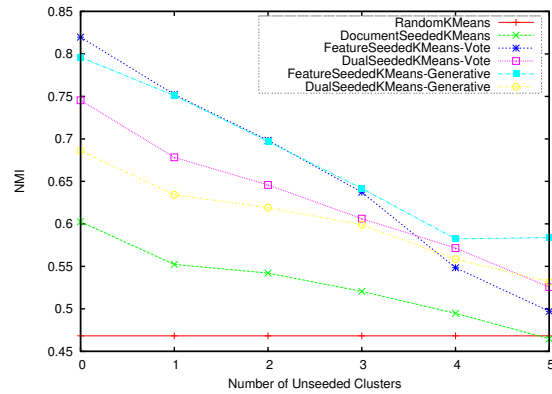
We did thorough pair comparisons (Table 2) to demonstrate

Table 2: Supervised K Means compared to peer algorithms refined by intermediate clusters. 10 documents are labeled for each cluster and features are labeled by feature oracle from the labeled documents. We did two-tailed paired t-test with $p = 0.05$ for comparing pairs of algorithms. In this table, we compare algorithms by pairs, i.e., DocumentSeeded K Means vs. DocumentSupervised K Means, FeatureSeeded K Means vs. FeatureSupervised K Means using Feature-Vote-Model and Feature-Generative-Model. DualSeeded K Means vs. DualSupervised K Means using Feature-Vote-Model and Feature-Generative-Model. All algorithms refined by intermediate clusters works significantly better than peer Supervised K Means algorithm except FeatureSeeded K Means and FeatureSupervised K Means using Feature-Vote-Model on D3 (news-multi-10-100) and DualSeeded K Means and DualSupervised K Means using Feature-Vote-Model on D1 (news-similar-3-100) indicated by *.

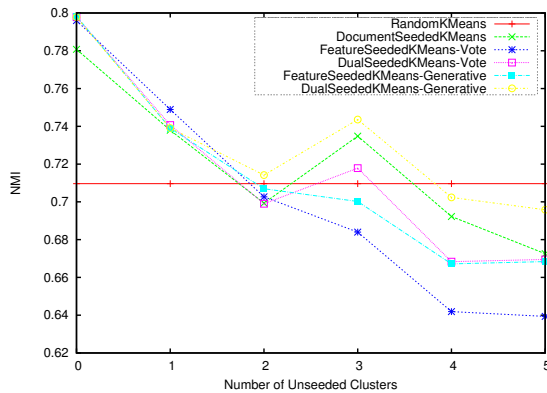
Supervision	Algorithm		D1	D2	D3	D4	D5	D6
No Supervision	Basic K Means		0.069	0.523	0.468	0.341	0.710	0.350
Document Only	DocumentSeeded K Means		0.276	0.692	0.686	0.397	0.815	0.637
	DocumentSupervised K Means		0.266	0.625	0.624	0.319	0.786	0.581
Feature Only	Feature-Vote-Model	FeatureSeeded K Means	0.551	0.770	0.820*	0.464	0.795	0.649
		FeatureSupervised K Means	0.548	0.766	0.820*	0.428	0.791	0.637
	Feature-Generative-Model	FeatureSeeded K Means	0.515	0.724	0.791	0.470	0.805	0.692
		FeatureSupervised K Means	0.512	0.681	0.747	0.413	0.734	0.660
Dual Supervision	Feature-Vote-Model	DualSeeded K Means	0.482*	0.757	0.783	0.421	0.822	0.687
		DualSupervised K Means	0.482*	0.745	0.765	0.372	0.815	0.660
	Feature-Generative-Model	DualSeeded K Means	0.423	0.732	0.738	0.443	0.824	0.684
		DualSupervised K Means	0.421	0.703	0.700	0.391	0.812	0.642



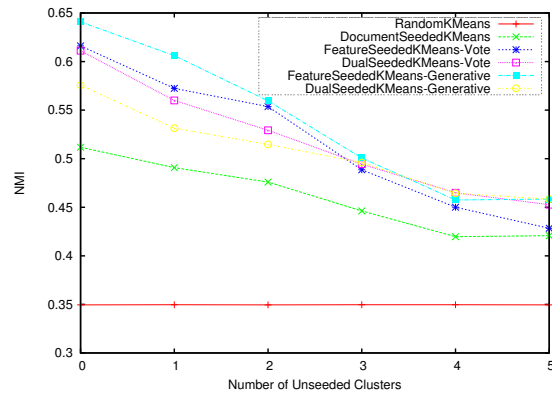
(a) news-multi-7-100



(b) news-multi-10-100



(c) sector-multi-10-100

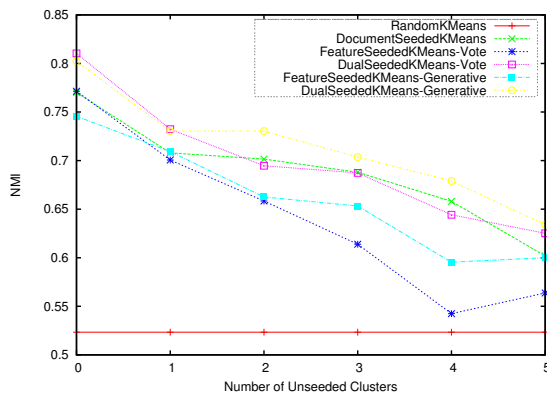


(d) reuters-multi-10-100

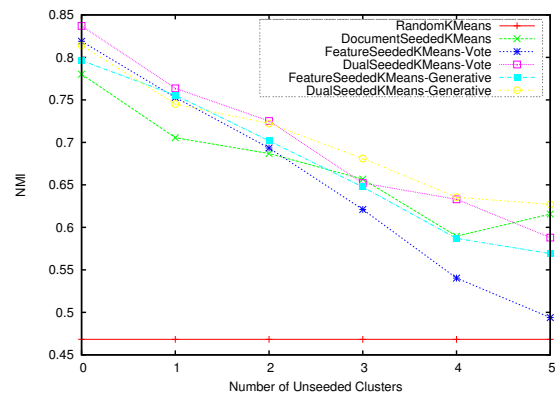
Figure 2: Performance as a function of the Number of Unseeded Clusters. 5 Documents Are Labeled for Each Seeded Cluster where FeatureSeeded K Means works better than DocumentSeeded K Means and DualSeeded K Means

Table 3: Comparison of algorithms with dual supervision to algorithms with any single supervision. 20 documents are labeled for each cluster and features are labeled by feature oracle from those labeled documents. We did two-tailed paired t-test with $p = 0.05$ for comparing pairs of algorithms. In this table, we compared *DualSeededKMeans* with *DocumentSeededKMeans*, *DualSeededKMeans* with *FeatureSeededKMeans* using *Feature-Vote-Model* or *Feature-Generative-Model*. *DualSeededKMeans* works better than *DocumentSeededKMeans* on all datasets. *DualSeededKMeans* works better than *FeatureSeededKMeans* on all datasets except D1 (*news-similar-3-100*) and D4 (*webkb-sfcp-4-250*) with *Feature-Generative-Model* indicated by *.

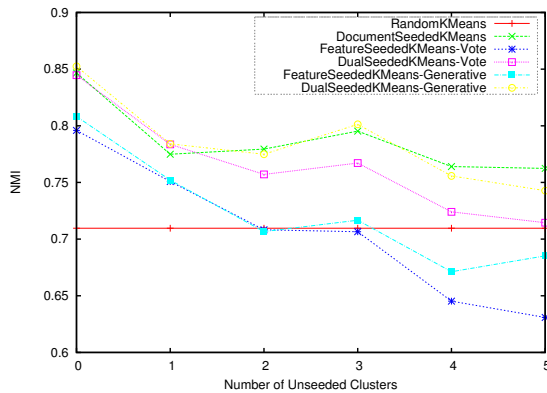
Feature Model	Algorithm		D1	D2	D3	D4	D5	D6
No Supervision	Basic <i>KMeans</i>		0.069	0.523	0.468	0.341	0.710	0.350
Document Only	<i>DocumentSeededKMeans</i>		0.416	0.770	0.780	0.466	0.847	0.767
Feature-Vote-Model	Feature Only	<i>FeatureSeededKMeans</i>	0.560	0.771	0.819	0.468	0.796	0.679
	Dual Supervision	<i>DualSeededKMeans</i>	0.561	0.810	0.837	0.484	0.845	0.786
Feature-Generative-Model	Feature Only	<i>FeatureSeededKMeans</i>	0.515*	0.746	0.796	0.504*	0.808	0.736
	Dual Supervision	<i>DualSeededKMeans</i>	0.507*	0.802	0.814	0.502*	0.852	0.797



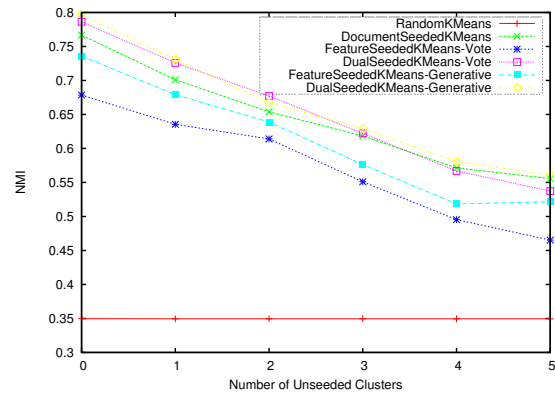
(a) news-multi-7-100



(b) news-multi-10-100

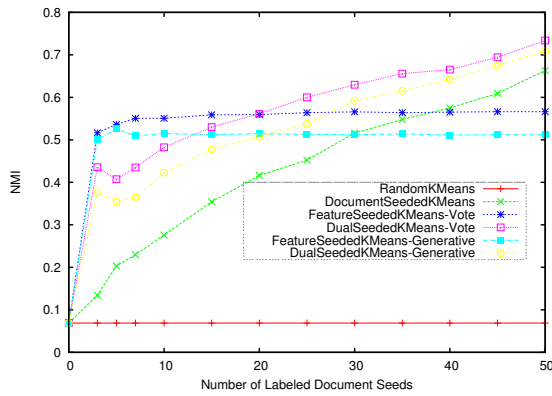


(c) sector-multi-10-100

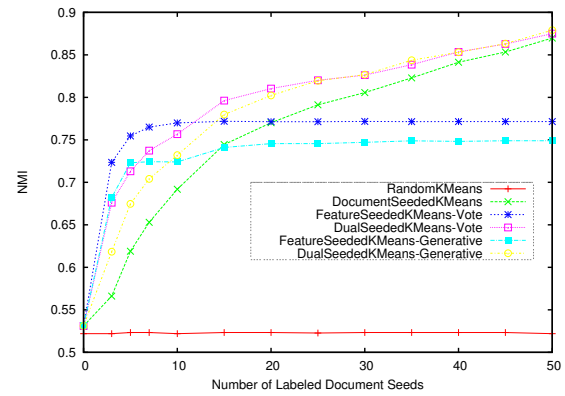


(d) reuters-multi-10-100

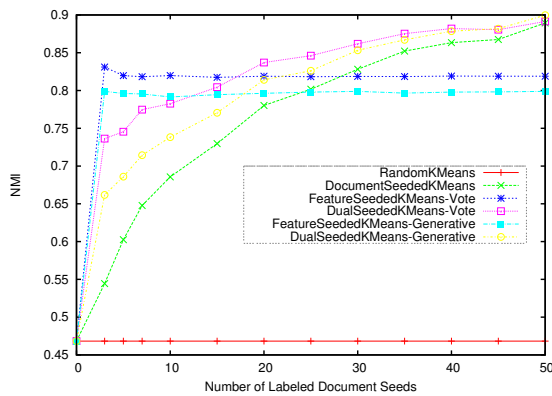
Figure 3: Performance as a Function of the Number of Unseeded Clusters. 20 Documents Are Labeled for Each Seeded Cluster where *DualSeededKMeans* works better than *DocumentSeededKMeans* and *FeatureSeededKMeans*



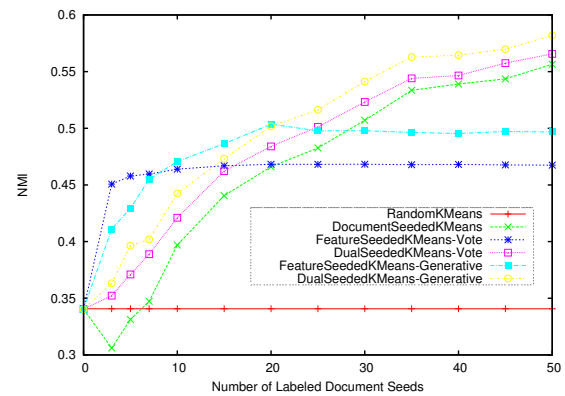
(a) news-similar-3-100



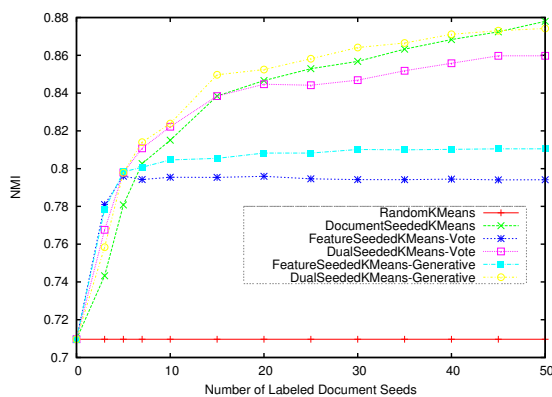
(b) news-multi-7-100



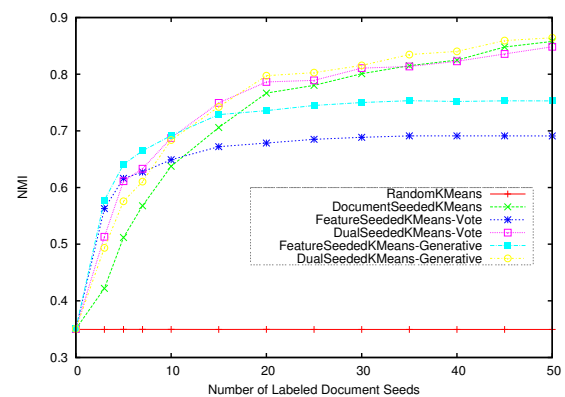
(c) news-multi-10-100



(d) webkb-sfcp-4-250



(e) sector-multi-10-100



(f) reuters-multi-10-100

Figure 4: Performance as a Function of the Number of Labeled Documents. The more documents labeled, the more features labeled and the better performance. The usefulness of labeled features are more obvious when there are only a few documents labeled, e.g., < 10 . In fact, the feature supervision even works better than dual supervision at the beginning of the curves, indicating that feature supervision is more reliable when only few documents are labeled.

that incorporating unlabeled documents can refine the information provided by the user and produce better clusters than only using labeled information. Concretely, we compared the following pairs of algorithms:

- DocumentSeedsedKMeans vs. DocumentSupervisedKMeans
- FeatureSeededKMeans vs. FeatureSupervisedKMeans using Feature-Vote-Model or Feature-Generative-Model.
- DualSeededKMeans vs. DualSupervisedKMeans using Feature-Vote-Model or Feature-Generative-Model.

From Table 2, we can tell that all algorithms with refinement by intermediate clusters improve its clustering performance over the peer algorithms of SupervisedKMeans except when Feature-Vote-Model with only feature supervision works on dataset D3 (news-multi-10-100) and DualSeededKMeans and DualSupervisedKMeans using Feature-Vote-Model on D1 (news-similar-3-100) (indicated by * in Table 2), which is only 2 out of 30 comparisons. Therefore, intermediate clusters are helpful in improving clustering performance in addition to labeled information.

The second set of comparisons is designed to see whether dual supervision performs better than any single supervision. Thus, we compare *DualSeededKMeans* with *DocumentSeededKMeans*, and *FeatureSeededKMeans*. Again, we have two variants when feature seed set is involved. From Table 3, we can tell that dual supervision with both document and feature generally improve the clustering performance over any single supervision except with feature-generative-model on D1 (news-similar-3-100) and D4 (webkb-sfcp-4-250) indicated by * in Table 3, which is only 2 out of 24 comparisons. Note that algorithms with dual supervision works better than document only supervision on all datasets. Therefore, it is worth labeling features.

Second, we ran experiments with incomplete seeding, namely, only a fraction of categories are seeded by labeled documents, or labeled features, or both (Fig. 2 and Fig. 3). It can be seen that the performances decreases with increase number of unseeded clusters. However, the performances do not decrease substantially, showing that *DualSeededKMeans* can extend the seeded clusters and generate more clusters to fit the regularities in the dataset. Therefore, not all clusters have to be seeded by labeled information.

Finally, we study the behaviors of the *DualSeededKMeans* with different numbers of document seeds. Note that the more document seeds labeled, the more feature seeds labeled because the feature seeds are labeled while a document is being labeled. We have the following observations from Fig. 4.

- *DualSeededKMeans* always works better than *DocumentSeededKMeans*. However, the performances of the two algorithms are getting close when more documents are provided. It suggests that the feature labeling is more useful when there are few documents labeled, i.e., little effort. One of the possible explanations is that few labeled documents can not represent the cluster

very well, which can be enhanced by the labeled features at the beginning of the learning curve. However, with enough documents labeled, the cluster structures can be represented pretty well with only documents so that the dual supervision has similar performance to document supervision only.

- When there are only few documents labeled, *FeatureSeededKMeans* (fewer feature seeds) performs better than *DualSeededKMeans* and *DocumentSeededKMeans*. It suggests that feature supervision is more reliable than document supervision when only little supervision can be provided. However, *DualSeededKMeans* and *DocumentSeededKMeans* improve their performances quickly than *FeatureSeededKMeans* when more document seeds labeled. When there are enough document seeds labeled, both *DualSeededKMeans* and *DocumentSeededKMeans* performs better than *FeatureSeededKMeans*. Our explanation is that a few labeled features can represent the cluster structures better than a few documents, which also contains other non-discriminating features. Therefore, it is better to label features than documents if only limited user supervision is available.
- Learning curves of *FeatureSeededKMeans* are steep at the beginning but become flat quickly. Our explanation is that enough feature seeds are labeled after a few document seeds labeled at the beginning. The number of feature seeds labeled does not change much when more document seeds are labeled later.

5. CONCLUSIONS AND FUTURE WORK

In this paper, we incorporate document supervision and feature supervision in the form of feature seeding. *DualSeededKMeans* is a unified framework to combine document supervision, feature supervision and unlabeled documents in the form of seeding. *DocumentSeededKMeans* and *FeatureSeededKMeans* are two specialized cases of *DualSeededKMeans*. Experimental results demonstrate that unlabeled documents can help to refine the information provided by the user and feature supervision is worth the effort to improve the clustering performance further compared to document supervision only and much more helpful when only few documents can be labeled due to manually cost.

The research presented in this paper is in the context of a document management system that support user-driven organization of document collections. Evaluation of the effectiveness of the system through user studies is available in [12].

6. ACKNOWLEDGMENTS

We would like to thank anonymous reviewers for their insightful comments. This research was supported in part by the NSERC (Natural Sciences and Engineering Research Council) Business Intelligence Network, and by the MITACS NCE.

References

- [1] J. Attenberg, P. Melville, and F. Provost. A Unified Approach to Active Dual Supervision for Labeling

- Features and Examples. In *ECML PKDD 2010 Part I, LNAI 6321*, pages 40–55. Springer, 2010.
- [2] A. Bar-Hillel, T. Hertz, N. Shental, and D. Weinshall. Learning distance functions using equivalence relations. In *International Conference on Machine Learning*, volume 20, page 11, 2003.
- [3] S. Basu, A. Banerjee, and R. Mooney. Semi-supervised clustering by seeding. In *International Conference on Machine Learning*, pages 19–26, 2002.
- [4] S. Basu, M. Bilenko, and R.J. Mooney. A probabilistic framework for semi-supervised clustering. In *Proceedings of the tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 59–68. ACM, 2004.
- [5] C.M. Bishop. *Pattern Recognition and Machine Learning*. Springer New York, 2006.
- [6] H. Cheng, K.A. Hua, and K. Vu. Constrained locally weighted clustering. *Proceedings of VLDB'08*, 1(1): 90–101, 2008.
- [7] I.S. Dhillon, S. Mallela, and D.S. Modha. Information-theoretic co-clustering. In *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 89–98. ACM, 2003. ISBN 1581137370.
- [8] B.E. Dom. An information-theoretic external cluster-validity measure. Technical Report RJ 10219, IBM Research Division, 2001.
- [9] G. Druck, G. Mann, and A. McCallum. Learning from labeled features using generalized expectation criteria. In *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 595–602. ACM, 2008.
- [10] Y. Hu, E.E. Milios, and J. Blustein. Interactive Document Clustering Using Iterative Class-Based Feature Selection. Technical report, CS-2010-04, Faculty of Computer Science, Dalhousie University, Canada, 2010.
- [11] Y. Hu, E.E. Milios, and J. Blustein. Interactive feature selection for document clustering. In *Proceedings of the 26th Symposium On Applied Computing, On Track “Information Access and Retrieval”*, pages 1148–1155. ACM Special Interest Group on Applied Computing, 2011.
- [12] Y. Hu, E.E. Milios, and J. Blustein. Personalized document clustering with dual supervision. In *Proceedings of the 12th ACM Symposium on Document Engineering*. ACM, 2012.
- [13] Y. Hu, E.E. Milios, and J. Blustein. Enhancing Semi-supervised Document Clustering with Feature Supervision. In *Proceedings of the 27th ACM Symposium Applied Computing, On Track “Information Access and Retrieval”*, pages 950–957. ACM, 2012.
- [14] Y. Hu, E.E. Milios, and J. Blustein. Semi-supervised Document Clustering with Dual Supervision through Seeding. In *Proceedings of the 27th ACM Symposium Applied Computing, On Track “Data Mining”*, pages 463–470. ACM, 2012.
- [15] R. Huang and W. Lam. An active learning framework for semi-supervised document clustering with language modeling. *Data & Knowledge Engineering*, 68(1): 49–67, 2009.
- [16] R. Huang, W. Lam, and Z. Zhang. Active learning of constraints for semi-supervised text clustering. In *Proceedings of the SIAM International Conference on Data Mining*, pages 113–124, 2007.
- [17] Y. Huang and T.M. Mitchell. Text clustering with extended user feedback. In *Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, page 420. ACM, 2006.
- [18] X. Ji and W. Xu. Document clustering with prior knowledge. In *Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, page 412. ACM, 2006.
- [19] J. Lamantia. Text Clouds: A New Form of Tag Cloud? <http://www.joelamantia.com/tag-clouds/text-clouds-a-new-form-of-tag-cloud>, 2007. Accessed on April 12, 2012.
- [20] B. Liu, X. Li, W.S. Lee, and P.S. Yu. Text classification by labeling words. In *Proceedings of the National Conference on Artificial Intelligence*, pages 425–430, 2004.
- [21] P. Melville, W. Gryc, and R.D. Lawrence. Sentiment analysis of blogs by combining lexical knowledge with text classification. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1275–1284, 2009.
- [22] H. Raghavan, O. Madani, and R. Jones. Interactive feature selection. In *Proceedings of IJCAI 05: The 19th International Joint Conference on Artificial Intelligence*, pages 841–846, 2005.
- [23] V. Sindhwani and P. Melville. Document-word co-regularization for semi-supervised sentiment analysis. In *Proceedings of the Eighth IEEE International Conference on Data Mining*, pages 1025–1030. IEEE, 2009.
- [24] W. Tang, H. Xiong, S. Zhong, and J. Wu. Enhancing semi-supervised clustering: a feature projection perspective. In *Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 707–716. ACM, 2007.
- [25] K. Wagstaff. *Intelligent Clustering with Instance-Level Constraints*. PhD thesis, Cornell University, 2002.
- [26] K. Wagstaff, C. Cardie, S. Rogers, and S. Schrödl. Constrained k-means clustering with background knowledge. In *Proceedings of the Eighteenth International Conference on Machine Learning*, pages 577–584, 2001.
- [27] X. Wu and R. Srihari. Incorporating prior knowledge with weighted margin support vector machines. In *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 326–333. ACM, 2004. ISBN 1581138881.

ABOUT THE AUTHORS:



Yeming Hu is a PhD candidate at Faculty of Computer Science, Dalhousie University (Canada). He got his Master of Engineering from Dalhousie University in 2006 and his B.Sc. from Sun Yat-sen University (China) in 2004. His PhD thesis focuses on interactive document clustering for producing personalized clusters for users.



Evangelos E. Milios received a diploma in Electrical Engineering from the NTUA, Athens, Greece, and Master's and Ph.D. degrees from the Massachusetts Institute of Technology. Since July of 1998 he has been with the Faculty of Computer Science, Dalhousie University, Halifax, Nova Scotia. He served as Director of the Graduate Program (1999-2002) and he is currently Associate Dean, Research. He is a Senior Member of the IEEE. He was a member of the ACM Dissertation Award committee (1990-1992), a member of the AAAI/SIGART Doctoral Consortium Committee (1997-2001) and he is co-editor-in-chief of Computational Intelligence. At Dalhousie, he held a Killam Chair of Computer Science. He has published on the interpretation of visual and range signals for landmark-based navigation and map construction in robotics. He currently works on modelling and mining of content and link structure of Networked Information Spaces.

No photo

Dr. James Blustein is an Associate professor in both Dalhousie's Faculty of Computer Science and School of Information Management. His overall goal is to help people find and use information more effectively. His doctorate in Computer Science was co-supervised by the late Jean Tague-Sutcliffe.