

Term Weighting Evaluation in Bipartite Partitioning for Text Clustering

Chao Qu^{1,2}, Yong Li¹, Jun Zhu^{1,2}, Peican Huang¹,
Ruifen Yuan¹, and Tianming Hu^{1,3}

¹ Dongguan University of Technology, China

² Zhongshan University, China

³ East China Normal University, China

tmhu05@gmail.com

Abstract. To alleviate the problem of high dimensions in text clustering, an alternative to conventional methods is bipartite partitioning, where terms and documents are modeled as vertices on two sides respectively. Term weighting schemes, which assign weights to the edges linking terms and documents, are vital for the final clustering performance. In this paper, we conducted a comprehensive evaluation of six variants of *tf/idf* factor as term weighting schemes in bipartite partitioning. With various external validation measures, we found *tfidf* most effective in our experiments. Besides, our experimental results also indicated that *df* factor generally leads to better performance than *tf* factor at moderate partitioning size.

1 Introduction

The high dimension in text clustering is a major difficulty for most probabilistic methods such as Naive Bayes [1] and AutoClass [2]. To circumvent this problem, graph-theoretic techniques have been considered for clustering [3]. They model the document similarity by a graph whose vertices correspond to documents and weighted edges give the similarity between vertices. Graphs can also model terms as vertices and similarity between terms is based on documents in which they co-occur. Partitioning the graph yields a clustering of terms, which is assumed to be associated with similar concepts [4]. The duality between document and term clustering can also be naturally modeled using a bipartite, where documents and terms are modeled as vertices on two sides respectively and only edges linking different types of vertices are allowed in the graph [5]. Finding an optimal partitioning in such a bipartite gives a co-clustering of documents and terms, with the expectation that documents and terms in the same cluster are related to the same topic. In addition, bipartite graphs have also been used to model other relationships, such as (documents, concepts)[6] and (authors, publications)[7].

Since the general partitioning goal is to minimize the edge cut, term weighting schemes, which assign weights to edges linking term and document vertices, is a vital step to the final clustering performance. Many researchers have studied text clustering based on different term weighting schemes or different criterion

functions using conventional probabilistic methods [8,9,10,11]. For instance, the authors of [9] pointed out that it is the text representation schemes that dominate the clustering performance rather than the kernel functions of support vector machines (SVM). In other words, choosing an appropriate term weighting scheme is more important than choosing and tuning kernel functions of SVM for text categorization. However, to the best of our knowledge, there is little work on comparing weighting schemes for graph based text clustering, not to mention bipartite partitioning. For this purpose, we concentrate on the comparison of various term weighting schemes in bipartite based text clustering. Specifically, we provide a comprehensive experimental evaluation with real world document datasets from various sources, and with various external validation measures.

Overview. The rest of this paper is organized as follows. Section 2 describes the necessary background. Section 3 provides the results of extensive experiments. Finally, we draw conclusions in Section 4.

2 Background

In this section, we first describe the overview of bipartite based text clustering. Then we introduce the term weighting schemes that will be compared in this paper.

2.1 Bipartite Generation

To apply clustering algorithms, a document data set is usually represented by a matrix. First we extract from documents unique content-bearing words as features, which involves removing stopwords and those with extreme document frequencies. More sophisticated techniques use support or entropy to filter words further. Then each document is represented as a vector in this feature space. With rows for documents and columns for terms, the matrix A 's non-zero entry A_{ij} indicates the presence of term w_j in document d_i , while a zero entry indicates an absence.

A graph $G = (V, E)$ is composed of a vertex set $V = \{1, 2, \dots, |V|\}$ and an edge set $\{(i, j)\}$ each with edge weight E_{ij} . The graph can be stored in an adjacency matrix M , with entry $M_{ij} = E_{ij}$ if there is an edge (i, j) , $M_{ij} = 0$ otherwise. Given the $n \times m$ document-term matrix A , the bipartite graph $G = (V, E)$ is constructed as follows. First we order the vertices such that the first m vertices index the terms while the last n index the documents, so $V = V_W \cup V_D$, where V_W contains m vertices each for a term, and V_D contains n vertices each for a document. Edge set E only contains edges linking different kinds of vertices, so the adjacency matrix M may be written as $\begin{pmatrix} 0, A \\ A^T, 0 \end{pmatrix}$.

2.2 Bipartite Partitioning

Given a weighted graph $G = \{V, E\}$ with adjacency matrix M , clustering the graph into K parts means partitioning V into K disjoint clusters of vertices

V_1, V_2, \dots, V_K , by cutting the edges linking vertices in different parts. The general goal is to minimize the sum of the weights of those cut edges. Formally, the cut between two vertex groups V_1 and V_2 is defined as $cut(V_1, V_2) = \sum_{i \in V_1, j \in V_2} M_{ij}$. Thus the goal can be expressed as $\min_{\{V_1, V_2, \dots, V_K\}} \sum_{k=1}^K cut(V_k, V - V_k)$. To avoid trivial partitions, often the constraint is imposed that each part should be roughly balanced in terms of part weight $wgt(V_k)$, which is often defined as sum of its vertex weight. That is, $wgt(V_k) = \sum_{i \in V_k} wgt(i)$. The objective function to minimize becomes $\sum_{k=1}^K \frac{cut(V_k, V - V_k)}{wgt(V_k)}$. Given two different partitionings with the same cut value, the above objective function value is smaller for the more balanced partitioning.

In practice, different optimization criteria have been defined with different vertex weights. The ratio cut criterion [12], used for circuit partitioning, defines $wgt(i) = 1$ for all vertices i and favors equal sized clusters. The normalized cut criterion [13], used for image segmentation, defines $wgt(i) = \sum_j M_{ij}$. It favors clusters with equal sums of vertex degrees, where vertex degree refers to the sum of weights of edges incident on it.

Finding a globally optimal solution to such a graph partitioning problem is in general NP-complete [14], though different approaches have been developed for good solutions in practice [15,16]. Here we employ Graclus [16], a fast kernel based multilevel algorithm, which involves coarsening, initial partitioning and refinement phases. As for the graph partitioning criterion used in Graclus, we tried both the normalized cut criterion and the ratio cut criterion. We found the former always produces better results, possibly for the following reasons. First, our datasets are highly imbalanced, which makes unreasonable the constraint of equal sized clusters by the ratio cut criterion. Second we find that sometimes it yields clusters of pure word vertices, which makes it impossible to determine the number of document clusters (clusters containing the document vertices) beforehand. Those terms with low frequencies are likely to be isolated together, since few edges linking outside are cut. As for the normalized cut criterion that tries to balance sums of vertex degrees in each cluster, the resultant clusters tend to contain both document and term vertices. So in this paper we only report results from the normalized cut criterion.

2.3 Term Weighting

Term weighting schemes determine the value of non-zero entry A_{ij} in the document-term matrix when term w_j appears in document d_i . Two frequencies are commonly used. Term frequency tf_{ij} denotes the raw frequency of term w_j in document d_i . Inverse document frequency $idf_j = \log(n/n_j)$ considers the discriminating power of w_j , where n_j is the number of documents that contain w_j . In this paper, we compared six term weighting schemes listed in Table 1. Most of these term weighting schemes have been widely used in information retrieval and text categorization. The first four term weighting schemes are different variants of tf factor. The last two incorporate idf . According to a recent

Table 1. Term weighting schemes

name	description
<i>binary</i>	1 for presence and 0 for absence
<i>tf</i>	raw <i>tf</i>
<i>logtf</i>	$\log(1 + tf)$
<i>itf</i>	$1 - 1/(1 + tf)$
<i>idf</i>	$\log(n/n_j)$
<i>tfidf</i>	$tf \times idf$

study of text classification with SVM [11], although the first four schemes relate with term frequency alone, all of them show competitive performance with other sophisticated schemes except *binary*. The *idf* factor, taking the collection distribution into consideration, does not improve the terms discriminating power with SVM. The *tfidf* factor, combining both term and document frequencies, usually yields best results in query based document retrieval [17].

3 Experimental Evaluation

In this section, we present an extensive experimental evaluation of various term weighting schemes. First we introduce the experimental datasets and cluster validation criteria, then we report comparative results.

3.1 Experimental Datasets

For evaluation, we selected 10 real data sets from different domains used in [10]. The RE0 and RE1 datasets are from the Reuters-21578 text categorization test collection Distribution 1.0. The datasets K1a, K1b and WAP are from the WebACE project; each document corresponds to a web page listed in the subject hierarchy of Yahoo. In particular, K1a and K1b contain the same data but K1a's class labels are at a finer level. The datasets TR31 and TR41 were derived from the TREC collection. The LA1 and LA2 datasets were obtained from articles of the Los Angeles Times that was used in TREC-5. The FBIS dataset is from the Foreign Broadcast Information Service data of TREC-5. For all data sets, we used a stoplist to remove common words, stemmed the remaining words using Porter's suffix-stripping algorithm, and discard those with very low document frequencies. Some characteristics of them are shown in Table 2.

3.2 Validation Measures

Because the true class labels of documents are known, we can measure the quality of the clustering solutions using external criteria that measure the discrepancy between the structure defined by a clustering and what is defined by the class labels. First we compute the confusion matrix C with entry C_{ij} as the number of documents from true class j that are assigned to cluster i . Then we calculate the following four measures: normalized mutual information(NMI), conditional entropy(CE), error rate(ERR) and F-measure.

Table 2. Characteristics of data sets

data	re0	re1	k1a	k1b	wap	tr31	tr41	la1	la2	fbis
#doc	1504	1657	2340	2340	1560	927	878	3204	3075	2463
#word	2886	3758	4707	4707	8460	10128	7454	6188	6060	2000
#class	13	25	20	6	20	7	10	6	6	17
MinClass	11	13	9	60	5	2	9	273	248	38
MaxClass	608	371	494	1389	341	352	243	943	905	506
min/max	0.018	0.035	0.018	0.043	0.015	0.006	0.037	0.29	0.274	0.075
source	Reuters-21578		WebACE			TREC-6,7		TREC-5		

NMI and CE are entropy based measures. The cluster label can be regarded as a random variable with the probability interpreted as the fraction of data in that cluster. Let T and C denote the random variables corresponding to the true class and the cluster label, respectively. The two entropy-based measures are defined as $NMI = \frac{H(T)+H(C)-H(T,C)}{\sqrt{H(T)H(C)}}$, $CE = H(T|C) = H(T, C) - H(C)$, where $H(X)$ denotes the entropy of X and \log_2 is used here in computing entropy. NMI measures the shared information between T and C and it reaches the maximal value of 1 when they are the same. CE tells the information remained in T after knowing C and it reaches the minimal value of 0 when they are identical. Error rate $ERR(T|C)$ computes the fraction of misclassified data when all data in each cluster is classified as the majority class in that cluster. It can be regarded as a simplified version of $H(T|C)$.

F-measure combines the precision and recall concepts from information retrieval [17]. We treat each cluster as if it were the result of a query and each class as if it were the desired set of documents for a query. We then calculate the recall and precision of that cluster for each given class as follows: $R_{ij} = C_{ij}/C_{+j}$, $P_{ij} = C_{ij}/C_{i+}$, where C_{+j}/C_{i+} is the sum of j th column/ i -th row, i.e., j -th class size / i -th cluster size. Note that C_{+j} could be larger than the true size of class j if some documents from it appear in more than one cluster. F-measure of cluster i and class j is then given by $F_{ij} = \frac{2R_{ij}P_{ij}}{P_{ij}+R_{ij}}$. The overall value for the F-measure is a weighted average for each class, $F = \frac{1}{n} \sum_j C_{+j} \max_i \{F_{ij}\}$, where n is the total sum of all elements of matrix C . F-measure reaches its maximal value of 1 when the clustering is the same as the true classification.

3.3 Clustering Results

Let c in cK denote the number of partitions we set. The six term weighting schemes are evaluated in terms of the four validation measures at $5K$, $10K$, $15K$ and $20K$. The detailed results are shown in Table 3. NMI and F are preferred large while ERR and CE are preferred small. For each setting, the best results are highlighted in bold numbers. One can see in most cases $tfidf$ gives the best results and $binary$ performs worst. Although tf , $logtf$ and itf perform best in certain settings, their gap between $tfidf$ is not significant. The last six rows of Table 3 give the number of wins over all measures at four levels of clustering granularity, respectively. The superiority of $tfidf$ is obvious at three levels.

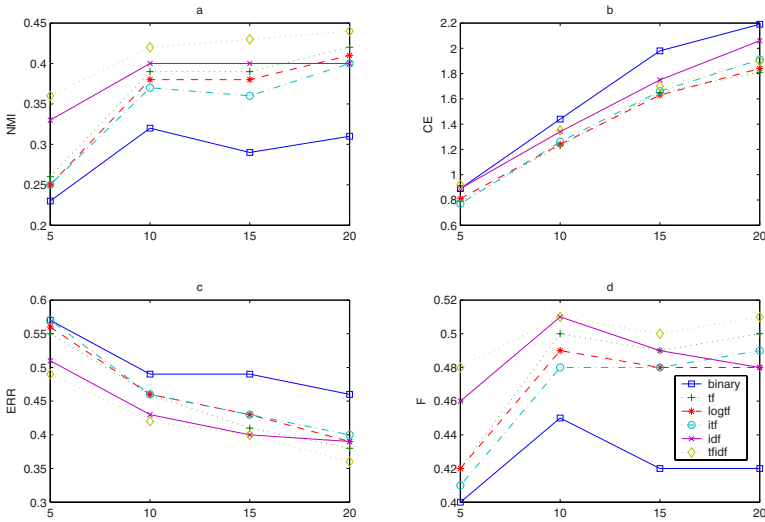


Fig. 1. The average results over the 10 datasets

To give a summary of performance, Figure 1 illustrates the average results over all datasets for each measure. Apparently, *tfidf* gives the best results on all the measures except *CE*. Recall that *CE* is an un-normalized measure, which makes averaging questionable. This indicates that the incorporation of the discriminating factor, *idf*, really makes a difference. It is confirmed again if we compare the performance by *idf* and *tf* alone. One can see at 5K, 10K, 15K, *idf* beats *tf*. It shows that at relatively moderate clustering granularity, the term relevance for the whole document set is more important than that for each individual document. However, as the number of clusters gets larger (e.g., 20K) and each cluster gets smaller, the term frequency within each document matters.

4 Conclusions

Since the goal of bipartite partitioning is to minimize the edge cut, term weighting schemes are essential for the final clustering performance. In this paper, we provided an extensive comparison of six commonly used schemes. Our experimental results show that *tfidf* generally yields better performance than other term weighting schemes in terms of various external validation measures. Besides, at moderate clustering granularity, *idf* is more important than *tf*. Because the graph partitioning is always subject to the balance constraint, the vertex weighting also plays an important role. For the future work, we plan to investigate the impact of vertex weighting schemes. To capture the full semantics that cannot be represented with single words, another direction is to augment the raw vocabulary with word-sets based on frequent itemsets [18,19] or hypercliques [20,21].

Acknowledgments. This work was partially supported by the Scientific Research Foundation for the Returned Overseas Chinese Scholars, State Education Ministry of China, and the Dongguan Foundation of Scientific and Technological Plans.

References

1. Rennie, J.D., Shih, L., Teevan, J., Karger, D.R.: Tackling the poor assumptions of naive bayes text classifiers. In: ICML, pp. 616–623 (2003)
2. Cheeseman, P., Stutz, J.: Bayesian classification (AutoClass): Theory and results. In: KDD, pp. 153–180 (1996)
3. Strehl, A., Ghosh, J., Mooney, R.: Impact of similarity measures on web-page clustering. In: Proc. AAAI: Workshop of Artificial Intelligence for Web Search, pp. 58–64 (2000)
4. Baker, L.D., McCallum, A.: Distributional clustering of words for text classification. In: SIGIR, pp. 96–103 (1998)
5. Dhillon, I.S.: Co-clustering documents and words using bipartite spectral graph partitioning. In: SIGKDD, pp. 269–274 (2001)
6. Yoo, I., Hu, X., Song, I.Y.: Integration of semantic-based bipartite graph representation and mutual refinement strategy for biomedical literature clustering. In: SIGKDD, pp. 791–796 (2006)
7. Sun, J., Qu, H., Chakrabarti, D., Faloutsos, C.: Relevance search and anomaly detection in bipartite graphs. *ACM SIGKDD Explorations* 7(2), 48–55 (2005)
8. Salton, G., Buckley, C.: Term-weighting approaches in automatic text retrieval. *Information Processing and Management* 24(5), 513–523 (1988)
9. Leopold, E., Kindermann, J.: Text categorization with support vector machines: How to represent texts in input space? *Machine Learning* 46(1-3), 423–444 (2002)
10. Zhao, Y., Karypis, G.: Empirical and theoretical comparisons of selected criterion functions for document clustering. *Machine Learning* 55(3), 311–331 (2004)
11. Lan, M., Tan, C.L., Low, H.B.: Proposing a new term weighting scheme for text categorization. In: AAAI (2006)
12. Hagen, L., Kahng, A.: New spectral methods for ratio cut partitioning and clustering. *IEEE Trans. CAD* 11, 1074–1085 (1992)
13. Shi, J., Malik, J.: Normalized cuts and image segmentation. *IEEE Trans. PAMI* 22, 888–905 (2000)
14. Garey, M.R., Johnson, D.S.: *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman & Company, New York (1979)
15. Karypis, G., Kumar, V.: A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM J. Scientific Computing* 20(1), 359–392 (1998)
16. Dhillon, I.S., Guan, Y., Kulis, B.: A fast kernel-based multilevel algorithm for graph clustering. In: SIGKDD, pp. 629–634 (2005)
17. Baeza-Yates, R., Ribeiro-Neto, B.: *Modern Information Retrieval*. Addison-Wesley, Reading (1999)
18. Agrawal, R., Imielinski, T., Swami, A.N.: Mining association rules between sets of items in large databases. In: SIGMOD, pp. 207–216 (1993)
19. Agrawal, R., Srikant, R.: Fast algorithms for mining association rules. In: VLDB, pp. 487–499 (1994)
20. Xiong, H., Tan, P.N., Kumar, V.: Mining strong affinity association patterns in data sets with skewed support distribution. In: ICDM, pp. 387–394 (2003)
21. Xiong, H., Tan, P.N., Kumar, V.: Hyperclique pattern discovery. *Data Mining and Knowledge Discovery Journal* 13(2), 219–242 (2006)