

Exploiting User Queries for Search Result Clustering

Abdul Wahid, Xiaoying Gao, and Peter Andreae

School of Engineering and Computer Science,
Victoria University of Wellington,
19 Kelburn Parade 6012. Wellington, New Zealand
{abdul.wahid,xgao,pondy}@ecs.vuw.ac.nz
<http://ecs.victoria.ac.nz>

Abstract. Search Result Clustering (SRC) groups the results of a user query in such a way that each cluster represents a set of related results. To be useful to the user, the different cluster should contain the results corresponding to different possible meanings of the user query and the cluster labels should reflect these meanings. However, existing SRC algorithms often ignore the user query and group the results based just on the similarity of search results. This can lead to two problems: *low quality cluster*, where the results within a single cluster are related to different meanings of the query; and *poor cluster labels*, where the label of the cluster does not reflect the query meaning associated with the results in the cluster.

This paper presents a new SRC algorithm called QSC that exploits the user query and uses both syntactic and semantic features of the search results to construct clusters and labels. Experiments show that the query senses are good candidates for the cluster labels and the algorithm can lead to high quality cluster and more semantically meaningful labels than other state-of-the-art algorithms.

Keywords: Web Clustering Engine, Search Result Clustering, Query Senses, Document Clustering.

1 Introduction

The goal of Search Result Clustering is not only to cluster search results but also to provide semantically meaningful cluster labels. A Cluster label is a one-phrase description of all the documents in a cluster enabling users to decide whether to browse the list of documents in a cluster by looking at the cluster label. It is a common practice to use the most common keywords shared by all the documents in a cluster as a cluster label. Documents can have common keywords that might represent either more than one sense or might not represent any sense of the user query. Therefore, cluster labels based on common keywords are not always useful to the user. Also the clusters will be more useful to the user if all the documents in a cluster represent only one particular sense of the user query.

Traditional Search Result Clustering algorithms which ignore the user query are more vulnerable to the problems of *low quality cluster* and *poor cluster labels*. *Low quality cluster* is having documents in a cluster that represent more than one senses of the user query and *poor cluster labels* are cluster labels that do not represent any senses of the user query.

The similarity between two documents is often measured using word frequency. Such similarity measures are regarded as syntactic measures because they only consider counts of words. In order to minimize the problem of *low quality cluster*, this paper uses both syntactic and semantic features (topics) of the documents.

This paper presents a new algorithm Query Sense Clustering (QSC) that exploits the user query and combines semantic and syntactic features of a document for the clustering solution. The paper is organized as follows: section 2 highlights the related work; section 3 discusses the representation and similarity measures of the documents and the query senses; section 4 describes the algorithm; section 5 focuses on the evaluation and analysis of the results and section 6 concludes the paper.

2 Related Work

Search Result Clustering (SRC) methods can be classified into three categories: data-centric, description-aware and description-centric [4].

The data-centric category contains traditional clustering algorithms (hierarchical, partitioning) and the focus is on the clustering process. The Scatter/Gather algorithm [10,19] is the pioneer example of the data-centric category. The main drawback of this category is the *poor cluster labels* which are often generated from the text and are often meaningless.

The description-aware methods carefully select one or more features to construct meaningful cluster labels. Suffix Tree Clustering (STC) [23] was the first algorithm that used suffix trees to build cluster labels and perform clustering on search results. The issue with description-aware methods is that the cluster labeling procedure dominates the clustering process and the overall quality of the clusters is compromised.

The description-centric methods are specialized clustering methods that not only focus on cluster labels but also try to provide quality clusters. Examples in this category include LINGO [17]. Our algorithm QSC also belongs to this category.

3 Representation and Similarity Measure

This work uses query senses to generate initial clusters and then uses a new document similarity measure to refine the initial clusters. The new document similarity measure is based on a new document representation using both syntactic and semantic features (topics). The following subsections introduce the new document representation, the document similarity measure, the query sense

representation and the sense similarity measure. The algorithm is presented in section 4.

3.1 Document Representation and Document Similarity Measure

The traditional bag-of-words model is widely used in document clustering to represent documents in Vector Space. Terms are commonly weighted using the tf-idf weighting scheme [21]. A document d in term-space is represented as

$$Tm(d) = \{tfidf(t_1, d), tfidf(t_2, d), tfidf(t_3, d), \dots, tfidf(t_n, d)\} \quad (1)$$

where n is the total number of terms and $tfidf$ is the tf-idf function defined as

$$tfidf(t, d) = tf(t, d) \times \log \frac{|D|}{df(t)} \quad (2)$$

where $tf(t, d)$ is the frequency of term t in the document d , $|D|$ is the total number of documents and $df(t)$ is the number of documents containing term t . A criticism of this model is that it only uses a syntactic representation of the document and ignores semantic representation of the document. One semantic representation is based on topics representing the subjects or concepts that a document is about. If we can identify all the topics of a documents, then we can represent a document as a vector in topic space with weights for each topic representing the importance of the topic to the document. We propose a new document representation in which a document d containing topics $\tau_1 \dots \tau_m$ in topic-space is represented as

$$Tp(d) = \{w(\tau_1, d), w(\tau_2, d), w(\tau_3, d), \dots, w(\tau_m, d)\} \quad (3)$$

where m is total number of topics and $w(\tau, d)$ is a weight of a topic τ , generated using topic detector of Wikiminer Toolkit¹ [15], in document d .

The most common and well known similarity measure for comparing documents is cosine similarity function [18]. We define the combined cosine similarity that includes semantic and syntactic features of document d_i and d_j as

$$Sim(d_i, d_j) = \lambda Cosine(Tp(d_i), Tp(d_j)) + (1 - \lambda) Cosine(Tm(d_i), Tm(d_j)) \quad (4)$$

where λ is a scaling variable and the value of λ is 0.1 based on the preliminary experiments, $Tp(d)$ is document vector in topic-space and $Tm(d)$ is document vector in term space.

3.2 Query Sense Representation and Sense Similarity Measure

We represent a query using a set of senses $S = \{s_1, s_2, s_3 \dots s_n\}$ of the query which is generated using Wikiminer² [15] word disambiguation. These raw senses are

¹ The topic detector is comparable to state-of-the-art LDA based topic detectors.

² Wikiminer parses Wikipedia disambiguation pages to get different senses of a word.

filtered and noise is removed by using tokenization, stemming and stop word removal techniques. Tokens generated from these senses are mostly bi-grams such as *jaguar car*, *sepecat jaguar*, *fender jaguar*, *mac os*. Other examples of senses are *panthera* and *south alabama jaguar football*.

We define the similarity score between a document d_i and a sense s_j as a weighted sum of six different criteria:

$$SimSense(d_i, s_j) = \frac{|s_j|}{|d_i|} \sum_{k=1}^6 w_k \cdot cmp_k(d_i, s_j) \quad (5)$$

The six criteria for cmp are exact sequence matching, semantic matching, partial matching in both term space and topic space of the document d_i for sense s_j . The exact sequence matching counts the number of occurrence; the semantic matching counts overlap of either exact or synonyms; and partial matching counts overlap of individual words in sense s_j and document d_i .

4 QSC Algorithm

We had developed a new algorithm called QSC³ that uses our new document representation and similarity measures. It includes three main steps: the first step is to group all the documents according to their similarity to the different senses of the user query; the second step is to iteratively optimize clusters by relocating documents from one cluster to another cluster based on the similarity between documents and the clusters; the third step is to rank the documents and clusters based on similarity with the user query.

Step 1: Initial Cluster Generation. The initial clusters are formed by calculating the similarity of each document with each user query sense and assigning each document to each cluster associated with the maximally similar sense. Each cluster is labeled with its associated sense. Documents that are not sufficiently similar to any sense are placed in a cluster labeled *general*. The set of initial clusters C consists of all the clusters that contains at least one document.

Step 2: Cluster Optimization. Initial clusters were based on the similarity between documents and the senses. Base cluster labels can provide quality labeling of clusters. However the clusters, especially the general cluster may contain a mixed group of documents that might not be similar. We developed an iterative method to reassign some documents in order to improve cluster quality by increasing intra-cluster coherence and inter-cluster distinctiveness.

Step 3: Cluster Ranking. Users are interested in only those documents that are most closely related to the query. Therefore the ranking of clusters and documents are computed with respect to the query.

³ The full paper contains the pseudo code with more details and is available on request.

All the clusters were sorted, by calculating the relatedness score between the user query and the cluster label, using the term similarity measure *WikiSim* [12]. The *WikiSim* is Wikipedia based similarity measure that computes relatedness between two terms. Documents in its own cluster were also sorted by calculating the similarity of a document to its mean in its own cluster. The ranked result list is then sent to user for browsing.

5 Results

The QSC was evaluated on two datasets, AMBIENT [5] and MORESQUE [16].

Comparison 1: The results on the larger dataset, which consist of all queries of AMBIENT and MORESQUE, based on purity and entropy were not given in [14]. However we found another recent paper [8] that compared nine algorithms using F1-measure on this large dataset. Therefore we compared our algorithm QSC with these nine algorithms using F1-measure calculated by taking the harmonic mean of precision and recall of the cluster [6]. The comparison was made between STC, LINGO, KeySRC [1], Curvature [9], SquaT++ [16,8], B-MST [7], HyperLex [22], Chinese Whispers [2] and QSC.

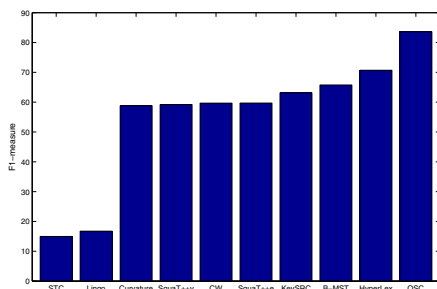


Fig. 1. Comparison of SRC methods

Figure 1 shows the percentage values of F1-measure of 10 methods on combined dataset of AMBIENT and MORESQUE taken from the paper [8] and the computed value of QSC. Clearly the QSC performed significantly better than others and have the highest value 83.62 (percentage) of F1-measure.

Other evaluation criteria Adjusted Rand Index(ARI) and Jaccard Index(JI) are also used for comparing the clustering algorithms in paper [8]. However we believe that they are not suitable for these two datasets. More discussion is provided in the last part of section 5.4.

Comparison 2: The search results needs to be diverse and top ranked results should represent different senses of the user query. In order to determine the diversification of this work, the search results were evaluated based on **S-recall@K** (Subtopic recall at rank K) and **S-precision@r** (Subtopic precision at

recall) [24] on combined dataset of AMBIENT and MORESQUE. The former evaluates the performance of the system based on K top-ranked results for number of topics of query q . S-precision@ r measures the ratio of subtopics covered by minimum set of results at given recall r .

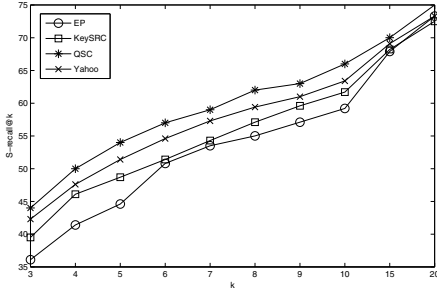


Fig. 2. S-recall@ k on all queries

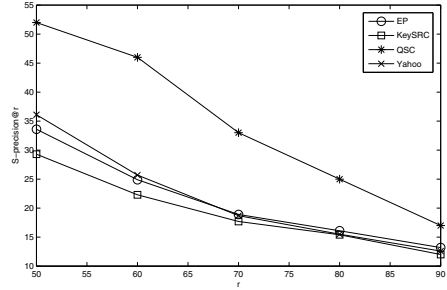


Fig. 3. S-precision@ r on all queries

These two measures are used to compare search engines (Yahoo! and Essential Pages) that return ranked list of search results. The results returned by QSC were compared by flattening the clusters. The result list was formed by iterating through clusters and selecting top results. The clusters that only had one document were appended at the end to avoid noise.

Figure 2 and Figure 3 shows the S-recall@ k and S-precision@ r respectively for search results of Yahoo, Essential Pages(EP), KeySRC and QSC. The QSC performs relatively better in terms of S-recall@ k and significantly outperformed others in terms of S-precision@ r for the given values of k and r . This shows that QSC produced more diverse results than currently available search engines.

5.1 Further Analysis

The detailed analysis consists of three sub sections: the first discusses the cluster labels; the second discusses the processing time of the QSC, and the third discusses the cluster numbers and some observations about final clusters.

Cluster Label Analysis: The goal of the QSC algorithm is to generate a useful set of distinct clusters with informative labels.

Table 1 shows the cluster labels of the clusters generated for the query *Jaguar* by STC, LINGO and QSC (the cluster labels are not in ranked order). The labels for STC and LINGO were generated using the Carrot2 framework by adjusting the parameter of maximum clusters number to 8. Table 1 shows that the cluster labels generated by QSC provides more precise, intuitive and distinct labels than the cluster labels from STC and LINGO.

Processing Time: The QSC was evaluated on standalone workstation using Linux (64 bit) with Intel(R) Core(TM) i7-3770 CPU @ 3.40GHZ, 8GB RAM and

Table 1. Cluster labels of STC, LINGO and QSC of the query Jaguar

STC	LINGO	QSC
Jaguar Car	Auto Show	Jaguar Car
S-Type, Used Jaguar	Jaguar Parts	Jaguar E-Type
XK, 2006 2007, Price- Quotes and Reviews	Dealer Price Quotes and- Reviews	Jaguar XK
Ford Motor Company- Division	Ford Motor Company- Division	
Jaguar Cars		
Jaguar Panthera Onca	Jaguar Panthera Onca	Panthera
Jaguar Animal		
	Website of Fender Musical- Instrument	Fender Jaguar
Information	Jaguar Video	Mac OS X
New		SEPECAT Jaguar South Alabama Jaguar- Football

1TB HD. Figure 4 shows the processing time of all the queries in AMBIENT Dataset. The average time required for processing the query is under 1.0 second for both AMBIENT and MORESQUE datasets. Most of the queries were processed under one second with few exceptional cases. The maximum processing time was 6.3 seconds on a query *jaguar* because it had 54 senses to be processed. This processing time was reduced to 1 second by eliminating overlapping senses and processing only 10 distinct senses.

Strictly speaking, we cannot directly compare the processing time of other algorithms due to different machines and platforms. However we would like to give indications that word sense induction based algorithms (Curvature, Squat++, B-MST, HyperLex and Chinese Whispers) need to construct the graph to identify the senses from the huge corpus, whereas QSC extract the senses from the Wikipedia. Therefore the word sense induction based algorithms might require more processing time than QSC. The processing time of clustering, without considering the time spent on graph construction, for all algorithms is under 1 second except for SquaT++ algorithm. The SquaT++v and SquaT++e spent around 28 and 21 seconds respectively for clustering results as described in their paper [8].

Cluster Analysis: The average number of clusters for all queries in the AMBIENT dataset was 7.84 i.e on average 7-8 clusters are formed for each query. The average number of clusters for all queries in the AMBIENT and MORESQUE datasets was 5.4. There were a few queries with a high number of clusters and the maximum number of clusters was 18 for the query *Monte Carlo*. In contrast the

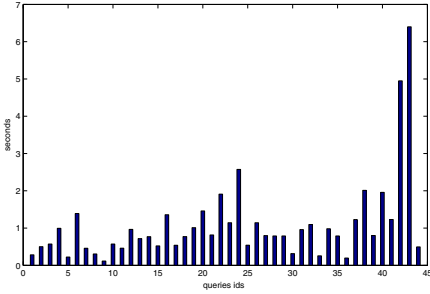


Fig. 4. Processing Time for All Queries

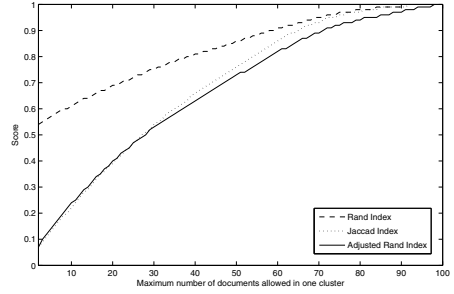


Fig. 5. RI, ARI and JI Analysis

query *Life on Mars* just had 1 big cluster. The reason for having many clusters was the large number of distinct query senses. The query *Life on Mars* had very few senses and they were overlapping with each other, e.g. *Life on Mars (TV series)*, *Life on Mars (U.S TV series)*, that causes single cluster for the query.

The QSC provided a more fine-grained clustering solution than the gold standard (manually labeled search results). The gold standard for the query *jaguar* had 7 clusters but QSC solution provided 9 clusters. The three clusters *jaguar car*, *jaguar e-type* and *jaguar xk* in QSC were sub clusters of gold standard *jaguar car*.

The QSC was not compared with other algorithms using index based evaluation measures (ARI and JI) because these measures have many issues [20,13]. One of the problems is that they do not handle fine-grained clustering solutions. If a gold standard G has a cluster g_i that contains 90 documents and clustering C has clusters c_j, c_{j+1} and c_{j+2} that contain all 90 documents then ARI and JI will penalize the clustering solution heavily. However the fine-grain clustering solution is consistent with the coarser solution and should not be penalized heavily. In fact it may even be better solution because it provides the distinctiveness that are not provided by the gold standard. ARI and JI do not measure this.

Figure 5 shows the phenomena of heavy penalty of ARI and JI as compared to Rand Index (RI) [11] on sub clusters. This experiment was performed on the AMBIENT and MORESQUE datasets by evaluating the perfect sub-clusters that gradually increased the limit of the maximum number of documents allowed in a cluster from 2 to 98. All the documents were perfectly assigned to the clusters and the values of RI, ARI and JI were computed at each iteration. The lowest value of RI, ARI and JI were 0.54, 0.05 and 0.04 respectively when the maximum allowed number of documents in sub clusters were 2. Figure 5 shows ARI and JI penalize small clusters and small sub-clusters heavily. The gold standard in our dataset had very unbalanced number of clusters. A few clusters were very small, and had 2 documents in a cluster and other were very large and had more than 90 documents. It was observed that the comparison based on ARI and JI is suitable only when the gold standard do not have sub clusters and all the clusters have almost the same number of members.

6 Conclusion

This paper presents a new description-centric search result clustering algorithm QSC that exploits query senses to generate meaningful cluster labels and use syntactic and semantic features of documents to generate quality clusters.

This paper shows that QSC outperforms existing algorithms. QSC is computationally inexpensive and provides better quality clusters with meaningful labels as compared to other algorithms, hence it has the potential to be applied to real time search result clustering applications.

The future direction for this work is to use Google WebIT and ukWac corpus along with Wikipedia to enhance the quality of query senses. The similarity measure and documentation representation are the key factors and a better similarity measure could bring more improvement. The greedy search in step 2 of the QSC could be improved to avoid local optima, by using the query senses in addition to document similarity. The currently used topic detection technique is not as good as state-of-the-art topic detection techniques such as LDA [3]; using LDA to detect topics from search results by considering query senses may further improve this work.

References

1. Bernardini, A., Carpineto, C., D'Amico, M.: Full-subtopic retrieval with keyphrase-based search results clustering. In: IEEE/WIC/ACM International Joint Conferences on Web Intelligence and Intelligent Agent Technologies, WI-IAT 2009, vol. 1, pp. 206–213. IET (2009)
2. Biemann, C.: Chinese whispers: an efficient graph clustering algorithm and its application to natural language processing problems. In: Proceedings of the First Workshop on Graph Based Methods for Natural Language Processing, pp. 73–80. Association for Computational Linguistics (2006)
3. Blei, D., Ng, A., Jordan, M.: Latent dirichlet allocation. *The Journal of Machine Learning Research* 3, 993–1022 (2003)
4. Carpineto, C., Osiniski, S., Romano, G., Weiss, D.: A survey of web clustering engines. *ACM Computing Surveys (CSUR)* 41(3), 17 (2009)
5. Carpineto, C., Romano, G.: Ambient dataset (2008)
6. Crabtree, D., Gao, X., Andreae, P.: Improving web clustering by cluster selection. In: Proceedings of the 2005 IEEE/WIC/ACM International Conference on Web Intelligence, pp. 172–178. IEEE (2005)
7. Di Marco, A., Navigli, R.: Clustering web search results with maximum spanning trees. In: Pirrone, R., Sorbello, F. (eds.) *AI*IA 2011. LNCS*, vol. 6934, pp. 201–212. Springer, Heidelberg (2011)
8. Di Marco, A., Navigli, R.: Clustering and diversifying web search results with graph-based word sense induction. *Computational Linguistics*, 1–76 (just accepted, 2013)
9. Dorow, B., Widdows, D., Ling, K., Eckmann, J.-P., Sergi, D., Moses, E.: Using curvature and markov clustering in graphs for lexical acquisition and word sense discrimination. *arXiv preprint cond-mat/0403693* (2004)

10. Hearst, M., Pedersen, J.: Reexamining the cluster hypothesis: scatter/gather on retrieval results. In: Proceedings of the 19th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 76–84. ACM (1996)
11. Hubert, L., Arabie, P.: Comparing partitions. *Journal of Classification* 2(1), 193–218 (1985)
12. Jabeen, S., Gao, X., Andreae, P.: Harnessing wikipedia semantics for computing contextual relatedness. In: Anthony, P., Ishizuka, M., Lukose, D. (eds.) PRICAI 2012. LNCS, vol. 7458, pp. 861–865. Springer, Heidelberg (2012)
13. Meilă, M.: Comparing clusterings—an information based distance. *Journal of Multivariate Analysis* 98(5), 873–895 (2007)
14. Meiyappan, Y., Iyengar, N.C.S.N., Kannan, A., Suyoto, Y.D., Suselo, T., Prasetyaningrum, T., Tlili, R., Slimani, Y., Dufreche, S., Zappi, M., et al.: Scluster: Web clustering engine based on wikipedia. *International Journal of Advanced Science and Technology* 39(1), 1–18 (2012)
15. Milne, D., Witten, I.H.: An open-source toolkit for mining wikipedia. *Artificial Intelligence* (2012)
16. Navigli, R., Crisafulli, G.: Inducing word senses to improve web search result clustering. In: Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing, pp. 116–126. Association for Computational Linguistics (2010)
17. Osiriski, S., Stefanowski, J., Weiss, D.: Lingo: Search results clustering algorithm based on singular value decomposition. In: Intelligent Information Processing and Web Mining: Proceedings of the International IIS: IIPWM 2004 Conference held in Zakopane, Poland, p. 359 (2004)
18. Pang-Ning, T., Steinbach, M., Kumar, V.: Introduction to data mining. WP Co. (2006)
19. Pirolli, P., Schank, P., Hearst, M., Diehl, C.: Scatter/gather browsing communicates the topic structure of a very large text collection. In: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, pp. 213–220. ACM (1996)
20. Rosenberg, A., Hirschberg, J.: V-measure: A conditional entropy-based external cluster evaluation measure. In: Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL), vol. 410, p. 420 (2007)
21. Salton, G., McGill, M.J.: Introduction to modern information retrieval (1986)
22. Véronis, J.: Hyperlex: lexical cartography for information retrieval. *Computer Speech & Language* 18(3), 223–252 (2004)
23. Zamir, O., Etzioni, O., Madani, O., Karp, R.: Fast and intuitive clustering of web documents. In: Proceedings of the 3rd International Conference on Knowledge Discovery and Data Mining, pp. 287–290. MIT Press (1997)
24. Zhai, C.X., Cohen, W.W., Lafferty, J.: Beyond independent relevance: methods and evaluation metrics for subtopic retrieval. In: Proceedings of the 26th Annual International ACM SIGIR Conference on Research and Development in Informaion Retrieval, pp. 10–17. ACM (2003)