

Exploiting Internal and External Semantics for the Clustering of Short Texts Using World Knowledge

Xia Hu,^{1,2} Nan Sun,¹ Chao Zhang,¹ Tat-Seng Chua¹

¹School of Computing
National University of Singapore

²School of Computer Science and Engineering
BeiHang University

November 2, 2009

Improve the Clustering of Short Texts

Xia Hu

outline

Introduction

Proposed
Framework

Evaluation

Conclusion
and Future
Work

- 1 Introduction
- 2 Proposed Framework
- 3 Evaluation
- 4 Conclusion and Future Work

Improve the Clustering of Short Texts

Xia Hu

outline

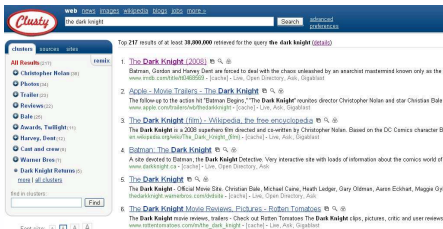
Introduction

Proposed Framework

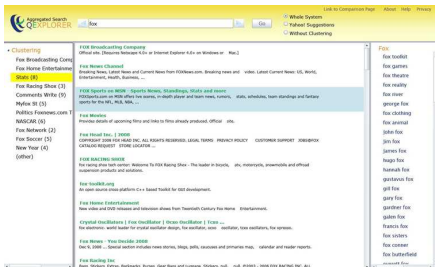
Evaluation

Conclusion and Future Work

- The form of browsing search results.



The screenshot shows the Clusty search engine interface. At the top, there is a search bar with the text "the dark knight" and a search button. Below the search bar, there are navigation links for "web", "news", "images", "wikipedia", "blogs", "jobs", and "more...". The main content area displays "Top 217 results of at least 30,000,000 retrieved for the query the dark knight (Details)". The results are numbered 1 through 6, each with a title and a brief description. For example, result 1 is "The Dark Knight (2008)" with a description about Batman, Gordon, and Harvey Dent. Result 2 is "Apple - Movie Trailers - The Dark Knight" with a description about the actor Val Kilmer. Result 3 is "The Dark Knight (film) - Wikipedia, the free encyclopedia" with a description of the 2008 superhero film. Result 4 is "Batman: The Dark Knight" with a description of a site devoted to Batman. Result 5 is "The Dark Knight" with a description of the official movie site. Result 6 is "The Dark Knight Movie Reviews, Pictures - Rotten Tomatoes" with a description of movie reviews and trailers. At the bottom of the results, there are options for font size and a "Find" button.



The screenshot shows a search engine results page for the query "fox". The search bar at the top contains the text "fox" and has buttons for "Go" and "Web". Below the search bar, there are several search filters: "Whole System", "Fuzzy Suggestions", and "Without Clustering". The main content area is divided into two columns. The left column lists various search results, including "FOX Broadcasting Company", "Fox News Channel", "Fox Racing Sheds", "Fox Sports on NBC", "Fox WinStar", "Fox World Inc.", "FOX RACING SHEDS", "Fox-Web.org", "Fox Video Entertainment", "Crystal Oscillators", and "Fox News - Fox Decides 2008". The right column lists a list of "Fox" related terms, including "fox toonies", "fox games", "fox magazine", "fox reality", "fox rider", "george fox", "fox clothing", "fox journal", "gale fox", "jim fox", "james fox", "hugo fox", "hannah fox", "graciela fox", "gill fox", "gary fox", "garth fox", "golden fox", "francis fox", "fox sisters", "fox coaster", "fox butterfly", and "foxmatt fox".

Improve the Clustering of Short Texts

Xia Hu

outline

Introduction

Proposed
Framework

Evaluation

Conclusion
and Future
Work

- Short texts, such as the snippets, product descriptions, QA passages and image captions etc., have played important roles in current Web and IR applications.
- Unlike standard texts with lots of words in length, short texts, which only consist of a few phrases or 2–3 sentences, especially present great challenges in clustering.
- Problems : “data sparseness” & “semantic gap”.

Improve the Clustering of Short Texts

Xia Hu

outline

Introduction

Proposed Framework

Evaluation

Conclusion and Future Work

- Many methods have been proposed to improve the representation of standard text for clustering and classification, including “surface representation” [3,19] and “integrating world knowledge” [14].
- Several clustering techniques were employed to place the search engine snippets to their highly relevant topic-coherent groups [5,29].
- World knowledge bases have been found useful in improving the short text representation [1,23].

Improve the
Clustering of
Short Texts

Xia Hu

outline

Introduction

Proposed
Framework

Evaluation

Conclusion
and Future
Work

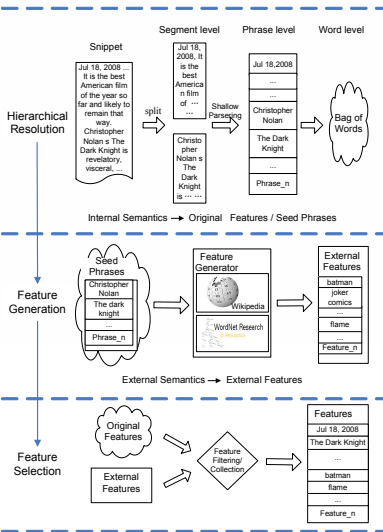


Fig: Framework for feature constructor

“Jul 18, 2008 ... It is the best American film of the year so far and likely to remain that way. Christopher Nolan’s The Dark Knight is revelatory, visceral ...”

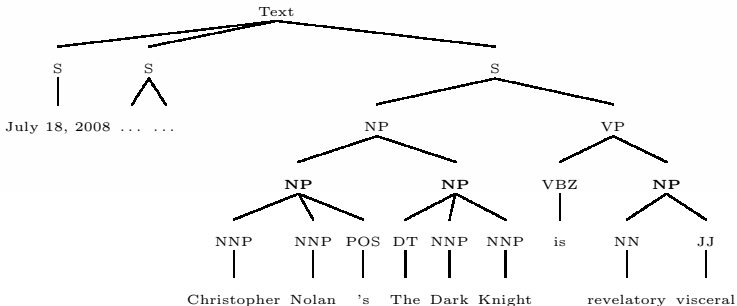


Fig: Syntax tree of the snippet

- Segment-level features.
- Phrase-level features.
 - *Sentence1* : [NP July 18 2008]
 - *Sentence2* : [NP It] [VP is] [NP the best American film] [PP of] [NP the year] [ADVP so far] and/CC [ADJP likely] [VP to remain] [NP that way]
 - *Sentence3* : [NP Christopher Nolan 's] [NP The Dark Knight] [VP is] [NP revelatory visceral]
- Word-level features.

Improve the
Clustering of
Short Texts

Xia Hu

outline

Introduction

Proposed
Framework

Evaluation

Conclusion
and Future
Work

Two steps:

- the construction of basic features
 - seed phrases from internal semantics.
- the generation of external features.
 - external features from world knowledge bases.

- There are redundancies between *phrase level* features and *segment level* features.
- We propose to measure the semantic similarity between the two kinds of feature to eliminate information redundancy.
- For Wikipedia we download the XML corpus, remove xml tags and create a Solr index of all XML articles.

- Let P denotes a *segment level* feature, $P = \{p_1, p_2, \dots, p_n\}$.
- We calculate the semantic similarity between p_i and $\{p_1, p_2, \dots, p_n\}$ as $InfoScore(p_i)$.
- The p^* which has the largest similarity with other features in P will be removed as the redundant feature.

- Given two phrases p_i and p_j , the variants of three popular co-occurrence measures[6] are defined as below:

$$\begin{aligned}
 & \text{WikiDice}(p_i, p_j) \\
 = & \begin{cases} 0 & \text{if } f(p_i | p_j) = 0 \\ & \text{or } f(p_j | p_i) = 0 \\ \frac{f(p_i|p_j)+f(p_j|p_i)}{f(p_i)+f(p_j)} & \text{otherwise} \end{cases}, \quad (1)
 \end{aligned}$$

where WikiDice is a variant of the Dice coefficient.

$$\begin{aligned}
 & \text{WikiJaccard}(p_i, p_j) \\
 = & \frac{\min(f(p_i | p_j), f(p_j | p_i))}{f(p_i) + f(p_j) - \max(f(p_i | p_j), f(p_j | p_i))}, \quad (2)
 \end{aligned}$$

where WikiJaccard is a variant of the Jaccard coefficient.

$$WikiOverlap(p_i, p_j) = \frac{\min(f(p_i | p_j), f(p_j | p_i))}{\min(f(p_i), f(p_j))}, \quad (3)$$

where WikiOverlap is a variant of the Overlap(Simpson) coefficient.

Linear normalization formula is defined below:

$$WD_{ij} = \frac{WikiDice_{ij} - \min(WikiDice_k)}{\max(WikiDice_k) - \min(WikiDice_k)}, \quad (4)$$

A linear combination is then used to incorporate the three similarity measures into an overall semantic similarity between two phrases p_i and p_j , as follows:

$$WikiSem(p_i, p_j) = (1 - \alpha - \beta)WD_{ij} + \alpha WJ_{ij} + \beta WO_{ij}, \quad (5)$$

where α and β weight the importance of the three similarity measures.

For each *segment level* feature, we rank the information score defined in Equation 5 for its child node features at *phrase level*.

$$InfoScore(p_i) = \sum_{j=1, j \neq i}^n WikiSem(p_i, p_j). \quad (6)$$

Finally, we remove the *phrase level* feature p^* , which delegates the most information duplicate to the *segment level* feature P , and it is defined as:

$$p^* = \arg \max_{p_i \in \{p_1, p_2, \dots, p_n\}} InfoScore(p_i). \quad (7)$$

Improve the Clustering of Short Texts

Xia Hu

outline

Introduction

Proposed
Framework

Evaluation

Conclusion
and Future
Work

- Wikipedia, as background knowledge, has a wider knowledge coverage than WordNet and is regularly updated to reflect recent events.
- On the other hand, as the construction of WordNet follows theoretical model or corpus evidence, it contains rich lexical semantic knowledge.

Algorithm 1: GenerateFeatures(S)

input : a set S of seed phrases

output: external features EF

$EF \leftarrow null$

for seed phrase $s \in S$ **do**

if $s.non-stop > 1$ **then**

if $s \in Segment\ level$ **then**

$s.Query \leftarrow SolrSyntax(s, OR)$

else

$s.Query \leftarrow SolrSyntax(s, AND)$

$WikiPages \leftarrow Retrive(s.Query)$

$EF \leftarrow EF + Analyze(WikiPages)$

else

$EF \leftarrow EF + WordNet.Synsets(s)$

return EF

Fig: External feature generation scheme

Feature filtering for unstructured features:

- Remove features generated from too general *seed phrase* that returns a large number (more than 10,000) of articles from the index corpus.
- Transform features used for Wikipedia management or administration, e.g. “List of hotels” → “hotels”, “List of twins” → “twins”.
- Apply phrase sense stemming using Porter stemmer[24], e.g. “fictional books” → “fiction book”.
- Remove features related to chronology, e.g. “year”, “decade” and “centuries”.

To avoid “curse of dimensionality”:

- The number of *external features* we need to collect is determined by:

$$n_2 = \frac{n_1 \times \theta}{1 - \theta}. \quad (8)$$

- Select one external feature for each *seed phrase*.

$$f_i^* = \arg \max_{f_{ij} \in \{p_{i1}, p_{i2}, \dots, p_{ik}\}} tf_idf(f_{ij}). \quad (9)$$

- The top $n_2 - m$ features are extracted from the remaining *external features* based on their frequency.

Reuters-21578 :

- We remove the texts which contain more than 50 words and filter those clusters with less than 5 texts or more than 500 texts.
- Thus it leaves 19 clusters comprising 879 texts. The number of texts in each cluster ranges from 6 (the cluster “income”) to 438 (the cluster “acq”).

Web Dataset is built to simulate a real web application.

- As the users' interests are varied, we choose queries of different length according to the statistics of Google Trends during Nov. 26th 2007 – Nov. 25th 2008.

query length	One	Two	Three	more
count	4552	19762	6992	5290
percentage	12.4%	54.0%	19.1%	14.5%

- Ten hot queries are selected.

Tab: The selected hot queries in Web Dataset

NFL	Amazing Grace
Green Bay	Fox News Channel
60 Minutes	New York Giants
Total Eclipse	The Dark Knight
Black Friday	National Economic Council

K-means and *EM* are employed in this study. Six different text representation methods, as defined below:

- *BOW* (baseline 1) : Traditional “bag of words” model with the *tf-idf* weighting schema.
- *BOW+WN* (baseline 2) : *BOW* integrated with additional features from WordNet as presented in [14].
- *BOW+Wiki* (baseline 3) : *BOW* integrated with additional features from Wikipedia as presented in [1].
- *BOW+Know* (baseline 4) : *BOW* integrated with additional features from Wikipedia and WordNet as in baselines 2 and 3.
- *BOF* : The bag of *original features* extracted with the hierarchical view.
- *SemKnow* : Our proposed framework.

We evaluate performance of the methods using F_1 measure and *Average Accuracy*.

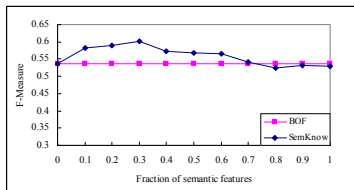
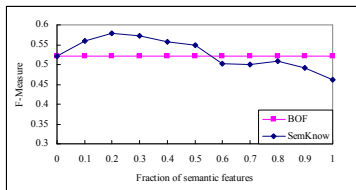
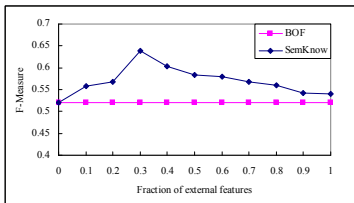
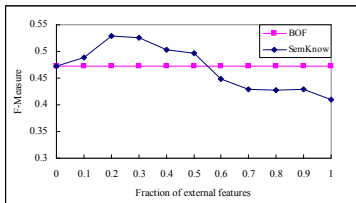
Tab: Results using *k-means* algorithm

	Reuters-21578		Web Dataset	
	F_1 measure (Impr)	AveAccuracy (Impr)	F_1 measure (Impr)	AveAccuracy (Impr)
<i>BOW</i>	0.471 (N.A.)	0.550 (N.A.)	0.491 (N.A.)	0.563 (N.A.)
<i>BOW + WN</i>	0.473 (+0.43%)	0.552 (+0.26%)	0.530 (+8.01%)	0.576 (+2.30%)
<i>BOW + Wiki</i>	0.481 (+2.03%)	0.563 (+2.18%)	0.556 (+13.38%)	0.584 (+3.85%)
<i>BOW + Know</i>	0.489 (+3.75%)	0.566 (+2.86%)	0.558 (+13.79%)	0.583 (+3.70%)
<i>BOF</i>	0.473 (+0.33%)	0.551 (+0.19%)	0.520 (+5.95%)	0.570 (+1.24%)
<i>SemKnow</i>	0.497 (+5.41%)	0.572 (+3.98%)	0.583(+18.81%)	0.586(+4.11%)

Tab: Results using *EM* algorithm

	Reuters-21578		Web Dataset	
	F_1 measure (Impr)	AveAccuracy (Impr)	F_1 measure (Impr)	AveAccuracy (Impr)
<i>BOW</i>	0.516 (N.A.)	0.579 (N.A.)	0.521 (N.A.)	0.608 (N.A.)
<i>BOW + WN</i>	0.525 (+1.72%)	0.585 (+0.99%)	0.540 (+3.59%)	0.626 (+3.02%)
<i>BOW + Wiki</i>	0.540 (+4.74%)	0.598 (+3.39%)	0.550 (+5.50%)	0.629 (+3.44%)
<i>BOW + Know</i>	0.542 (+5.13%)	0.607 (+4.54%)	0.556 (+6.74%)	0.635 (+4.41%)
<i>BOF</i>	0.520 (+0.82%)	0.594 (+2.63%)	0.536 (+2.73%)	0.624 (+2.55%)
<i>SemKnow</i>	0.548 (+6.28%)	0.622 (+7.51%)	0.569 (+9.07%)	0.670 (+10.20%)

Impact of the parameter θ on Reuters and Web Dataset using K – means and EM respectively.



Tab: Optimal results using two algorithms

<i>kmeans</i>	Reuters	F_1 meas(Impr)	AveAcc(Impr)
	<i>BOW</i>	0.471(N.A.)	0.550(N.A.)
	<i>Optimal</i>	0.530(+12.35%)	0.604(+9.72%)
	Webdata		
<i>EM</i>	<i>BOW</i>	0.491(N.A.)	0.563(N.A.)
	<i>Optimal</i>	0.640(+30.39%)	0.607(+7.83%)
	Reuters	F_1 meas(Impr)	AveAcc(Impr)
	<i>BOW</i>	0.516(N.A.)	0.579(N.A.)
<i>EM</i>	<i>Optimal</i>	0.578(+12.02%)	0.672(+15.40%)
	Webdata		
	<i>BOW</i>	0.521(N.A.)	0.608(N.A.)
	<i>Optimal</i>	0.602(+16.14%)	0.709(+16.56%)

- In this study, we proposed a novel framework to augment the clustering accuracy of short texts by exploiting the internal and external semantics.
- The combination of internal and external semantics well tackled the problems of data sparseness and semantic gap in short texts.
- Empirical evaluations demonstrated that our framework significantly outperformed all the baselines including previously proposed knowledge-based short text clustering methods on two datasets.

Improve the Clustering of Short Texts

Xia Hu

outline

Introduction

Proposed
Framework

Evaluation

Conclusion
and Future
Work

- As this work is for aggregated search, the efficiency of the whole framework should be optimized for real applications.
- Moreover, we will explore more tasks in NLP and information retrieval using the internal and external semantics generated by our proposed framework.

Improve the
Clustering of
Short Texts

Xia Hu

outline

Introduction

Proposed
Framework

Evaluation

Conclusion
and Future
Work

Thank you!