

ABSTRACT:

Contributions:

- (1) Introduce a topic-oriented query expansion model based on the Information Bottleneck theory
- (2) Define a term-term similarity matrix
- (3) Propose two measures, intracluster and intercluster similarities

Purposes:

- In (1), we classify terms into distinct topical clusters in order to find out candidate terms for the query expansion
- With (2), we are available to improve the term ambiguous problem
- Two measures in (3) are based on proximity between the topics represented by two clusters in order to evaluate the retrieval effectiveness

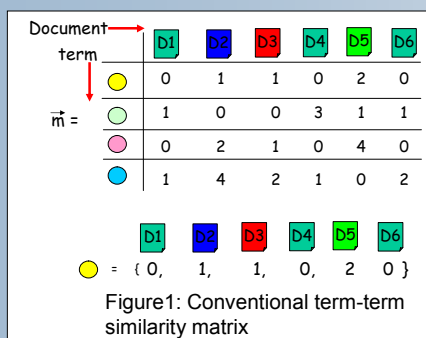
INTRODUCTION:

Term ambiguous example:

The original text mining algorithm was created by a university student named Taylor. In our interview, we found out that he is currently in his second year as an undergraduate student of the Mining Engineering Department of this university and has many other interests beside of programming.

Remarks:

- (1) The first "mining" and the second "Mining" describe two different senses.
- (2) Co-occurrence terms are different, i.e., "algorithm" and "Engineering".
- (3) However, in the conventional term-term similarity matrix, "mining" is counted twice times (Figure 1) and considered as one concept.



Remarks:

Term $mining$ have different senses in documents d_2 , d_3 and d_5 . However, in the conventional term-term similarity matrix, we only considered $mining$ as one concept and count its numbers of amount in the documents d_2 , d_3 and d_5 . Using this type of similarity matrix to look for the co-occurrence words of the original query will lead to the problem of term ambiguous.

RELATED WORKS:

Information Bottleneck Theory [5]

A compressed variable is X and its relevant variable is Y . The purpose is to search a middle variable T (clusters of X) in which the mutual information between T and Y , i.e., $I(T, Y)$ is maximized while the mutual information $I(X, T)$ is minimized.

$$\mathcal{L}[P(T|X)] = I(X, T) - \beta I(T, Y) \quad (1)$$

three self-consistent equations

$$\begin{cases} p(t|x) = \frac{p(t)}{\sum_x p(t|x)} e^{-\beta D_{KL}(p(y|x)||p(y|t))} \\ p(y|t) = \frac{1}{p(t)} \sum_x p(t|x)p(x)p(y|x) \\ p(t) = \sum_x p(t|x)p(x) \end{cases} \quad (2)$$

TOPIC-ORIENTED QUERY EXPANSION MODEL:

2. Topic-oriented query expansion model:

- (1) to collect pages that come from incoming links.
- (2) to extract extended anchor texts.
- (3) to construct a term-term similarity matrix.
- (4) to establish a topic-oriented cluster.
- (5) to rank terms.

2.1 Constructing a term-term similarity matrix

Refer to [1], we redefine the term-term similarity matrix

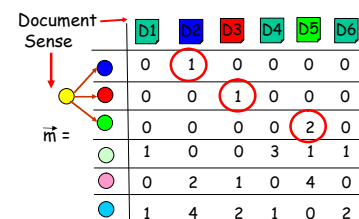


Figure 3: Our proposed term-term similarity matrix

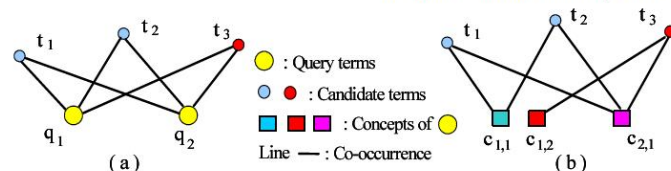


Figure 4: Correlation among candidate terms, query terms, and query concepts

Example:

Query $q = \{q_1, q_2\}$
 Document $d_1: d_1 = \{t_1, q_1, t_2\}$
 Document $d_2: d_2 = \{t_1, t_2, q_2, t_3\}$
 Document $d_3: d_3 = \{q_1, t_3\}$
 Document $d_4: d_4 = \{t_2, t_3\}$

Using the conventional similarity matrix

$$m = \begin{matrix} & d_1 & d_2 & d_3 & d_4 \\ t_1 & 1 & 1 & 0 & 0 \\ t_2 & 1 & 1 & 0 & 1 \\ t_3 & 0 & 1 & 1 & 1 \\ q_1 & 1 & 0 & 1 & 0 \\ q_2 & 0 & 1 & 0 & 0 \end{matrix}$$

Using our proposed similarity matrix

$$m = \begin{matrix} & d_1 & d_2 & d_3 & d_4 \\ t_1 & 1 & 1 & 0 & 0 \\ t_2 & 1 & 1 & 0 & 1 \\ t_3 & 0 & 1 & 1 & 1 \\ c_{1,1} & 1 & 0 & 0 & 0 \\ c_{1,2} & 0 & 0 & 1 & 0 \\ c_{2,1} & 0 & 1 & 0 & 0 \end{matrix}$$

2.2 Establishing a topic-oriented cluster

We employ sIB (sequential Information Bottleneck)[4] to cluster the extracted terms.

2.3 Ranking terms

$$e(y|t) = -\sum_x p(t)p(y|t)\log_2 p(y|t) \quad (3)$$

The variables $X(x \in X)$ and $Y(y \in Y)$ denote the same term collection.

INTRA-CLUSTER AND INTER-CLUSTER SIMILARITY MEASURES:

We expand the candidates together with original query to an n-dimensional vector which are postulated to proximity of a topic.

$$\text{Term vectors: } V_q^t = \{w_{q1}, w_{q2}, \dots, w_{qn}\}^t \quad (t \in C) \quad (4)$$

$$\text{Document vectors: } V_{d_k}^j = \{w_{d_{k1}}, w_{d_{k2}}, \dots, w_{d_{kn}}\}^j \quad (j \in C, k \in N) \quad (5)$$

$$\text{Average intracluster similarity: } \frac{1}{|C|} \sum_{i,j \in C} \sum_{k \in N} \cos(\angle V_q^i, V_{d_k}^j) \quad (6)$$

$$\text{Average intercluster similarity: } \frac{1}{|C| \cdot (|C| - 1)} \sum_{i,j \in C, i \neq j} \sum_{k \in N} \cos(\angle V_q^i, V_{d_k}^j) \quad (7)$$

C is the set of clusters, N denotes the collection of the documents that are retrieved by the new expanded query.

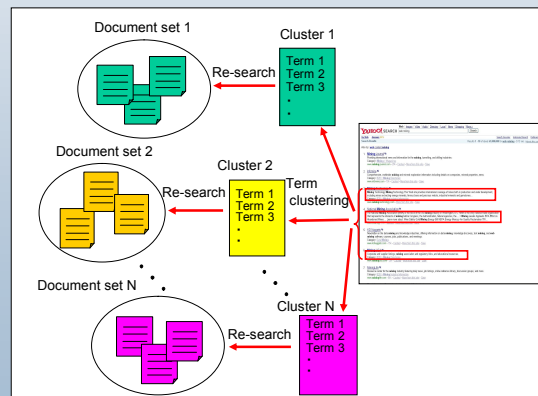


Figure 5: A topic-proximity measure

EXPERIMENT AND EVALUATION:

(1) Experiment in January, 2005

Original query: "Web AND mining"; Search Engine: Google; Term source: Anchor texts that link to the top 10 valid pages.

Terms clustering by using the conventional definition and sIB

Cluster 1	Cluster 2	Cluster 3	Cluster 4	Cluster 5
data	Web	index	art	design
discovery	mining	company	clip	intelligence
knowledge	text	major	intelligence	
resource	reference	major	usage	clip
tech	information	list	usage	clip
warehouse	retrieval	portal	usage	clip
guide	extraction		usage	clip
center	institute		usage	clip
	technology		usage	clip
	workshop		usage	clip
	usage		usage	clip
	analysis		usage	clip
	proceeding		usage	clip

Terms clustering by using our proposed definition and sIB

Cluster 1	Cluster 2	Cluster 3	Cluster 4	Cluster 5
mining(9)	mining(13)	Web(6)	mining(8)	Web(12)
Web(6)	Web(12)	mining(5)	Web(6)	mining(3)
data	text	workshop	institute	design
resource	reference	company	technology	intelligence
tech	information	index	art	designer
map	retrieval	usage	clip	designer
discovery	extraction	analysis	collection	designer
knowledge	list	proceeding	set	designer
warehouse	portal	listing	links	designer
guide	major	major	link	designer
center				designer

Note: There are three topical terms including "Web mining", "mining institute" and "workshop" in the same cluster (cluster 2).

(2) Experiment and Evaluation in August, 2005

- Experiment:** Original query: "Web AND mining"; Search Engine: Google; Term source: 100 extended anchor texts that link to each one of the top 10 valid pages.
- Evaluation:** Collect the top 30 documents by using the extended query term sequences corresponding to each clusters. Compute intracluster and intercluster similarities by using equations (4)-(7).

Experiment:

Term clustering (Old definition)					Term clustering (New definition)				
Cluster 1	Cluster 2	Cluster 3	Cluster 4	Cluster 5	Cluster 1	Cluster 2	Cluster 3	Cluster 4	Cluster 5
business	offer	committee	people	mining	mining(9)	mining(13)	mining(5)	mining(8)	mining(3)
knowledge	clustering	chairman	discovery	data	Web(6)	Web(12)	Web(6)	Web(12)	Web(6)
research	development	thesis	tip	stern	data	offer	business	perspective	ma
team	localization	member	workshop	web	software	clustering	stand	route	track
directory	interface	advisory	project	link	research	people	technology	find	collaboration
speculation	risk	management	procurement	model	yes	acquisition	knowledge	workshop	yes
media	marketing	student	expert	site	intelligence	tip	chairman	de	retrieval
mission	government	conference	library	intelligence	server	database	committee	tree	review
ma	tree	consulting	technique	software	expert	end	advisory	stern	support
record	decision	job	those	traffic	system	localization	member	risk	format

Evaluation:

Intracluster and intercluster similarities (Old definition)					Intracluster and intercluster similarities (New definition)				
Cluster 1	Cluster 2	Cluster 3	Cluster 4	Cluster 5	Cluster 1	Cluster 2	Cluster 3	Cluster 4	Cluster 5
Stability	0.81	0.79	0.83	0.76	Stability	0.81	0.79	0.83	0.76
Cluster 1	0.84	0.80	0.82	0.80	Cluster 1	0.84	0.80	0.82	0.80
Cluster 2	0.80	0.80	0.82	0.80	Cluster 2	0.80	0.80	0.82	0.80
Cluster 3	0.80	0.80	0.82	0.80	Cluster 3	0.80	0.80	0.82	0.80
Cluster 4	0.80	0.80	0.82	0.80	Cluster 4	0.80	0.80	0.82	0.80
Cluster 5	0.80	0.80	0.82	0.80	Cluster 5	0.80	0.80	0.82	0.80

Average intracluster and intercluster similarities		
Old Definition	Intracluster similarity	Intercluster similarity
New definition	1.101	1.038
Improvement	1.572	0.783
	79.1% ↑	36.0% ↓

CONCLUSIONS:

- (1) We proposed a topic-oriented query expansion model and employed IB in this model.
 - (2) In order to treat the multiplicity of the query term meanings at the conceptual level, we have proposed a new definition of the term-term similarity matrix.
 - (3) We have proposed intracluster and intercluster similarity measures to evaluate the relevance between a topic of a cluster and the retrieved documents in the cluster itself as well as the documents in the other clusters.
- Experimental results and evaluations showed that the obtained candidate terms using the new definition of the term-term similarity matrix in each cluster are almost of higher topic relevance than the obtained ones using the old definition.

ACKNOWLEDGMENTS:

The work presented in this paper has been supported by the 21st Century COE (Center of Excellence) Program of Japan Society of the Promotion of Science (JSPS)

BIBLIOGRAPHY:

- [1] B.-Y. Ricardo and R.-N. Berthier. Modern Information Retrieval, 1999.
- [2] E.-N. Efthimiadis. Query expansion, 1996.
- [3] M. Sanderson. Word sense disambiguation and information retrieval, 1994.
- [4] N. Slonim et al. Unsupervised document classification using sequential information maximization, 2002.
- [5] N. Tishby et al. The information bottleneck method, 1999.
- [6] S. Hünrich. Automatic word sense discrimination, 1998.
- [7] T. Hofmann. Probabilistic latent semantic indexing, 1999.
- [8] V. Lavrenko and W. B. Croft. Relevance-based language models, 2001.