

# Learning Query and Document Similarities from Click-through Bipartite Graph with Metadata

Wei Wu<sup>a</sup>, Hang Li<sup>b</sup>, Jun Xu<sup>b</sup>

<sup>a</sup>*Department of Probability and Statistics, Peking University*

<sup>b</sup>*Microsoft Research Asia*

---

## Abstract

We consider learning query and document similarities from a click-through bipartite graph with metadata on the nodes. The metadata consists of multiple types of features of queries and documents. We aim to leverage both the click-through bipartite and the features to learn similarity functions that can effectively measure query-document, document-document, and query-query similarities. The challenges include how to model the similarity functions based on the graph data, and how to accurately and efficiently learn the similarity functions. We propose a method that can solve the problems in a principled way. Specifically, we use two different linear mappings to project queries and documents in two different feature spaces into the same latent space, and take the dot product in the latent space as the similarity function between queries and documents. Query-query and document-document similarities can also be naturally defined as dot products in the latent space. We formalize the learning of similarity functions as learning of the mappings that maximize the similarities of the observed query-document pairs on the enriched click-through bipartite. When queries and documents have multiple types of features, the similarity function is defined as a linear combination of multiple similarity functions, each based on one type of features. We further solve the learning problem by using a new technique called Multi-view Partial Least Squares (M-PLS) developed by us. We prove that the learning problem has a global optimal solution. The global optimum can be efficiently obtained through Singular Value Decomposition (SVD). We also theoretically demonstrate that the proposed method is also capable of finding high quality similar queries. We conducted large scale experiments on enterprise search data and web search data. The experimental results on relevance ranking and similar query finding demonstrate that the proposed method works significantly better than the baseline methods.

*Key words:* Similarity Learning, Multi-view Partial Least Squares, Click-through Bipartite

---

## 1. Introduction

Many tasks in Information Retrieval (IR) rely on similarities between pairs of objects. In relevance ranking, given a query, one retrieves the most relevant documents and ranks them based on their degrees of relevance to the query. The relevance between a query and a document can be viewed as a kind of similarity. In query reformulation or rewriting, queries that convey similar search intents but in different forms are created to refine the original query, so that the documents better meeting the users' information need can be properly retrieved and ranked. The original query and refined queries are in fact similar queries. In query suggestion, queries with

related intents are recommended to the user, to help the user to search other useful information. Those queries are also similar queries. In all these tasks, we need to measure similarities between two objects, either query-document or query-query.

Similarity functions are usually defined based on the features of queries and documents. The relevance models in IR, including Vector Space Model (VSM) [23], BM25 [19], and Language Models for Information Retrieval (LMIR) [18, 36] can all be viewed as similarity functions between query and document feature vectors [34, 32]. Similarity between a query and a document is calculated as similarity between term vectors or n-gram vectors. Similarly, queries are represented as vectors in a term space or n-gram space, and the dot product or cosine is taken as a similarity function between them [37, 35]. All the methods utilize features to calculate similarity functions and we call them feature based methods.

Recently, mining query and document similarities from click-through bipartite graph has also been proposed (cf., [8, 17]). Click-through bipartite, which represents users' implicit judgments on query-query, document-document, and query-document relevance relations, has been proved to be a very valuable source for measuring the similarities. For example, queries which share many co-clicked documents may represent similar search intents, and they can be viewed as similar queries [5, 30]<sup>1</sup>. These methods only rely on the structure of bipartite graph. We call them graph based methods.

In this paper, we consider leveraging information from both click-through bipartite and features to measure query and document similarities. In other words, this is about how to combine the feature based methods and graph based methods. As far as we know, this problem was not well studied previously. We formalize the issue as that of learning query and document similarities from a click-through bipartite graph with metadata on the nodes representing multiple types of features. The features may come from multiple sources. For example, for queries, the content of queries, a taxonomy of semantic classes, and the user information associated with queries can be defined as features; for documents, features can be extracted from the URLs, the titles, the bodies, and the anchor texts of web documents. Formally, we assume that there is a query-document bipartite graph. The bipartite graph consists of triplets  $(q, d, t)$ , where  $q$  denotes a query,  $d$  denotes a document, and  $t$  denotes the number of clicks between  $q$  and  $d$ . Besides, we assume that queries and documents on the bipartite graph have multiple types of features. For each type  $i$  ( $i \geq 1$ ), queries and documents are represented as vectors in query space  $\mathcal{Q}_i$  and document space  $\mathcal{D}_i$ .  $\mathcal{Q}_i$  and  $\mathcal{D}_i$  are subspaces of the Euclidian spaces  $\mathbb{R}^{s_{qi}}$  and  $\mathbb{R}^{s_{di}}$ , where  $s_{qi}$  and  $s_{di}$  represent the dimensionalities of the Euclidean spaces.  $\mathcal{Q}_i$  and  $\mathcal{D}_i$  may or may not be the same space, which means that queries and documents may be either homogeneous or heterogeneous data. Figure 1 illustrates the relationships.

Our goal is to accurately and efficiently learn the similarity functions from the enriched click-through bipartite graph. There are several challenges: 1) how to model the similarity functions based on the complicated graph data, particularly when there are multiple types of features; 2) how to accurately learn the similarity functions; 3) how to efficiently perform the learning task. We propose a method that can solve all the problems in a theoretically sound way.

1) Specifically, for each type of features, we use two linear mappings to project query vectors in the query space and document vectors in the document space into the same latent space. We take the dot product of the images in the latent space as the similarity function between the query and document vectors. Figure 2 illustrates the method. The dot product in the latent

---

<sup>1</sup>There are different ways to define query similarity. In this paper query similarity and document similarity are defined from the viewpoint search intents.

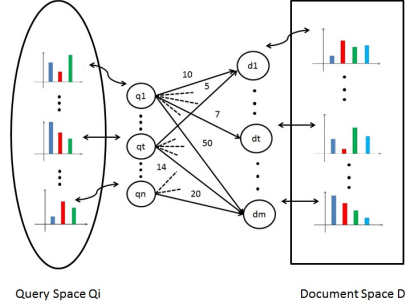


Figure 1: Click-through bipartite with metadata on nodes representing queries and documents in feature spaces.

space is also taken as the similarity function between query and query and the similarity function between document and document. The two mappings are supposed to be different, because of the difference between the query space and document space. The final similarity functions are defined as linear combinations of similarity functions from different feature types.

2) We learn the mappings and combination weights using the enriched click-through bipartite data. The numbers of clicks between queries and documents represent similarities between queries and documents. We formalize the learning method as an optimization problem, in which the objective is to maximize the similarities of the observed query-document pairs on the click-through bipartite. We regularize the combination weights with  $L^2$  norm and we make orthogonal assumptions on the mappings. We propose a new machine learning technique called Multi-view Partial Least Squares (M-PLS) to learn the linear mappings. M-PLS is an extension of the Partial Least Squares technique into multi-view (multiple feature types). When there is only one type of features (one view), then M-PLS degenerates to the conventional PLS. Moreover, if there is no metadata used, then M-PLS becomes equivalent to Latent Semantic Indexing (LSI) and thus our method can also be regarded as an extension of LSI. We prove that although the optimization problem is not convex, the solution of M-PLS is global optimal. We also conduct theoretical analysis on the method, and point out that it has two properties that enable the method to capture query-query similarity as well (this is also true for document-document similarity), although the learning of it is not explicitly incorporated into the formulation.

3) The learning task can be efficiently performed through Singular Value Decomposition (SVD). First, we employ the power method to build an SVD solver. After that, we solve a simple quadratic program to learn the optimal combination weights. We show that the quadratic program has a close form solution.

We conducted experiments on large scale enterprise search data and web search data. The results on relevance ranking and similar query finding show that our method significantly outperforms the baseline methods. Specifically, in relevance ranking, we compare our method with the state of the art feature based methods such as BM25, graph based methods such as random walk and their linear combinations. Our method not only significantly outperforms the feature based methods and graph based methods, but also significantly performs better than their linear combinations. In similar query finding, we use examples to demonstrate the capability of our method on finding high quality similar queries. When compared with the baseline methods such as random walk and cosine similarity, our method significantly outperforms the feature based methods and graph based methods, and is comparable with the best linear combination. Particularly, we find our method performs better on tail queries, and the results also strongly support

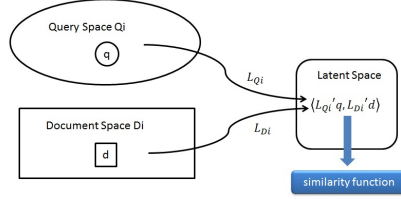


Figure 2: Projecting queries and documents from the query space and document space into a latent space using two linear mappings  $L_{Q_i}$  and  $L_{D_i}$ . The dot product in the latent space is taken as the similarity function.

the conclusions of our theoretical analysis.

Our contributions in this paper include: 1) proposal of a new problem on learning query and document similarities from an enriched click-through bipartite graph; 2) proposal of a new learning method called M-PLS to perform the learning task; 3) theoretical proof of the properties of the proposed method and empirical demonstration of the effectiveness of the method.

The rest of the paper is organized as follows: in Section 2, we conduct a survey on related work. Then, we define the similarity learning problem in Section 3. After that, we explain our method based on M-PLS in Section 4. Experimental results are reported in Section 5, and Section 6 concludes the paper and provides some future research directions. The proof of theorem is given in Appendix.

## 2. Related Work

Partial Least Squares (PLS) [20] refers to a class of algorithms in statistics that aims to model the relations between two or more sets of data by projecting them into a latent space. The underlying idea of PLS algorithms is to model collinearity between different data sets. Since the first work by Wold [31], many variants of PLS have been developed and used in many tasks such as regression [26], classification [4] and dimension reduction [25]. In practice, PLS has been successfully applied in chemometrics [20], biometrics [7], computer vision [25] and graph mining [22]. In this paper, we extend PLS to M-PLS and employ M-PLS in web search. If there is only one type of feature in the query as well as in the document space, our similarity learning problem degenerates to a problem which can be directly solved by PLS.

Canonical Correlation Analysis (CCA) [14] or its kernelized version KCCA [12, 13] is an alternative method to PLS. Both attempt to learn linear mapping functions to project objects in two spaces into the same latent space. The difference between CCA and PLS is that CCA learns cosine as the similarity function and PLS learns dot product as the similarity function. In our work, we choose PLS instead of CCA, because it is easier to enhance the efficiency for the former. In CCA it is necessary to compute the inverse of large matrices<sup>2</sup>, which is computationally expensive.

Measuring query and document similarities is always an important research topic in IR. Existing work on query and document similarities can be divided into two groups: feature based methods and graph based methods. In the former group, Vector Space Model (VSM) [23], BM25 [19], and Language Models for Information Retrieval (LMIR) [18, 36] make use of features,

<sup>2</sup>[http://en.wikipedia.org/wiki/Canonical\\_correlation](http://en.wikipedia.org/wiki/Canonical_correlation)

particularly, n-gram features to measure query-document similarity. As pointed out by Xu et al. [34] and others that these models can be viewed as models using the dot product between a query vector and a document vector as the query-document similarity function. Similarly, queries can also be represented as n-grams, and the cosine or dot product can be utilized as the similarity function between them [37, 35]. In [35], queries are represented as n-gram vectors, and a cosine similarity function is learned by using distance metric learning. In [6], the authors propose calculating query similarity with term and n-gram features enriched with a taxonomy of semantic classes. In [3], queries are represented as vectors in a high dimensional space with each dimension corresponding to a URL. The click frequency on a URL is used as the value of the corresponding dimension. In the latter group, graph based methods exploit the structure of click-through bipartite to learn the query-query, document-document, and query-document similarities. For example, Latent Semantic Indexing (LSI) [10] can be employed, which uses SVD to project queries and documents in a click-through bipartite into a latent space, and calculates query-document, query-query, and document-document similarities through the dot product of their images in the latent space. In [16, 1], the authors propose determining the similarity of a pair of objects based on the similarity of other pairs of objects, and the final similarity measure is iteratively calculated on a bipartite graph. A click-through bipartite can also be used in clustering of similar queries [5, 30]. Craswell and Szummer [8] extend the idea and propose adopting a backward random walk process on a click-through bipartite to propagate similarity through probabilistic transitions. Our method aims to leverage both features and click-through bipartite to more accurately learn the similarities.

Methods for learning similarities between objects by utilizing bipartite graphs built from multiple sources have also been studied. In [9], term relations and document relations are integrated into a document-term bipartite. In [28], the authors extend LSI when information from different types of bipartite graphs is available. In [33], the authors use a unified relationship matrix to represent different types of objects and their relations. In [17], matrix factorization is simultaneously conducted on two bipartites, a user-query bipartite and a query-document bipartite. In [11], similarity scores are first calculated by the LMIR model and then the scores are propagated on a click-through bipartite to find similar queries. The method proposed in this paper is different from these existing methods. In our method, we make use of both click-through bipartite and multiple types of features on the nodes. In that sense, we combine a feature based method and graph based method. In contrast, the existing methods are all graph-based methods and they use either more than one bipartite or a graph with more than one type of relation.

The method proposed in this paper is also related to multi-view learning [21] in machine learning. In multi-view learning, instances are assumed to have multiple types of features and the goal is to exploit the different types of features to learn a model. In our work, we assume that queries and documents on a click-through bipartite graph have multiple types of features, which is similar to the assumption in multi-view learning. Our work is unique in that we perform multi-view similarity learning on an enriched click-through bipartite graph. Our method can be regarded as Multi-view PLS, an extension of PLS to multi-view learning.

The similarity learning problem studied in this paper is not limited to search, and it can be potentially extended to other applications such as recommendation system and online advertisement. Recently, Wang et al. [27] propose a method for advertisement in a setting similar to ours. We learn query and document similarities from a bipartite graph with metadata, while they predict possible clicks between users and advertisements on a network with metadata. Note that there is significant difference between our method and their method. First, they only measure similarity between users and advertisements (corresponding to queries and documents), while we

also measure query-query similarity and document-document similarity. Second, their method assumes that the latent space is predefined and is independent from the similarity function, while in our method the latent space as well as the similarity function are learned from data.

### 3. Problem Formulation

Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  denote a click-through bipartite.  $\mathcal{V} = \mathcal{Q} \cup \mathcal{D}$  is the set of vertices, which consists of a query set  $\mathcal{Q} = \{q_i\}_{i=1}^m$  and a document set  $\mathcal{D} = \{d_i\}_{i=1}^n$ .  $\mathcal{E}$  is the set of edges between query vertices and document vertices. The edge  $e_{ij} \in \mathcal{E}$  between query vertex  $q_i$  and document vertex  $d_j$  is weighted by the click number  $t_{ij}$ . We assume that  $\mathcal{G}$  is an undirected graph. Besides, we assume that there exists rich metadata on the vertices of the bipartite. The metadata consists of  $l$  types of features ( $l \geq 1$ ). The features may stand for the content of queries and documents and the clicks of queries and documents on the bipartite [3], as will be seen in Section 5. For each type  $i$  ( $1 \leq i \leq l$ ), query  $q \in \mathcal{Q}$  and document  $d \in \mathcal{D}$  are represented as vectors in space  $\mathcal{Q}_i$  and space  $\mathcal{D}_i$ , respectively, where  $\mathcal{Q}_i$  and  $\mathcal{D}_i$  are subspaces of the Euclidean spaces  $\mathbb{R}^{s_{qi}}$  and  $\mathbb{R}^{s_{di}}$ . Figure 1 illustrates the relationships.

Our goal is to leverage both the click-through bipartite and the features to learn query and document similarities. The similarities may include query-document similarity, query-query similarity, and document-document similarity. In this paper we only study query-document similarity and query-query similarity, and the same method can be applied to learning of document-document similarity. Formally, we learn two similarity functions  $f(q, d)$  and  $g(q, q')$  given  $\mathcal{G}$  and  $\{(\mathcal{Q}_i, \mathcal{D}_i)\}_{i=1}^l$ . Similarity function  $f(q, d)$  measures the similarity between query  $q \in \mathcal{Q}$  and document  $d \in \mathcal{D}$ , and similarity function  $g(q, q')$  measures the similarity between queries  $q, q' \in \mathcal{Q}$ .

The question is then how to model the similarity functions  $f(q, d)$  and  $g(q, q')$ , and how to accurately and efficiently learn the similarity functions. We propose learning the similarity functions by linearly projecting the queries and documents in the query spaces and document spaces into latent spaces. For each type of feature spaces  $(\mathcal{Q}_i, \mathcal{D}_i)$ , we learn two linear mappings  $L_{\mathcal{Q}_i}$  and  $L_{\mathcal{D}_i}$ .  $L_{\mathcal{Q}_i}$  is an  $s_{qi} \times k_i$  dimensional matrix which can map a query  $q$  from  $\mathcal{Q}_i$  to the  $k_i$  dimensional latent space  $\mathcal{K}_i$ , and  $L_{\mathcal{D}_i}$  is an  $s_{di} \times k_i$  dimensional matrix which can map a document  $d$  from  $\mathcal{D}_i$  to the  $k_i$  dimensional latent space  $\mathcal{K}_i$ . We assume that  $\forall i, k_i \leq \min(s_{qi}, s_{di})$ .  $L_{\mathcal{Q}_i}$  and  $L_{\mathcal{D}_i}$  are different mapping functions due to the difference between  $\mathcal{Q}_i$  and  $\mathcal{D}_i$ . Given  $q \in \mathcal{Q}$  and  $d \in \mathcal{D}$ , the images in the latent space  $\mathcal{K}_i$  are  $L_{\mathcal{Q}_i}^\top q$  and  $L_{\mathcal{D}_i}^\top d$ , respectively. We define similarity functions  $f_i(q, d)$  and  $g_i(q, q')$  as  $f_i(q, d) = q^\top L_{\mathcal{Q}_i} L_{\mathcal{D}_i}^\top d$ ,  $g_i(q, q') = q^\top L_{\mathcal{Q}_i} L_{\mathcal{Q}_i}^\top q'$ . In other words, we take the dot products of images of queries and documents in the latent spaces as similarity functions. The final similarity functions  $f(q, d)$  and  $g(q, q')$  are defined as linear combinations of  $\{f_i(q, d)\}_{i=1}^l$  and  $\{g_i(q, q')\}_{i=1}^l$ :  $f(q, d) = \sum_{i=1}^l \alpha_i f_i(q, d)$ ,  $g(q, q') = \sum_{i=1}^l \alpha_i g_i(q, q')$ , where  $\alpha_i \geq 0$  is a combination weight.

We learn the mappings  $\{(L_{\mathcal{Q}_i}, L_{\mathcal{D}_i})\}_{i=1}^l$  and combination weights  $\{\alpha_i\}_{i=1}^l$  from the click-through bipartite  $\mathcal{G}$  and the features  $\{(\mathcal{Q}_i, \mathcal{D}_i)\}_{i=1}^l$ . Specifically, we view the click number of a query-document pair as an indicator of their similarity. We learn  $\{(L_{\mathcal{Q}_i}, L_{\mathcal{D}_i})\}_{i=1}^l$  and  $\{\alpha_i\}_{i=1}^l$  by maximizing the similarities of the observed query-document pairs on the click-through bipartite. The underlying assumption is that the higher the click number is, the more similar the query and the document are in the latent spaces. In graph based methods (cf., [8]), click numbers are usually directly used for query and document similarity calculation. In this paper, we take a transformation of click numbers and use the logarithm of them, which is verified to be effective.

Finally, we consider the following learning problem:

$$\arg \max_{\{(L_{Q_i}, L_{\mathcal{D}_i})\}_{i=1}^l, \{\alpha_i\}_{i=1}^l} \sum_{e_{uv} \in \mathcal{E}} \sum_{i=1}^l \alpha_i \cdot (\log(t_{uv}) \cdot q_u^{i\top} L_{Q_i} L_{\mathcal{D}_i}^\top d_v^i), \quad (1)$$

where  $q_u^i$  and  $d_v^i$  represent the feature vectors of query  $q_u$  and document  $d_v$  in  $Q_i$  and  $\mathcal{D}_i$ , respectively.

Note that an alternative method for learning query and document similarities from  $\mathcal{G}$  and  $\{(Q_i, \mathcal{D}_i)\}_{i=1}^l$  is to concatenate different types of feature vectors of queries as well as different types of feature vectors of documents and learning two mappings for queries and documents with the two concatenated vectors. The method is a special case of the learning method (1). We compare the performances of our proposed method (1) and this alternative method in Section 5.

Objective function (1) may go to infinity, since there are no constraints on the mappings  $\{(L_{Q_i}, L_{\mathcal{D}_i})\}_{i=1}^l$  and the weights  $\{\alpha_i\}_{i=1}^l$ . The features  $\{(q_u^i, d_v^i) \mid e_{uv} \in \mathcal{E}, 1 \leq i \leq l\}$  are also not bounded. We consider adding proper constraints to the mapping functions  $\{(L_{Q_i}, L_{\mathcal{D}_i})\}_{i=1}^l$  and the weights  $\{\alpha_i\}_{i=1}^l$ , as shown in the next section.

#### 4. Our Method

We further formalize the learning problem in (1) as a constrained optimization problem. We propose a new learning technique called Multi-view Partial Least Squares (M-PLS) to solve the problem, which is a multi-view learning version of PLS. We prove that the problem has a global optimal solution. The optimal mappings can be obtained by Singular Value Decomposition (SVD) and the optimal weights can be obtained by quadratic programming. We present the algorithm and its complexity. We also conduct theoretical analysis to demonstrate that our method has the capability to find high quality similar queries, although query-query similarity is not directly represented in the formulation.

##### 4.1. Constrained Optimization Problem

First, we normalize the feature vectors such that  $\forall u, v, i, \|q_u^i\| = 1$  and  $\|d_v^i\| = 1$ . Second, we add orthogonal constraints on the mapping matrices  $\{(L_{Q_i}, L_{\mathcal{D}_i})\}_{i=1}^l$ . Finally, we introduce  $L^2$  regularization on the weights  $\{\alpha_i\}_{i=1}^l$ .

The similarity learning method is re-formalized as

$$\begin{aligned} & \arg \max_{\{(L_{Q_i}, L_{\mathcal{D}_i})\}_{i=1}^l, \{\alpha_i\}_{i=1}^l} \sum_{e_{uv} \in \mathcal{E}} \sum_{i=1}^l \alpha_i \cdot (\log(t_{uv}) \cdot q_u^{i\top} L_{Q_i} L_{\mathcal{D}_i}^\top d_v^i) \\ & \text{subject to} \quad L_{Q_i}^\top L_{Q_i} = I_{k_i \times k_i}, L_{\mathcal{D}_i}^\top L_{\mathcal{D}_i} = I_{k_i \times k_i}, \alpha_i \geq 0, \sum_{i=1}^l \alpha_i^2 \leq 1, \end{aligned} \quad (2)$$

where  $k_i \leq \min(s_{qi}, s_{di})$  is a parameter and  $I_{k_i \times k_i}$  is an identity matrix. A larger  $k_i$  means preserving more information in the projection for the type of features. Although problem (2) is not convex, we can prove that the global optimum exists. Moreover, the global optimal mappings  $\{(L_{Q_i}, L_{\mathcal{D}_i})\}_{i=1}^l$  can be obtained by Singular Value Decomposition (SVD), and the global optimal weights  $\{\alpha_i\}_{i=1}^l$  can be obtained by quadratic programming having a close form solution.

#### 4.2. Global Optimal Solution

Basically, there are two steps in finding the global optimum. First, for each type of features, we find the optimal mappings through solving SVD. Second, we determine the optimal combination weights.

Specifically, the objective function (2) can be re-written as

$$\begin{aligned}
& \sum_{e_{uv} \in \mathcal{E}} \sum_{i=1}^l \alpha_i \cdot (\log(t_{uv}) \cdot q_u^{i\top} L_{Q_i} L_{\mathcal{D}_i}^\top d_v^i) \\
&= \sum_{i=1}^l \alpha_i \cdot \text{Trace} \left( \sum_{e_{uv} \in \mathcal{E}} \log(t_{uv}) \cdot q_u^{i\top} L_{Q_i} L_{\mathcal{D}_i}^\top d_v^i \right) \\
&= \sum_{i=1}^l \alpha_i \cdot \text{Trace} \left( L_{\mathcal{D}_i}^\top \left( \sum_{e_{uv} \in \mathcal{E}} \log(t_{uv}) d_v^i q_u^{i\top} \right) L_{Q_i} \right) \\
&= \sum_{i=1}^l \alpha_i \cdot \text{Trace} (L_{\mathcal{D}_i}^\top M_i L_{Q_i}),
\end{aligned}$$

where  $M_i$  is defined as  $\sum_{e_{uv} \in \mathcal{E}} \log(t_{uv}) d_v^i q_u^{i\top}$ .

Suppose that  $\forall i, k_i \leq \min(s_{q_i}, s_{d_i})$ , the following theorem indicates that the optimization problem

$$\begin{aligned}
& \arg \max_{L_{Q_i}, L_{\mathcal{D}_i}} \text{Trace} (L_{\mathcal{D}_i}^\top M_i L_{Q_i}) \\
& \text{subject to } L_{Q_i}^\top L_{Q_i} = I_{k_i \times k_i}, L_{\mathcal{D}_i}^\top L_{\mathcal{D}_i} = I_{k_i \times k_i}
\end{aligned} \tag{3}$$

has a global optimum and the global optimal solution can be obtained through SVD of  $M_i$ :

**Theorem 4.1.**  $\forall k_i \leq \min(s_{q_i}, s_{d_i})$ , the global optimal solution of problem (3) exists. Furthermore, suppose that  $M_i = U_i \Sigma_i V_i^\top$ , where  $\Sigma_i$  is an  $s_{d_i} \times s_{q_i}$  diagonal matrix with singular values  $\lambda_1^i \geq \lambda_2^i \geq \dots \lambda_p^i \geq 0$ ,  $p = \min(s_{q_i}, s_{d_i})$ ,  $U_i = (u_1^i, u_2^i, \dots, u_{s_{d_i}}^i)$  where  $\{u_j^i\}$  are left singular vectors, and  $V_i = (v_1^i, v_2^i, \dots, v_{s_{q_i}}^i)$  where  $\{v_j^i\}$  are right singular vectors. The global maximum is given by  $\sum_{j=1}^{k_i} \lambda_j^i$  and the global optimal  $\hat{L}_{Q_i}$  and  $\hat{L}_{\mathcal{D}_i}$  are given by  $\hat{L}_{Q_i} = (v_1^i, v_2^i, \dots, v_{k_i}^i)$  and  $\hat{L}_{\mathcal{D}_i} = (u_1^i, u_2^i, \dots, u_{k_i}^i)$ , respectively.

The proof is given in Appendix. With Theorem 4.1, if we define  $\Lambda_i = \sum_{j=1}^{k_i} \lambda_j^i$ , then we can re-write problem (2) as

$$\begin{aligned}
& \max_{\substack{\{(L_{Q_i}, L_{\mathcal{D}_i})\}_{i=1}^l, \{\alpha_i\}_{i=1}^l \\ L_{Q_i}^\top L_{Q_i} = L_{\mathcal{D}_i}^\top L_{\mathcal{D}_i} = I_{k_i \times k_i}, \alpha_i \geq 0, \sum_{i=1}^l \alpha_i^2 \leq 1}} \sum_{e_{uv} \in \mathcal{E}} \sum_{i=1}^l \alpha_i \cdot (\log(t_{uv}) \cdot q_u^{i\top} L_{Q_i} L_{\mathcal{D}_i}^\top d_v^i) \\
&= \max_{\substack{\{\alpha_i\}_{i=1}^l \\ \alpha_i \geq 0, \sum_{i=1}^l \alpha_i^2 \leq 1}} \sum_{i=1}^l \alpha_i \max_{\substack{\{(L_{Q_i}, L_{\mathcal{D}_i})\}_{i=1}^l \\ L_{Q_i}^\top L_{Q_i} = L_{\mathcal{D}_i}^\top L_{\mathcal{D}_i} = I_{k_i \times k_i}}} \sum_{e_{uv} \in \mathcal{E}} (\log(t_{uv}) \cdot q_u^{i\top} L_{Q_i} L_{\mathcal{D}_i}^\top d_v^i) \\
&= \max_{\substack{\{\alpha_i\}_{i=1}^l \\ \alpha_i \geq 0, \sum_{i=1}^l \alpha_i^2 \leq 1}} \sum_{i=1}^l \alpha_i \Lambda_i.
\end{aligned} \tag{4}$$



It is easy to prove that problem (4) has a close form global optimum. The optimal weights are given by

$$\hat{\alpha}_i = \frac{\Lambda_i}{\sqrt{\sum_{i=1}^l \Lambda_i^2}}, \quad 1 \leq i \leq l. \quad (5)$$

It is easy to verify that when  $l = 1$ , problem (2) becomes a problem solvable by Partial Least Squares (PLS) [20, 24]. The orthogonal assumptions on the mapping matrices make it feasible to employ PLS to solve the problem. Therefore, our method of Multi-view PLS is an extension of PLS.

If we assume that only the click-through bipartite is used, i.e., no feature is used, then problem (2) becomes equivalent to Latent Semantic Indexing (LSI) [10]<sup>3</sup>. In other words, LSI is a special case of our method. Specifically, in such case  $l = 1$ , both query vectors and document vectors are orthogonal, and in each query vector and each document vector, only one element equals to one that indicates the index of a query or a document and the other elements are zero. As a result, matrix  $M_i$  for SVD is exactly the one in LSI.

In problem (2), we consider using  $L^2$  norm to regularize the weights  $\{\alpha_i\}_{i=1}^l$ . An alternative method would be to use  $\sum_{i=1}^l \alpha_i \leq 1$  (i.e.,  $L^1$  norm) to regularize them. However, such a regularization will make the final solution become  $\alpha_i = 1$ , if  $\Lambda_i = \max_{1 \leq j \leq l} \Lambda_j$ , and  $\alpha_i = 0$ , otherwise. This is a degenerative solution and is not desirable. The  $L^2$  regularization in our method does not suffer from the problem.

#### 4.3. Algorithm

The learning algorithm is described in Algorithm 1. The algorithm contains two steps. First, for each type of features, it calculates  $M_i$ , and solves SVD of  $M_i$  to learn the linear mappings. Then, it calculates the combination weights using (5).

At Step 2.a, it is necessary to calculate  $M_i$ . Suppose that there are  $n_q$  queries on the click-through bipartite  $\mathcal{G}$ . Each query has on average  $\kappa_q$  clicked documents. Then for each type of features, the worst case time complexity of calculating  $M_i$  is of order  $O(n_q \cdot \kappa_q \cdot s_{qi} \cdot s_{di})$ , where  $s_{qi}$  and  $s_{di}$  denote the dimensionalities of query space  $\mathcal{Q}_i$  and document space  $\mathcal{D}_i$ , respectively. Although  $s_{qi}$  and  $s_{di}$  can be very large, the query vectors and document vectors are usually very sparse. This can make the average time complexity much smaller than the worst case time complexity. Suppose that each dimension of query vectors has on average  $c_{qi}$  non-zero values, and each document vector has on average  $c_{di}$  non-zero values. The average time complexity of calculating  $M_i$  becomes of order  $O(s_{qi} \cdot c_{qi} \cdot \kappa_q \cdot c_{di})$ . Since  $c_{qi}$  and  $c_{di}$  are much smaller than  $n_q$  and  $s_{di}$ ,  $M_i$  can be calculated very efficiently.

We empirically find in our experiments that sparse query and document feature vectors also make the matrix  $M_i$  sparse. In our experiments, the ratio of non-zero elements in the matrix  $M_i$  is not larger than 0.5%. Therefore, we can employ the power method (cf., [29]) to build an SVD solver. It is an iterative algorithm, and after  $i$  iterations, the  $i^{th}$  largest singular values can be obtained. Suppose that there are  $C$  non-zero elements in matrix  $M_i$ . The time complexity for calculating each singular value is  $O(C + k_i \cdot \max(s_{qi}, s_{di}))$ , and the total time complexity is  $O(k_i \cdot C + k_i^2 \cdot \max(s_{qi}, s_{di}))$ . Since  $M_i$  is sparse and  $C$  is much smaller than  $s_{qi} \cdot s_{di}$ , when  $k_i$  is small, SVD of  $M_i$  can be solved efficiently.

---

<sup>3</sup>LSI is usually used for learning the similarity between term and document from a term-document bipartite graph.

---

**Algorithm 1**

---

- 1: Input: click-through bipartite  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , feature spaces  $\{(Q_i, \mathcal{D}_i)\}_{i=1}^l$ , parameters  $\{k_i\}_{i=1}^l$ .
  - 2: For each type of feature spaces  $(Q_i, \mathcal{D}_i)$ 
    - a. Calculate  $M_i$  using  $\sum_{e_{uv} \in \mathcal{E}} \log(t_{uv}) d_v^i q_u^{i\top}$ .
    - b. Calculate SVD of  $M_i$ .
    - c. Take the first  $k_i$  left singular vectors  $(u_1^i, u_2^i, \dots, u_{k_i}^i)$  as  $\hat{L}_{\mathcal{D}_i}$ , and first  $k_i$  right singular vectors  $(v_1^i, v_2^i, \dots, v_{k_i}^i)$  as  $\hat{L}_{Q_i}$ .
    - d. Calculate  $\Lambda_i$  using  $\sum_{j=1}^{k_i} \lambda_j^i$ .
  - 3: Calculate combination weight  $\hat{\alpha}_i$  using equation (5).
  - 4: Output: similarity functions:
    - a.  $\forall q, d, f(q, d) = \sum_{i=1}^l \hat{\alpha}_i \cdot q^{\top} \hat{L}_{Q_i} \hat{L}_{\mathcal{D}_i}^{\top} d^i$ .
    - b.  $\forall q, q', g(q, q') = \sum_{i=1}^l \hat{\alpha}_i \cdot q^{\top} \hat{L}_{Q_i} \hat{L}_{Q_i}^{\top} q'^i$ .
- 

#### 4.4. Query Similarity

Problem (2) is formalized as learning of query-document similarity. Actually the optimization can also help learning of query-query similarity (equivalently document-document similarity). Here, we show that our method is able to find high quality similar queries as well. Under the orthogonal assumptions, the results of problem (2) have two properties which can guarantee that the performance of the method on similar query finding is high as well.

For a specific feature type  $i$ , if the similarity between  $q, q' \in Q_i$  is high, then the similarity between images of  $q$  and  $q'$  in the latent space is also high, provided that  $k_i$  is sufficiently large. Formally, we assume that there is a small number  $\varepsilon/2 > 0$  such that  $q^{\top} q' \geq 1 - \varepsilon/2$  (i.e., they have high similarity because  $q^{\top} q' \leq 1$ ). Since  $\|q\| = \|q'\| = 1$ , we have  $\|q - q'\|^2 = \|q\|^2 + \|q'\|^2 - 2q^{\top} q' \leq \varepsilon$ . Suppose that  $L_{Q_i} = (l_1^{Q_i}, \dots, l_{k_i}^{Q_i})$ . We have

$$\|L_{Q_i}^{\top} q - L_{Q_i}^{\top} q'\|^2 = (q - q')^{\top} \sum_{j=1}^{k_i} l_j^{Q_i} l_j^{Q_i\top} (q - q') = \sum_{j=1}^{k_i} |\langle l_j^{Q_i}, q - q' \rangle|^2.$$

Since  $L_{Q_i}^{\top} L_{Q_i} = I_{k_i \times k_i}$ , we have

$$\|L_{Q_i}^{\top} q - L_{Q_i}^{\top} q'\|^2 = \sum_{j=1}^{k_i} |\langle l_j^{Q_i}, q - q' \rangle|^2 \leq \|q - q'\|^2 \leq \varepsilon.$$

Thus, we have  $q^{\top} L_{Q_i} L_{Q_i}^{\top} q' \geq (\|L_{Q_i}^{\top} q\|^2 + \|L_{Q_i}^{\top} q'\|^2 - \varepsilon)/2$ . This inequality indicates that for a specific feature type  $i$ , if the similarity between two queries  $q$  and  $q'$  is high in the feature space  $Q_i$ , then the similarity between their images in the latent space is determined by their norms in the space. The square of the norm of query  $q$  in the latent space can be calculated as

$$\|L_{Q_i}^{\top} q\|^2 = q^{\top} L_{Q_i} L_{Q_i}^{\top} q = q^{\top} \sum_{j=1}^{k_i} l_j^{Q_i} l_j^{Q_i\top} q = \sum_{j=1}^{k_i} |\langle l_j^{Q_i}, q \rangle|^2. \quad (6)$$

Since usually  $k_i \leq s_{qi}$ , we have  $\|L_{Q_i}^{\top} q\|^2 \leq \|q\|^2 = 1$ . With a sufficiently large  $k_i$ , we can assume that  $\|L_{Q_i}^{\top} q\|^2 \geq 1 - \delta$ , and  $\|L_{Q_i}^{\top} q'\|^2 \geq 1 - \delta$ , where  $\delta \geq 0$  is a small number. Thus, we obtain

$q^\top L_{Q_i} L_{Q_i}^\top q' \geq 1 - \delta - \frac{\varepsilon}{2}$ . Note that  $\delta$  monotonically decreases when  $k_i$  increases. In an extreme case, when  $k_i = s_{qi}$ ,  $\delta = 0$  holds. We call this property *inheritance*. This property ensures that when sufficient information is preserved in the latent space (i.e.,  $k_i$  is sufficiently large), similar queries in the original feature space will also be similar in the latent space.

Second, suppose that  $q_1, q_2 \in Q_i$ ,  $d_1, d_2 \in \mathcal{D}_i$ , and  $q'_1 = L_{Q_i}^\top q_1$ ,  $q'_2 = L_{Q_i}^\top q_2$ ,  $d'_1 = L_{\mathcal{D}_i}^\top d_1$ , and  $d'_2 = L_{\mathcal{D}_i}^\top d_2$ . If we assume that  $q'_1$  and  $d'_1$ ,  $q'_2$  and  $d'_2$ , and  $d'_1$  and  $d'_2$  are similar pairs with high similarity, then when  $k_i$  is sufficiently large, we can obtain high similarity between  $q'_1$  and  $q'_2$ . Formally, we assume that there is a small number  $\varepsilon > 0$  such that  $q'_1{}^\top d'_1 \geq 1 - \varepsilon/18$ ,  $q'_2{}^\top d'_2 \geq 1 - \varepsilon/18$ , and  $d'_1{}^\top d'_2 \geq 1 - \varepsilon/18$ . We have  $\|q'_1 - q'_2\| \leq \|q'_1 - d'_1\| + \|d'_1 - d'_2\| + \|d'_2 - q'_2\|$ . Since  $q'_1 = L_{Q_i}^\top q_1$  and  $d'_1 = L_{\mathcal{D}_i}^\top d_1$ , similar to the analysis of equation (6), we have  $\|q'_1\| \leq \|q_1\| = 1$  and  $\|d'_1\| \leq \|d_1\| = 1$ . We know that  $\|q'_1 - d'_1\|^2 = \|q'_1\|^2 + \|d'_1\|^2 - 2q'_1{}^\top d'_1 \leq \varepsilon/9$ . Similarly, we can get  $\|q'_2 - d'_2\|^2 \leq \varepsilon/9$  and  $\|d'_1 - d'_2\|^2 \leq \varepsilon/9$ . Thus, we obtain  $\|q'_1 - q'_2\|^2 \leq \varepsilon$ . Similarly to the analysis above, we can say that with a sufficiently large  $k_i$ , we can have the similarity between  $q'_1$  and  $q'_2$  large enough. We call this property *transitivity*. The property ensures that when sufficient information is preserved in the latent space, we can derive similar query pairs from similar query-document pairs. The property needs high similarity between documents in the latent space. The condition can be easily met, because document similarity can be preserved with the “inheritance” property.

In our experiments, we find that the inheritance property and the transitivity property can really help us on finding high quality similar queries. We also give an example to support our theoretical analysis.

## 5. Experiments

We conducted experiments to test the performance of the proposed method on relevance ranking and similar query finding. We used two data sets: enterprise search data and web search data.

### 5.1. Data Sets

We collected one year click-through data from an enterprise search engine of an IT company and one week click-through data from a commercial web search engine. We built two click-through bipartite graphs with the data. If there are more than 3 clicks between a query and a URL, we added a link between them. There are 51,699 queries and 81,186 URLs on the bipartite of enterprise search data. There are 94,022 queries and 111,631 URLs on the bipartite of web search data. That means that we discarded the links between queries and URLs whose click frequency is lower than or equal to 3. Each query has on average 2.44 clicked URLs in the enterprise bipartite and each query has on average 1.74 clicked URLs in the web bipartite.

We extracted features for both data sets. Two types of the features were extracted and attached to the nodes of the bipartite graphs. First, we took the words in queries and the words in URLs as features, referred to as “word features”. Words were stemmed and stop words were removed. With word features, each query is represented by a tf-idf vector in the query space, and each URL is also represented by a tf-idf vector in the document space. Each dimension of the query space corresponds to a unique term and so does each dimension of the document space. Note that the two spaces are of high dimension and very sparse. For the enterprise search data, there are 9,958 unique terms. For the web search data, the dimensionality of term space is 10,791. Next, we took the numbers of clicks of URLs as the features of queries and the numbers of clicks of

queries as features of URLs. We call the features “graph features”, because they are derived from the bipartite graph. Each query is represented by a vector of log click-numbers in the query space and each document is represented by a vector of log click-numbers in the document space. Each dimension of the query space corresponds to a URL on the click-through bipartite, and similarly each dimension of the document space corresponds to a query on the click-through bipartite. The dimensionalities of the two spaces are also very high, while the densities of the two spaces are very low.

We randomly sampled queries and their associated URLs from the click-through bipartites of the two data sets. The relevance between the queries and associated URLs were made by human experts. There are five level judgments, including “Perfect”, “Excellent”, “Good”, “Fair”, and “Bad”. For the enterprise search data, 1,701 queries and their associated URLs have judgments. Each query has on average 16 judged URLs. For the web search data, the number of judged queries is 4,445, and each query has on average 11.34 judged URLs.

## 5.2. Experiment Setup

We consider four alternatives to learn a similarity function using our method: 1) Only word features are used. We denote the model as  $M\text{-PLS}_{\text{Word}}$ ; 2) Only graph features are used. We refer to the model as  $M\text{-PLS}_{\text{Bipar}}$ ; 3) The vectors from the word feature space and the vectors from the graph feature space are concatenated to create long vectors. The corresponding space is the Cartesian product of the word feature space and the graph feature space. We call the model  $M\text{-PLS}_{\text{Conca}}$ ; 4) We apply our method to learn the similarity function assuming that queries and documents have two types of features and we learn a linearly combined similarity function. We call the model  $M\text{-PLS}_{\text{Com}}$ .

As baselines, we consider feature based methods, graph based methods, and their linear combinations. In relevance ranking, we take BM25 as an example of feature based methods. We choose LSI on click-through bipartite graph and random walk (RW for short) [8] as examples of graph based methods. We also use a linear combination of BM25 and LSI as well as a linear combination of BM25 and RW. In similar query finding, besides LSI and RW, we adopt as baseline methods cosine similarity of two query vectors represented with graph features and cosine similarity of two query vectors represented with word features. We denote them as  $\text{Cos}_B$  and  $\text{Cos}_W$ , respectively. We also linearly combine  $\text{Cos}_B$ , LSI, and RW with  $\text{Cos}_W$ .

To evaluate the performances of different methods in relevance ranking, we employ MAP [2] and NDCG [15] at positions of 1, 3, and 5 as evaluation measures. For similar query finding, qualitative evaluation and quantitative evaluation were made on the discovered similar queries. In the qualitative evaluation, we randomly sampled examples of similar queries found by  $M\text{-PLS}_{\text{Com}}$  and evaluated their quality. In quantitative evaluation, we randomly sampled 500 queries from each data set, retrieved the top 3 most similar queries for each query, and manually judged the quality of the similar queries found by each method. The final results are presented in pos-com-neg graphs, where “pos” represents the ratio of queries for which our method provides higher quality similar queries, “com” represents the ratio of queries on which the performances of our method and baseline methods are comparable, and “neg” represents the ratio of queries on which the baseline methods do a better job.

Note that all the methods in our experiments are ‘unsupervised methods’. We tested the results for all the methods on the data sets at different parameter settings, and report the best performances of the methods.

Table 1: Weights in combination methods on two data sets

Model	Components	Combination weights	
		Enterprise	Web
Relevance ranking			
LSI+BM25	(LSI,BM25)	(0.9, 0.1)	(0.8, 0.2)
RW+BM25	(RW,BM25)	(0.9, 0.1)	(0.8, 0.2)
Similar query finding			
Cos <sub>B</sub> +Cos <sub>W</sub>	(Cos <sub>B</sub> ,Cos <sub>W</sub> )	(0.9, 0.1)	(0.9, 0.1)
LSI+Cos <sub>W</sub>	(LSI,Cos <sub>W</sub> )	(0.9, 0.1)	(0.9, 0.1)
RW+Cos <sub>W</sub>	(RW,Cos <sub>W</sub> )	(0.9, 0.1)	(0.9, 0.1)

### 5.3. Parameter Setting

We set the parameters of the methods in the following way. In BM25, the default setting is used. There are two parameters in random walk (RW for short): the self-transition probability and the number of transition steps. Following the conclusion in [8], we fix the self-transition probability as 0.9 and choose the number of transition steps from  $\{1, \dots, 10\}$ . We found that RW reaches a “stable” state with a few steps walk. In the experiments on both data sets, after 5 steps we saw no improvement on all the evaluation measures. Therefore, we set 5 as the number of transition steps of RW on both data sets.

In LSI,  $\text{M-PLS}_{\text{Word}}$ ,  $\text{M-PLS}_{\text{Bipar}}$ ,  $\text{M-PLS}_{\text{Conca}}$ , and  $\text{M-PLS}_{\text{Com}}$ , the parameter is the dimensionality of the latent space. We set the dimensionalities in the range of  $\{100, 200, \dots, 1000\}$  for the enterprise search data and in the range of  $\{100, 200, \dots, 3000\}$  for the web search data. We found that when we increase the dimensionality of the latent space, the performances of LSI,  $\text{M-PLS}_{\text{Word}}$ ,  $\text{M-PLS}_{\text{Bipar}}$ ,  $\text{M-PLS}_{\text{Conca}}$ , and  $\text{M-PLS}_{\text{Com}}$  are all improved. The larger the dimensionality is, the better the performance is. On the other hand, a large dimensionality means that we need to calculate more singular values and use more computation power. Therefore, we finally choose 1000 as the dimensionalities of the latent spaces for LSI,  $\text{M-PLS}_{\text{Word}}$ ,  $\text{M-PLS}_{\text{Bipar}}$ ,  $\text{M-PLS}_{\text{Conca}}$ ,  $\text{M-PLS}_{\text{Com}}$  for the enterprise data. We choose 3000 as the dimensionalities of the latent spaces for the web data.

In the combination models, the weights are also parameters. In LSI+BM25 an RW+BM25 for relevance ranking, LSI+ $\text{Cos}_W$ , RW+  $\text{Cos}_W$ ,  $\text{Cos}_B + \text{Cos}_W$  for similar query finding, we choose the combination weights within  $\{0.1, 0.2, \dots, 0.9\}$ . We report the best performance of the combination methods. Table 1 shows the weights in the best performing combination models. In  $\text{M-PLS}_{\text{Com}}$ , the combination weights are chosen automatically using equation (5).

Table 2 shows the properties of matrices for SVD in each method. We can see that although the matrices have high dimensionalities, they are really sparse, and we can efficiently conduct SVD on them.

### 5.4. Experimental Results

#### 5.4.1. Relevance Ranking Results

Table 3 and Table 4 give the results on relevance ranking on two data sets. We can see that on both data sets, our method  $\text{M-PLS}_{\text{Com}}$  not only outperforms the state of the art feature based methods such as BM25 and graph based methods such as RW, but also performs better than their linear combinations. We conducted sign test on the improvement of  $\text{M-PLS}_{\text{Com}}$  over the baseline methods. The results show that the improvements of  $\text{M-PLS}_{\text{Com}}$  are statistically

Table 2: Properties of SVD matrices in each method

Enterprise search data		
	Dimension	Density
M-PLS <sub>Com</sub>	$9958 \times 9958, 81186 \times 51699$	0.4%, 0.08%
M-PLS <sub>Word</sub>	$9958 \times 9958$	0.4%
M-PLS <sub>Bipar</sub>	$81186 \times 51699$	0.08%
M-PLS <sub>Conca</sub>	$91144 \times 61657$	0.5%
LSI	$81186 \times 51699$	0.003%
Web search data		
	Dimension	Density
M-PLS <sub>Com</sub>	$10791 \times 10791, 111631 \times 94022$	0.3%, 0.01%
M-PLS <sub>Word</sub>	$10791 \times 10791$	0.3%
M-PLS <sub>Bipar</sub>	$111631 \times 94022$	0.01%
M-PLS <sub>Conca</sub>	$122422 \times 104813$	0.03%
LSI	$111631 \times 94022$	0.002%

Table 3: Relevance ranking result on enterprise search data

	MAP	NDCG@1	NDCG@3	NDCG@5
M-PLS <sub>Com</sub>	<b>0.644</b>	<b>0.730</b>	<b>0.742</b>	<b>0.754</b>
M-PLS <sub>Conca</sub>	0.624	0.717	0.734	0.748
M-PLS <sub>Word</sub>	0.628	0.711	0.725	0.742
M-PLS <sub>Bipar</sub>	0.589	0.673	0.689	0.707
BM25	0.543	0.646	0.650	0.655
RW	0.603	0.657	0.684	0.702
RW+BM25	0.608	0.666	0.686	0.702
LSI	0.577	0.666	0.681	0.699
LSI+BM25	0.605	0.688	0.702	0.712

significant ( $p < 0.01$ ), except NDCG@1 compared with RW + BM25 on the web search data. Improvement of M-PLS<sub>Com</sub> on the web search data is slightly lower than the improvement of M-PLS<sub>Com</sub> on the enterprise search data. This may be because the bipartite of web search data is more sparse than that of enterprise search data (cf., Table 2).

M-PLS<sub>Conca</sub> also performs well on both data sets. The higher complexity of it makes it less attractive than M-PLS<sub>Com</sub> (Note that the matrix of M-PLS<sub>Conca</sub> for SVD has a high dimensionality and is dense). M-PLS<sub>Bipar</sub> performs worst among the alternatives of our method. This may be because 1) the features of M-PLS<sub>Bipar</sub> are more sparse (cf., Table 2); 2) there is some overlap between the features of M-PLS<sub>Bipar</sub> and the similarities between queries and documents in learning.

Figure 3 and Figure 4 show the performances of M-PLS<sub>Com</sub> with respect to different dimensionalities of the latent spaces on the two data sets. We can see that the performance of M-PLS<sub>Com</sub> will increase when the dimensionalities of the latent spaces increase. On the other hand, the efficiency of the method will decrease, because the time complexity of it is the order of the square of dimensionality ( $O(k_i \cdot C + k_i^2 \cdot \max(s_{qi}, s_{di}))$ ). We find that after the dimensionality reaches 1000 in the enterprise data and 3000 in the web data, the improvement of performance becomes slower. That is why we chose 1000 and 3000 as the dimensionalities of M-PLS<sub>Com</sub> for

Table 4: Relevance ranking result on web search data

	MAP	NDCG@1	NDCG@3	NDCG@5
M-PLS <sub>Com</sub>	<b>0.533</b>	<b>0.667</b>	<b>0.724</b>	<b>0.730</b>
M-PLS <sub>Conca</sub>	0.521	0.658	0.718	0.725
M-PLS <sub>Word</sub>	0.519	0.654	0.712	0.719
M-PLS <sub>Bipar</sub>	0.461	0.593	0.669	0.683
BM25	0.461	0.626	0.682	0.684
RW	0.491	0.656	0.705	0.706
RW+BM25	0.498	0.661	0.712	0.711
LSI	0.448	0.584	0.663	0.675
LSI+BM25	0.484	0.644	0.699	0.702

Table 5: Performances of combination of M-PLS<sub>Com</sub> and BM25 on relevance ranking

	MAP	NDCG@1	NDCG@3	NDCG@5
Enterprise search data	0.651	0.732	0.741	0.753
Web search data	0.537	0.678	0.729	0.732

the two data sets in the experiments. The same thing can be said to the other methods.

Finally, an interesting thing is that when we linearly combine our method M-PLS<sub>Com</sub> with BM25, the performances of our method can be further improved, indicating that our method M-PLS<sub>Com</sub> is complementary to BM25. Table 5 shows the results.

#### 5.4.2. Similar Query Finding Results

We show the performance of the proposed method on similar query finding. We compare M-PLS<sub>Com</sub> with other baselines. We first show some examples, and then we use pos-com-neg graphs to make the comparisons.

*Qualitative Evaluation.* Table 6 gives some examples. For each query, we show the top 3 most similar queries. We can see that M-PLS<sub>Com</sub> is able to find high quality similar queries by effectively using the enriched click-through bipartite graph, even for some “difficult” queries, which contain typos, abbreviations, concatenations and rare words. For the tail queries consisting of rare words, usually it is very hard for the baseline methods to find their similar queries, because they only have a few clicks. Nonetheless, M-PLS<sub>Com</sub> is still able to find similar queries for them, which is very impressive.

*Quantitative Evaluation.* We compared the performances of M-PLS<sub>Com</sub> and the baseline methods on the two data sets. In each data set, we evaluated the quality of similar queries of 500 random queries found by each method.

Figure 5 and Figure 6 present the results. We can see that M-PLS<sub>Com</sub> significantly outperforms the feature based methods and graph based methods, including Cos<sub>B</sub>, Cos<sub>W</sub>, LSI and RW. Our method performs better on more than 15% queries from the enterprise data and on more than 25% queries from the web data. Only on less than 7% queries from the enterprise data and less than 3% queries from the web data, our method performs worse. We conducted sign test, and the results show that all the improvements are statistically significant ( $p < 0.01$ ). Among the combined models, Cos<sub>B</sub> + Cos<sub>W</sub> and RW + Cos<sub>W</sub> perform comparably well, and our method

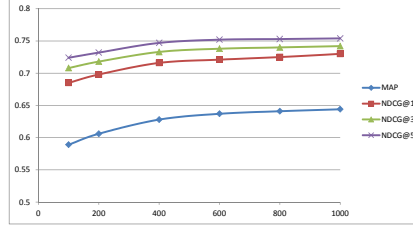


Figure 3: Performances of M-PLS<sub>Com</sub> with respect to different dimensionalities of the latent space on the enterprise data

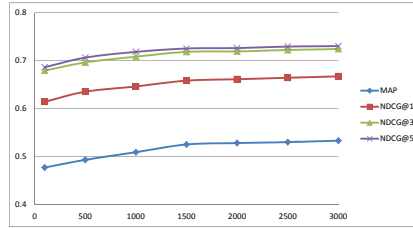


Figure 4: Performances of M-PLS<sub>Com</sub> with respect to different dimensionalities of the latent space on the web data

still outperforms them. However, the improvement of our method over those methods is not so significant.

We investigated the reasons that the methods can or cannot achieve high performance in similar query finding. First, we found that  $\text{Cos}_B$ , LSI and RW are good at handling head queries, since they all rely on co-clicks on the click-through bipartite to calculate query similarity. Among them, RW performs a little better than  $\text{Cos}_B$ . This may be because similarity on the bipartite can be propagated by RW. In contrast, for tail queries, we cannot expect these methods to have good performances. In an extreme case, if a query only has one clicked document and the document is only clicked by the query, then no similar queries can be found for the query. We call this kind of query “isolated island”. Second, no matter whether two queries have co-clicked documents, if they share some words, they are likely to be judged as similar queries by  $\text{Cos}_W$ . Therefore,  $\text{Cos}_W$  is good at handling queries sharing terms. This is especially true for tail queries, since tail queries tend to be longer. However, if two similar queries do not share any term, their similarity cannot be captured by  $\text{Cos}_W$ . The problem is called “term mismatch”, and becomes serious when the queries have spelling errors, abbreviations, and concatenations. The two types of methods (either use click graph or use terms) are complementary, and thus when combined together, it is possible to find similar queries with higher probability.

We found that on most queries our method M-PLS<sub>Com</sub> performs equally well with the combination baseline  $\text{Cos}_B + \text{Cos}_W$ . Sometimes, the two methods even give the same top 3 most similar queries. It indicates that click-through bipartite and features are really useful for similar query finding. On the other hand, we also found that our method can work very well for some really difficult queries on which graph based methods and word based methods fail. The result demonstrates that our method really has the capability to leverage the enriched click-through bipartite to learn query similarities.

Finally, we particularly investigated the performance of our method on tail queries on the web search data, since the web search data set is larger and more sparse than the enterprise search data set. First, Table 7 shows an example. The query ‘walmartmoneycard’ is a concatenation, and



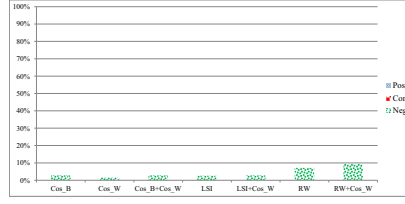


Figure 5: Similar query evaluation on enterprise search data.

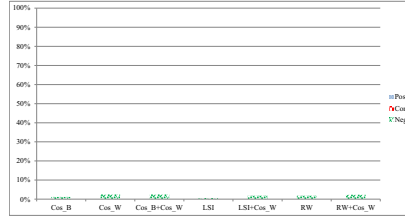


Figure 6: Similar query evaluation on web search data.

is also an isolated island. As a result, for  $\text{Cos}_B$ ,  $\text{LSI}$  and  $\text{RW}$ , no similar queries can be found.  $\text{Cos}_W$  suffers from term mismatch, and thus only those queries that can exactly match ‘walmartmoneycard’ are returned. The combination methods can only use  $\text{Cos}_W$ , and thus return the same results with  $\text{Cos}_W$ . These baseline methods cannot find similar queries for ‘walmartmoneycard’. In contrast,  $\text{M-PLS}_{\text{Com}}$  can work better than the baseline methods. We found that the key reason is that our method can effectively use the click bipartite graph, specifically similar URLs on the graph. The query ‘walmartmoneycard’ has a clicked document <https://www.walmartmoneycard.com/>, and its similar query ‘wal-mart prepaid visa activation’ also has a clicked document <https://www.walmartmoneycard.com/walmart/homepage.aspx>. Through our optimization, both the similarity between ‘walmartmoneycard’ and <https://www.walmartmoneycard.com/> and the similarity between ‘wal-mart prepaid visa activation’ and <https://www.walmartmoneycard.com/walmart/homepage.aspx> are maximized in the latent space (they are 0.71 and 0.53, respectively). Moreover, since the two documents share a common term ‘walmartmoneycard’ in the term space, their similarity is also captured in the latent space through the learning process of our method (with similarity 0.47). Therefore, with the two similar documents as a bridge, our method can connect the two queries together and treat them as similar queries.

We also quantitatively evaluated our method on the tail queries from the web search data. We treat queries with total click numbers less than 10 on the click-through bipartite as tail queries. There are totally 212 tail queries in the 500 samples from the web data. Figure 7 gives the “pos-com-neg” comparison results. We can see that on tail queries, our method  $\text{M-PLS}_{\text{Com}}$  performs even better than itself on the whole sample set. The results indicate the capability of our method on finding high quality similar queries in the tail.

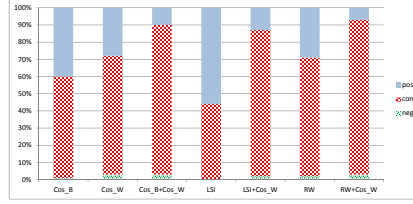


Figure 7: Similar query evaluation on the tail queries from web search data.

## 6. Conclusion and Future Work

In this paper, we have studied the issue of learning query and document similarities from a click-through bipartite with metadata. The click-through bipartite represents the click relations between queries and documents, while the metadata represents multiple types of features of queries and documents. We aim to leverage both the click-through bipartite and features to perform the learning task. We have proposed a method that can solve the problem in a principled way. Specifically, for each type of features, we use two different linear mappings to project queries and documents into a latent space. Then we take the dot product in the latent space as the similarity function between query-document pairs. Similarities between query-query and document-document pairs are simultaneously defined. The final similarity function is defined as a linear combination of similarity functions from different types of features. We learn the mappings and combination weights by maximizing the similarities of the observed query-document pairs on the click-through bipartite, and make orthogonal assumptions on the mappings and regularize the weights using  $L^2$  norm. We have proved that the learning problem has a global optimum. The mappings can be obtained efficiently through Singular Value Decomposition (SVD) and the weights can be obtained from a close form solution of a quadratic program. It turns out to be a new learning method as an extension of Partial Least Squares, referred to as Multi-view PLS. We have theoretically analyzed the proposed method, and demonstrated its capability on finding high quality similar queries (also similar documents). We have conducted experiments on large scale enterprise search and web search data to test the performance of our method. The results not only indicate that our method can significantly outperform the state of the art methods on relevance ranking and similar query finding, but also verify the correctness of our theoretical analysis.

As future work, we want to further enhance the efficiency of our method and test its performance on larger data sets. To achieve the goal, we may need to parallelize the learning algorithm. We also want to study the kernelization of our method, and thus our method is still applicable when kernel matrices instead of features are available.

## References

- [1] I. Antonellis, H.G. Molina, and C.C. Chang. Simrank++: Query rewriting through link analysis of the click graph. *VLDB*, 1(1):408–421, 2008.
- [2] R. Baeza-Yates and B. Ribeiro-Neto. *Modern Information Retrieval*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 1999.
- [3] R. Baeza-Yates and A. Tiberi. Extracting semantic relations from query logs. In *SIGKDD*, pages 76–85, 2007.
- [4] M. Barker and W. Rayens. Partial least squares for discrimination. *Journal of chemometrics*, 17(3):166–173, 2003.
- [5] D. Beeferman and A. L. Berger. Agglomerative clustering of a search engine query log. In *KDD*, pages 407–416, 2000.

- [6] A. Broder, P. Ciccolo, E. Gabrilovich, V. Josifovski, D. Metzler, L. Riedel, and J. Yuan. Online expansion of rare queries for sponsored search. In *WWW'09*, pages 511–520, 2009.
- [7] H. Chun and S. Kelecs. Sparse partial least squares regression for simultaneous dimension reduction and variable selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(1):3–25, 2010.
- [8] N. Craswell and M. Szummer. Random walks on the click graph. In *SIGIR*, pages 239–246, 2007.
- [9] B.D. Davison. Toward a unification of text and link analysis. In *SIGIR'03*, 2003.
- [10] S. Deerwester, S. Dumais, G. Furnas, T. Landauer, and R. Harshman. Indexing by latent semantic analysis. *JASIS*, 41:391–407, 1990.
- [11] H. Deng, M.R. Lyu, and I. King. A generalized co-hits algorithm and its application to bipartite graphs. In *KDD'09*, 2009.
- [12] D.R. Hardoon and J. Shawe-Taylor. Kcca for different level precision in content-based image retrieval. In *Proceedings of Third International Workshop on Content-Based Multimedia Indexing, IRISA, Rennes, France*, 2003.
- [13] D.R. Hardoon and J. Shawe-Taylor. Sparse canonical correlation analysis. *stat*, 1050:19, 2009.
- [14] D.R. Hardoon, S. Szedmak, and J. Shawe-Taylor. Canonical correlation analysis: An overview with application to learning methods. *Neural Computation*, 16(12):2639–2664, 2004.
- [15] K. Jarvelin and J. Kekalainen. Ir evaluation methods for retrieving highly relevant documents. In *SIGIR'00*, pages 41–48, 2000.
- [16] G. Jeh and J. Widom. Simrank: a measure of structural-context similarity. In *SIGKDD*, pages 538–543, 2002.
- [17] H. Ma, H.X. Yang, I. King, and M. R. Lyu. Learning latent semantic relations from clickthrough data for query suggestion. In *CIKM*, pages 709–718, 2008.
- [18] J.M. Ponte and W.B. Croft. A language modeling approach to information retrieval. In *SIGIR'98*, pages 275–281, 1998.
- [19] S. E. Robertson, S. Walker, S. Jones, M. Hancock-Beaulieu, and M. Gatford. Okapi at trec-3. In *TREC*, 1994.
- [20] R. Rosipal and N. Krämer. Overview and recent advances in partial least squares. *Subspace, Latent Structure and Feature Selection*, pages 34–51, 2006.
- [21] S. Rüping and T. Scheffer. Learning with multiple views. In *Proc. ICML Workshop on Learning with Multiple Views*, 2005.
- [22] H. Saigo, N. Krämer, and K. Tsuda. Partial least squares regression for graph mining. In *SIGKDD*, pages 578–586, 2008.
- [23] G. Salton and M. McGill. *Introduction to Modern Information Retrieval*. McGraw-Hill, Inc., New York, NY, USA, 1986.
- [24] P.J. Schreier. A unifying discussion of correlation analysis for complex random vectors. *Signal Processing, IEEE Transactions on*, 56(4):1327–1336, 2008.
- [25] W.R. Schwartz, A. Kembhavi, D. Harwood, and L.S. Davis. Human detection using partial least squares analysis. In *ICCV*, pages 24–31, 2009.
- [26] R.D. Tobias. An introduction to partial least squares regression. In *Proceedings of the Twentieth Annual SAS Users Group International Conference*, pages 1250–1257, 1995.
- [27] C. Wang, R. Raina, D. Fong, D. Zhou, J. Han, and G. Badros. Learning relevance from heterogeneous social network and its application in online targeting. In *SIGIR'11*, 2011.
- [28] X. Wang, J.T. Sun, Z. Chen, and C.X. Zhai. Latent semantic analysis for multiple-type interrelated data objects. In *SIGIR*, pages 236–243, 2006.
- [29] J.A. Wegelin. A survey of partial least squares (pls) methods, with emphasis on the two-block case. *Technical Report, No.371, Seattle: Department of Statistics, University of Washington*, 2000.
- [30] J.R. Wen, J.Y. Nie, and H.J. Zhang. Query clustering using user logs. *ACM Trans. Inf. Syst.*, 20(1):59–81, 2002.
- [31] H. Wold. Path models with latent variables: the nipals approach. *Quantitative sociology: International perspectives on mathematical and statistical modeling*, pages 307–357, 1975.
- [32] W. Wu, J. Xu, H. Li, and O. Satoshi. Learning a robust relevance model for search using kernel methods. *JMLR*, 12:1429–1458, 2011.
- [33] W. Xi, E.A. Fox, W. Fan, B. Zhang, Z. Chen, J. Yan, and D. Zhuang. Simfusion: measuring similarity using unified relationship matrix. In *SIGIR*, pages 130–137. ACM, 2005.
- [34] J. Xu, H. Li, and Z.L. Zhong. Relevance ranking using kernels. In *AIRS '10*, 2010.
- [35] J.F. Xu and G. Xu. Learning similarity function for rare queries. In *WSDM*, pages 615–624, 2011.
- [36] C.X. Zhai and J. Lafferty. A study of smoothing methods for language models applied to information retrieval. *ACM Trans. Inf. Syst.*, 22(2):179–214, 2004.
- [37] Z. Zhang and O. Nasraoui. Mining search engine query logs for query recommendation. In *WWW*, pages 1039–1040, 2006.

### A. Proof of Theorem 4.1

We show the proof of Theorem 4.1 here.

*Proof.* Suppose that  $L_{Q_i} = (l_1^{Q_i}, \dots, l_{k_i}^{Q_i})$ , and  $L_{\mathcal{D}_i} = (l_1^{\mathcal{D}_i}, \dots, l_{k_i}^{\mathcal{D}_i})$ . Objective function (3) can be re-written as  $\sum_{j=1}^{k_i} l_j^{\mathcal{D}_i \top} M_i l_j^{Q_i}$   
 $= \sum_{j=1}^{k_i} \langle l_j^{\mathcal{D}_i}, M_i l_j^{Q_i} \rangle$ . Since  $M_i = \sum_{w=1}^p \lambda_w^i u_w^i v_w^{i\top}$ , we have

$$\begin{aligned} \langle l_j^{\mathcal{D}_i}, M_i l_j^{Q_i} \rangle &= \langle l_j^{\mathcal{D}_i}, \sum_{w=1}^p \lambda_w^i u_w^i v_w^{i\top} l_j^{Q_i} \rangle = \sum_{w=1}^p \lambda_w^i \langle u_w^i l_j^{\mathcal{D}_i \top}, v_w^i l_j^{Q_i} \rangle \\ &\leq \sum_{w=1}^p \lambda_w^i |\langle u_w^i, l_j^{\mathcal{D}_i} \rangle| |\langle v_w^i, l_j^{Q_i} \rangle| \\ &= \lambda_{k_i}^i \sum_{w=1}^p |\langle u_w^i, l_j^{\mathcal{D}_i} \rangle| |\langle v_w^i, l_j^{Q_i} \rangle| + \sum_{w=1}^{k_i} (\lambda_w^i - \lambda_{k_i}^i) |\langle u_w^i, l_j^{\mathcal{D}_i} \rangle| |\langle v_w^i, l_j^{Q_i} \rangle| \\ &\quad + \sum_{w=k_i+1}^p (\lambda_w^i - \lambda_{k_i}^i) |\langle u_w^i, l_j^{\mathcal{D}_i} \rangle| |\langle v_w^i, l_j^{Q_i} \rangle| \\ &\leq \lambda_{k_i}^i \sum_{w=1}^p |\langle u_w^i, l_j^{\mathcal{D}_i} \rangle| |\langle v_w^i, l_j^{Q_i} \rangle| + \sum_{w=1}^{k_i} (\lambda_w^i - \lambda_{k_i}^i) |\langle u_w^i, l_j^{\mathcal{D}_i} \rangle| |\langle v_w^i, l_j^{Q_i} \rangle| \end{aligned}$$

Since  $\|l_j^{\mathcal{D}_i}\| = \|l_j^{Q_i}\| = 1$  and  $\{u_w^i\}_{w=1}^p$  and  $\{v_w^i\}_{w=1}^p$  are orthonormal, we have

$$\sum_{w=1}^p |\langle u_w^i, l_j^{\mathcal{D}_i} \rangle| |\langle v_w^i, l_j^{Q_i} \rangle| \leq \left[ \left( \sum_{w=1}^p \langle u_w^i, l_j^{\mathcal{D}_i} \rangle^2 \right) \left( \sum_{w=1}^p \langle v_w^i, l_j^{Q_i} \rangle^2 \right) \right]^{\frac{1}{2}} \leq \|l_j^{\mathcal{D}_i}\| \cdot \|l_j^{Q_i}\| = 1.$$

Thus, we know  $\langle l_j^{\mathcal{D}_i}, M_i l_j^{Q_i} \rangle \leq \lambda_{k_i}^i + \sum_{w=1}^{k_i} (\lambda_w^i - \lambda_{k_i}^i) |\langle u_w^i, l_j^{\mathcal{D}_i} \rangle| |\langle v_w^i, l_j^{Q_i} \rangle|$ . By taking summation on both sides, we obtain

$$\sum_{j=1}^{k_i} \langle l_j^{\mathcal{D}_i}, M_i l_j^{Q_i} \rangle \leq k_i \lambda_{k_i}^i + \sum_{w=1}^{k_i} (\lambda_w^i - \lambda_{k_i}^i) \left( \sum_{j=1}^{k_i} |\langle u_w^i, l_j^{\mathcal{D}_i} \rangle| |\langle v_w^i, l_j^{Q_i} \rangle| \right).$$

From this inequality, we obtain

$$\sum_{w=1}^{k_i} \lambda_w^i - \sum_{j=1}^{k_i} \langle l_j^{\mathcal{D}_i}, M_i l_j^{Q_i} \rangle \geq \sum_{w=1}^{k_i} (\lambda_w^i - \lambda_{k_i}^i) \left( 1 - \sum_{j=1}^{k_i} |\langle u_w^i, l_j^{\mathcal{D}_i} \rangle| |\langle v_w^i, l_j^{Q_i} \rangle| \right).$$

Since  $L_{Q_i}^\top L_{Q_i} = I_{k_i \times k_i}$  and  $L_{\mathcal{D}_i}^\top L_{\mathcal{D}_i} = I_{k_i \times k_i}$ , we have

$$\sum_{j=1}^{k_i} |\langle u_w^i, l_j^{\mathcal{D}_i} \rangle| |\langle v_w^i, l_j^{Q_i} \rangle| \leq \left[ \left( \sum_{j=1}^{k_i} \langle u_w^i, l_j^{\mathcal{D}_i} \rangle^2 \right) \left( \sum_{j=1}^{k_i} \langle v_w^i, l_j^{Q_i} \rangle^2 \right) \right]^{\frac{1}{2}} \leq \|u_w^i\| \cdot \|v_w^i\| = 1.$$

Thus,  $\sum_{j=1}^{k_i} \langle l_j^{\mathcal{D}_i}, M_i l_j^{Q_i} \rangle \leq \sum_{w=1}^{k_i} \lambda_w^i$ .

Particularly, letting  $l_j^{\mathcal{D}_i} = u_j^i$  and  $l_j^{Q_i} = v_j^i$ , we can obtain the global maximum.  $\square$

Table 6: Examples of similar queries found by M-PLS<sub>Com</sub>

Original Query	Similar Queries
Spelling Errors	
dictionary	web dictionary
	onlinedictionary
	dictionery
wickapedia	www.wikipedia.org
	wikepedia
	www.wikipedia.com
gooole.com	www.google.com
	www.gooogle
	www.goolgle.com
Abbreviations	
fl.lottery	lottery florida
	florida lottery numbers
	florida lottery results
ym	yahoo messenger
	yahoomessenger
	yahoo im messenger
bofa online banking	bank america online banking
	bank of america online
	www.bank of america online banking
Concatenations	
dickssportinggoods	dicks sporting goods
	dicks sporting good store
	dicks sporting goods coupons
googlenews	google news
	goggle news
	news.google.com
peoplesearch	search people
	people search
	yahoo people search
Tail Queries	
star wars anniversary edition lego darth vader fighter	www.star wars legos.com
	star wars legos
	star wars lego
american express online account summary	american express account online
	american express account
	american express online
read full books online free	read books online free
	free online books to read
	read books online

Table 7: Result on a difficult query

Query: walmartmoneycard	
M-PLS <sub>Com</sub>	www.walmartmoneycard.com
	walmartmoneycard.com
	wal-mart prepaid visa activation
Cos <sub>B</sub> , LSI, RW	N / A
	N / A
	N / A
Cos <sub>W</sub>	www.walmartmoneycard.com
	walmartmoneycard.com
	N / A
Cos <sub>B</sub> + Cos <sub>W</sub> , LSI + Cos <sub>W</sub> , RW+ Cos <sub>W</sub>	www.walmartmoneycard.com
	walmartmoneycard.com
	N / A